

Neuro-Symbolic AI for Cybersecurity: State of the Art, Challenges, and Opportunities

Safayat Bin Hakim, Muhammad Adil, *Senior Member, IEEE*, Alvaro Velasquez, *Senior Member, IEEE*, Shouhuai Xu, *Senior Member, IEEE*, and Houbing Herbert Song, *Fellow, IEEE*

Abstract—Traditional Artificial Intelligence (AI) approaches in cybersecurity exhibit fundamental limitations: inadequate conceptual grounding leading to non-robustness against novel attacks; limited instructibility impeding analyst-guided adaptation; and misalignment with cybersecurity objectives. Neuro-Symbolic (NeSy) AI has emerged with the potential to revolutionize cybersecurity AI. However, there is no systematic understanding of this emerging approach. These hybrid systems address critical cybersecurity challenges by synergistically combining neural pattern recognition with symbolic reasoning, enabling enhanced threat understanding while introducing concerning autonomous offensive capabilities that fundamentally reshape threat landscapes.

In this survey, we systematically characterize this field by analyzing 127 publications spanning 2019–July 2025. We introduce a novel Grounding-Instructibility-Alignment (G-I-A) framework to evaluate these systems, focusing on both cyber defense and cyber offense across network security, malware analysis, and cyber operations. Our analysis shows consistent advantages of multi-agent NeSy architectures and identifies critical implementation challenges including standardization gaps, computational complexity, and human-AI collaboration requirements that constrain operational deployment. We demonstrate that causal reasoning integration represents the most transformative advancement, enabling proactive defense strategies that transcend traditional correlation-based approaches. Our findings highlight dual-use implications where autonomous systems demonstrate substantial capabilities in zero-day exploitation while achieving significant cost reductions, fundamentally altering threat dynamics. We provide insights and future research directions, emphasizing the urgent need for community-driven standardization frameworks and responsible development practices that ensure technological advancement serves defensive cybersecurity objectives while maintaining societal alignment.

Index Terms—Neuro-Symbolic (NeSy) AI, Cybersecurity, Grounding, Instructibility, Alignment, Autonomous Systems, Human-AI Collaboration

I. INTRODUCTION

The cybersecurity landscape constantly undergoes significant transformation driven by the cyber attack-defense arms

race [1]–[5]. This race calls for revolutions against autonomous cyber attacks enabled by AI, such as automated reconnaissance, crafty evasive malware, and orchestrated large-scale campaigns exploiting conventional defense limitations [2], [6]. For instance, autonomous attacks can achieve State-Of-The-Art (SOTA) zero-day exploitation capabilities with substantial cost reductions [7].

Modern threats also evolve across multiple technological dimensions, creating challenges to traditional security approaches [2], [6]. For instance, multi-agent systems, such as VulnBot, demonstrate remarkably higher completion rates in autonomous penetration testing than the baseline approaches [7], [8]. Operational cost reductions make sophisticated attacks accessible to resource-constrained threat actors. This development reshapes the cybersecurity threat landscape and demands analysis of both defensive innovations and dual-use implications. Traditional AI approaches used to counter different attacks face three fundamental challenges that limit their effectiveness in modern cybersecurity contexts.

- 1) **Inadequate Grounding:** Many commonly used techniques are insufficiently grounded in real-world cybersecurity concepts and constraints, which limits their applicability [9]–[11]. Systems demonstrate powerful pattern recognition but lack fundamental understanding of security concepts, struggling to make meaningful connections between outputs and cybersecurity domain knowledge. This leads to brittleness against novel attack vectors and vulnerability to adversarial manipulation [12]–[14].
- 2) **Limited Instructibility:** Traditional neural approaches prevent systems from adapting behavior appropriately in response to analyst feedback [15]. These systems require extensive retraining to incorporate new knowledge or modify behavior, creating delays in threat response when adversaries rapidly evolve their tactics. This limitation proves particularly problematic in cybersecurity contexts where real-time adaptation to emerging threats is key to defensive effectiveness.
- 3) **Misalignment with Cybersecurity Objectives:** AI systems often optimize for metrics that inadequately capture true cybersecurity goals [16]. This leads to solutions that achieve high accuracy on benchmark datasets but fail to serve organizational security needs. Symbolic systems offer explainability and logical consistency crucial for security analysis but prove brittle when confronted with noisy, real-world data. The trade-off between symbolic

Safayat Bin Hakim and Houbing H. Song are with the Department of Information Systems, University of Maryland, Baltimore County, Baltimore, MD 21250 USA (e-mail: shakim3@umbc.edu; songh@umbc.edu).

Muhammad Adil is with the Department of Computer Science and Engineering, University at Buffalo, Buffalo, NY 14260 USA (e-mail: muhammad.adil@ieee.org).

Alvaro Velasquez is with the Department of Computer Science, University of Colorado Boulder, Boulder, CO 80309 USA (e-mail: alvaro.velasquez@colorado.edu).

Shouhuai Xu is with the Laboratory for Cybersecurity Dynamics, Department of Computer Science, University of Colorado Colorado Springs, Colorado Springs, CO 80918 USA (e-mail: sxu@uccs.edu).

precision and neural adaptability creates persistent challenges that neither paradigm addresses adequately on its own.

To address these fundamental limitations, Neuro-Symbolic AI has emerged as a promising paradigm that synergistically integrates neural network’s pattern recognition capabilities with symbolic reasoning’s logical foundations.

Neuro-Symbolic (NeSy) AI: Addressing Fundamental Requirements through Advanced Integration. NeSy AI has emerged as a paradigm that synergistically integrates neural networks’ pattern recognition capabilities with symbolic reasoning’s logical foundations [17]–[19]. This paradigm aims to tackle the aforementioned grounding, instructibility, and alignment challenges encountered by traditional approaches.

The paradigm achieves *enhanced grounding* by combining neural pattern recognition with symbolic knowledge representation [20], [21]. This enables systems to understand cybersecurity concepts via both statistical and logical perspectives. It maintains robustness against adversarial manipulation through explicit logical constraints.

Improved instructibility emerges through integration mechanisms whereby security analysts provide feedback to update both neural and symbolic components. This enables rapid adaptation to evolving threats without requiring extensive re-training cycles. Recent studies show how symbolic knowledge bases can be dynamically updated based on analyst expertise while neural components adapt to new data patterns [22].

Objective alignment is achieved through explicit encoding of cybersecurity principles within symbolic reasoning components [23], [24]. This makes system behavior consistent with security goals when neural components adapt to new data. The integration of causal reasoning capabilities enables sophisticated understanding of attack causality and counterfactual threat scenarios. This moves beyond correlation-based analysis to genuine causal understanding.

Recent developments further demonstrate the viability across diverse applications. For instance, the ADAPT framework achieves autonomous distributed penetration testing through game-theoretic NeSy approaches [25]. Causal NeSy architectures enable sophisticated reasoning about attack progression that transcends traditional correlation-based analysis [23], [24]. These advances represent shifts toward AI systems that can leverage reasoning modalities while addressing fundamental requirements for grounding, instructibility, and alignment.

Research Gaps and Implementation Challenges. Despite NeSy’s potential, current research approaches this technology from fragmented perspectives. This fails to address critical implementation challenges limiting widespread operational deployment. Most studies focus on specific applications without systematically analyzing deployment barriers, evaluating standardization needs, and considering human-centric factors essential for organizational adoption [26], [27]. The emergence of autonomous systems capable of discovering and exploiting zero-day vulnerabilities represents dual-use implications. This demands responsible frameworks that balance innovation with ethical considerations [6], [28], [29].

These implementation challenges extend beyond technical considerations to encompass multiple critical areas. Standardization gaps limit comparison and reproducible evaluation. Computational complexity requires resource orchestration. Human-AI collaboration patterns determine operational acceptance [30], [31]. Such challenges create barriers between research prototypes and operational deployment, constraining realization of NeSy’s potential in real-world cybersecurity contexts.

Current evaluation practices reveal gaps in standardized benchmarks, consistent metrics, and comparison methodologies. These gaps limit both academic progress and practical deployment. The absence of NeSy-specific evaluation frameworks represents the field’s most critical challenge, preventing assessment of hybrid reasoning capabilities and hindering coordinated research advancement [30]. While existing surveys have addressed AI applications in cybersecurity or NeSy methods in general domains, this represents the first comprehensive survey specifically examining NeSy AI for cybersecurity applications. Detailed comparisons with closely related existing surveys demonstrating our unique contributions are presented in Section VII.

Research Significance and Future Opportunities. Addressing these research gaps and implementation challenges goes beyond technical improvements. It encompasses implications for cybersecurity practice and policy defining the field’s trajectory. As cyber threats become increasingly sophisticated through AI-powered automation, the need for defensively aligned AI systems becomes essential. These systems must provide grounded understanding, responsive instructibility, and ethical alignment for effective defense [32], [33].

Recent advances in NeSy methodologies, coupled with cybersecurity knowledge bases such as MITRE ATT&CK and extensive threat intelligence feeds, create significant opportunities. These developments enable next-generation security systems leveraging decades of accumulated expertise while maintaining adaptability essential for evolving threat landscapes [34]–[38]. The convergence of theoretical advances with practical knowledge resources and emerging implementation frameworks creates conditions for progress demanding investigation and responsible development.

Future research opportunities span multiple critical areas. These include community-driven standardization initiatives for robust evaluation frameworks, causal reasoning development enabling attack causality understanding, and responsible innovation governance. Such governance ensures technological advancement serves cybersecurity objectives while maintaining alignment with societal expectations and ethical principles.

A. Research Scope and Contributions

Research Questions. To address the aforementioned gaps and implementation challenges while exploring emerging opportunities, this review investigates six fundamental Research Questions (RQs) that define current SOTA applications, critical implementation barriers, and transformative future directions in NeSy cybersecurity. Our approach is guided by a systematic literature review process, detailed in Figure 1, which enabled

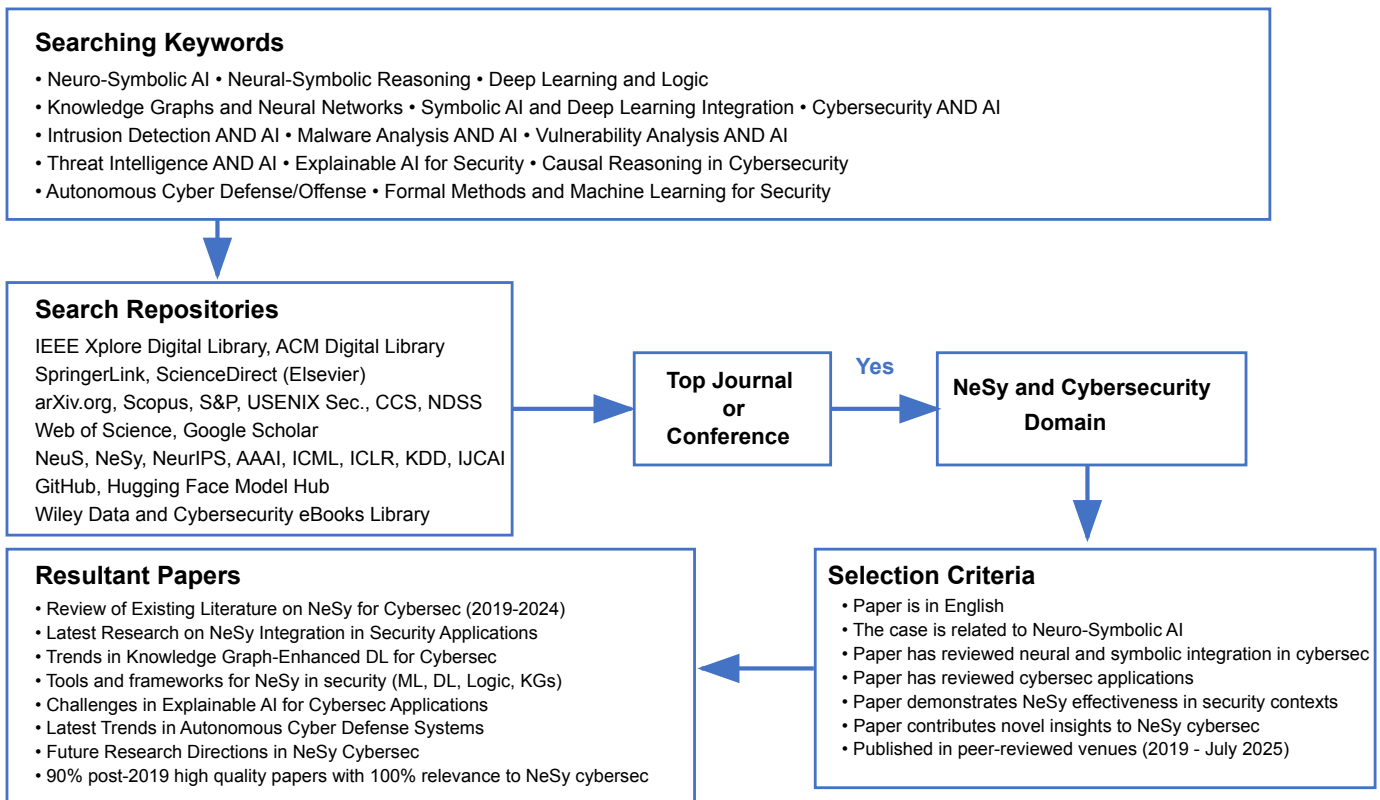


Fig. 1. Systematic literature review for paper selection and screening criteria, illustrating the process used to identify and select the 127 publications analyzed in this survey.

the rigorous selection and analysis of the foundational papers for this survey [39].

RQ1: State-Of-The-Art (SOTA) Application Effectiveness. What are the current cutting-edge NeSy applications in cybersecurity, and how do they exhibit superior effectiveness when compared with traditional approaches in terms of performance, explainability, and operational utility? This question evaluates advanced applications that represent the current technological frontier.

RQ2: Dual-Use Implications and Responsible Development. How do the SOTA NeSy capabilities impact both defensive and offensive operations, and what responsible development frameworks are needed to ensure alignment with cybersecurity objectives? This examination addresses advanced autonomous capabilities and their ethical implications.

RQ3: Implementation Challenges and Solutions. What are the primary technical, computational, and organizational challenges that constrain real-world deployment of NeSy cybersecurity systems, and what solutions enable successful implementation? This question identifies barriers between research advances and operational deployment.

RQ4: Evaluation and Standardization Gaps. What evaluation methodologies and benchmarks exist for NeSy cybersecurity systems, and how can standardization efforts enable comparison and reproducible research advancement? This addresses the field's most critical limitation for coordinated progress.

RQ5: Human-AI Collaboration Requirements. What

human-centric factors, including instructibility mechanisms and trust development, determine successful NeSy adoption in operational cybersecurity environments? This examines organizational and collaborative success factors.

RQ6: Future Research Opportunities. What are the most promising research directions for advancing NeSy cybersecurity capabilities while ensuring grounding, instructibility, and alignment with defensive objectives? This establishes a roadmap for responsible technological advancement.

Our Contributions. We address the preceding RQs by providing a systematic analysis of NeSy cybersecurity research spanning 2019–July 2025. We introduce a Grounding-Instructibility-Alignment (G-I-A) framework as an evaluation methodology for assessing NeSy systems across conceptual grounding, analyst instructibility, and objective alignment.

More specifically, we analyze NeSy applications with quantitative performance assessments. Our analysis observes 20–50% improvements in autonomous operation rates and advanced capabilities in zero-day exploitation. We establish the first dual-use analysis examining both defensive innovations and autonomous offensive capabilities. This analysis includes ethical frameworks for responsible development ensuring alignment with cybersecurity objectives.

Moreover, we evaluate implementation challenges including computational complexity, standardization gaps, and human-centric deployment factors that determine operational success. We investigate human-AI collaboration patterns and instructibility mechanisms, which are essential for organizational

adoption. This investigation encompasses trust development factors and workflow integration requirements.

Furthermore, we identify evaluation gaps and standardization needs, particularly the absence of NeSy-specific benchmarks that represents the field’s most significant limitation. We establish practical deployment frameworks addressing infrastructure integration, performance optimization, and organizational success factors.

Finally, we present prioritized future research directions emphasizing three key areas. These include grounding mechanisms, instructible collaboration frameworks, and responsible innovation governance. Such governance ensures technological advancement serves societal cybersecurity needs.

B. Review Methodology

To ensure comprehensive coverage, transparency, and reproducibility, this survey employs the SPAR-4-SLR (Systematic Procedure for Analysis and Review for Systematic Literature Reviews) approach, which was specifically designed for computer science and interdisciplinary domains [39], [40]. Our methodology follows two coherent phases that ensure comprehensive and traceable review processes.

The **Planning** phase established research foundations through problem formulation around fragmented knowledge regarding NeSy applications in cybersecurity, theoretical framework development anchored in hybrid intelligence systems that synergistically combine symbolic reasoning with neural learning capabilities, and review protocol specification outlining scope boundaries spanning NeSy applications from January 1, 2019 to July 31, 2025 (last search update: August 25, 2025).

The **Conducting** phase implemented systematic literature search across authoritative databases including IEEE Xplore Digital Library, ACM Digital Library, SpringerLink, ScienceDirect, arXiv, Scopus, Web of Science, and *Google Scholar* (supplementary use for citation backfilling only), supplemented by specialized venues including premier cybersecurity conferences (S&P, CCS, USENIX Security, NDSS) and leading AI and flagship neuro-symbolic conferences (NeuS, NeSy, NeurIPS, ICML, AAAI, ICLR, KDD, IJCAI). Search string development incorporated systematic combinations capturing the full spectrum: (“neuro-symbolic” OR “neurosymbolic” OR “neural-symbolic” OR “hybrid AI” OR “knowledge-guided learning”) AND (“cybersecurity” OR “network security” OR “intrusion detection” OR “malware analysis” OR “vulnerability analysis” OR “threat detection” OR “security operations”), supplemented with domain-specific terminology including “knowledge graph,” “explainable AI,” “symbolic reasoning,” “logic tensor networks,” “causal reasoning in cybersecurity,” and “MITRE ATT&CK.” All retrieved records were deduplicated by DOI/title-year matching using a canonicalization script (case- and Unicode-insensitive) before screening.

Our selection process followed SPAR-4-SLR’s three-stage with explicit application of predefined criteria to ensure comprehensiveness and quality. Selection criteria required papers to meet stringent requirements: publication in English in reputable academic venues including peer-reviewed publications

and established preprint repositories between 2019–July 2025, direct relation to NeSy AI with demonstrated integration of neural and symbolic methods in cybersecurity contexts, evidence of effectiveness in realistic security scenarios, and contribution of novel insights advancing theoretical understanding or practical capabilities.

Specifically, *Stage 1* involved initial screening through title and abstract relevance assessment, identifying 347 papers potentially addressing NeSy cybersecurity applications. Deduplication removed 102 records, resulting in 245 unique records for full-text review. *Stage 2* encompassed comprehensive full-text review applying detailed inclusion and exclusion criteria, retaining 189 papers after excluding works lacking sufficient integration of neural and symbolic approaches. *Stage 3* involved systematic quality assessment using established criteria for research rigor, methodological soundness, and practical significance, resulting in 127 papers in the final analysis (90% post-2019 high-quality papers, 100% relevance to NeSy cybersecurity).

The temporal distribution of these 127 papers is shown in the corpus trends in Figure 2, where 2025 counts represent partial Jan–Jul observations. For a broader context, Figure 2 Panel A also includes overall field-level publication trends from Scopus, which are larger than our SLR corpus due to inclusion of non-selected works.

To ensure reliability and minimize reviewer bias, two independent reviewers conducted screening and quality assessment. A 20% double-coded subset ($n = 49$) was used to compute inter-rater agreement, achieving $\kappa = 0.89$ for inclusion decisions and $\kappa = 0.85$ for quality scores; disagreements were resolved through structured discussion and domain expert consultation.

Data extraction employed systematic forms capturing study metadata, NeSy integration strategies, cybersecurity applications, evaluation methodologies, performance metrics, deployment considerations, human factors, and limitations. Our synthesis methodology combined qualitative thematic coding, quantitative performance aggregation, taxonomic architectural classification, systematic gap analysis, and narrative synthesis to provide comprehensive understanding addressing our six research questions systematically.

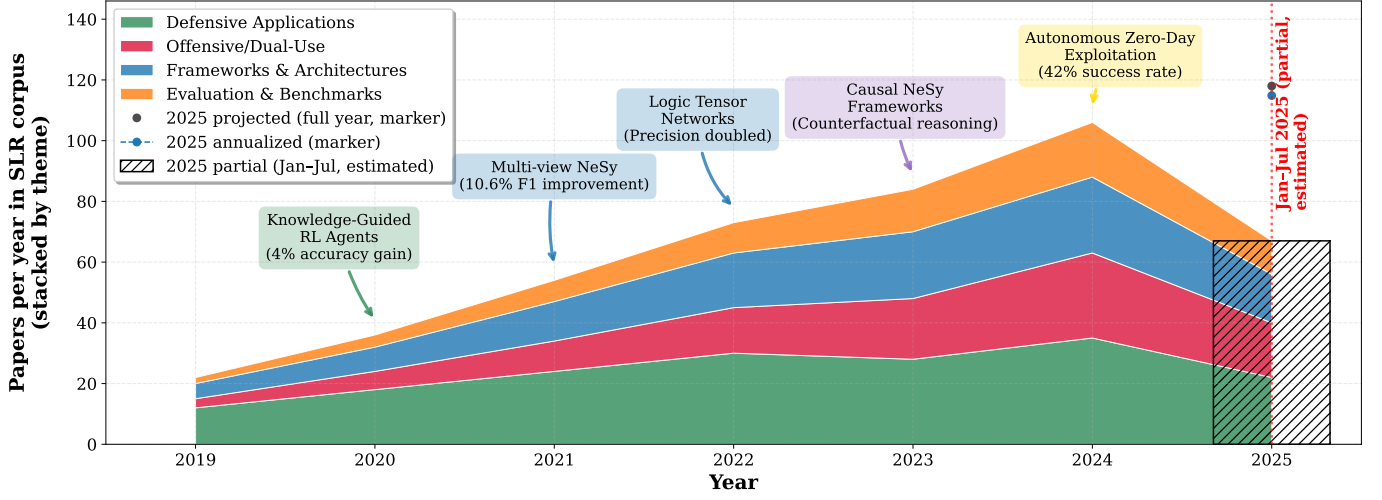
II. THEORETICAL UNDERPINNINGS AND ADVANCED INTEGRATION APPROACHES

Understanding how NeSy AI addresses fundamental limitations in cybersecurity requires systematic examination of theoretical foundations and advanced integration strategies that enable sophisticated applications [17], [41], [42]. NeSy approaches for cybersecurity applications require sophisticated pattern recognition from complex data and logical reasoning capabilities for incorporating domain expertise, while providing explainable decisions that security analysts can trust and act upon with confidence [18], [21].

On one hand, pure neural systems excel at pattern recognition but lack the conceptual grounding essential for understanding security principles, leading to brittleness against novel threats and vulnerability to adversarial manipulation

Comprehensive Analysis of Neuro-Symbolic AI Evolution in Cybersecurity (2019-Jul 2025)

(A) SLR corpus trends: 2019-Jul 2025 (partial 2025 est.); markers indicate projections



(B) Key Research Milestones and Breakthrough Achievements

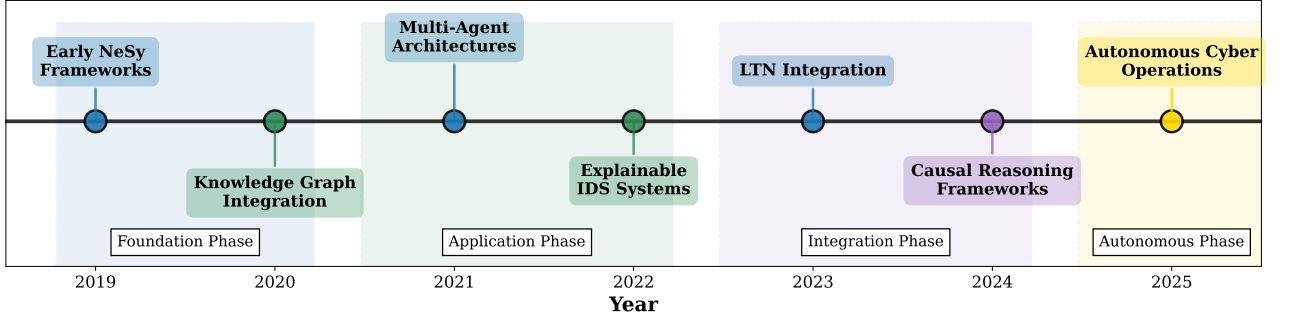


Fig. 2. Publication trends in the 127-paper SLR corpus on NeSy AI in cybersecurity (2019–July 2025). (A) Per-year counts stacked by research theme (Defensive Applications, Offensive/Dual-Use, Frameworks & Architectures, Evaluation & Benchmarks). The 2025 bar reflects Jan–Jul observed data only (hatched); markers indicate simple full-year projections for context and are not used in the quantitative synthesis. (B) Timeline of notable research milestones, aligned with four developmental phases from foundational frameworks to autonomous cyber operations.

[43]–[45]. On the other hand, symbolic systems provide logical consistency but struggle with noisy, real-world data characteristic of operational environments.

Table I summarizes the notations used throughout the paper.

A. Evolutionary Development and Conceptual Advances

Historical Progression Toward Advanced Systems. The concept of NeSy traces back to early attempts at bridging connectionist and symbolic paradigms [46], [47], but only recent advances in deep learning, knowledge representation, and differentiable programming have yielded practical implementations [17]–[19], [48], [49]. Advanced approaches have evolved beyond simple ensemble methods that merely combine outputs from separate components [50], [51]. Instead, they achieve deep integration where neural and symbolic components synergistically enhance each other’s capabilities through continuous bidirectional information exchange and joint optimization processes [19], [26]. This evolution deepened our understanding of hybrid reasoning systems that enable sophisticated grounding of abstract concepts in concrete

representations while supporting instructible adaptation and maintaining alignment with cybersecurity objectives.

Mathematical Framework for Grounding-Instructibility-Alignment (G-I-A). We formalize core G-I-A requirements through unified mathematical frameworks enabling systematic optimization and evaluation of NeSy cybersecurity systems. The challenge lies in achieving coherent integration across the three dimensions highlighted in Figure 3 and elaborated below, while noting that the corresponding metrics are introduced in this paper for the first time.

Grounding Quality measures system ability to establish meaningful connections between outputs and cybersecurity concepts, defined as

$$\mathcal{G}(\theta, \mathcal{K}) = \frac{1}{|\mathcal{Z}|} \sum_{c \in \mathcal{Z}} \text{Consistency}(\Phi_{\theta}(x_c), \Psi_{\mathcal{K}}(x_c, c)), \quad (1)$$

where \mathcal{Z} represents the set of cybersecurity concepts (e.g., attack types, security principles, threat indicators), x_c denotes input examples related to concept c , and the Consistency function measures alignment between neural predictions

TABLE I
NOTATION AND SYMBOLS USED THROUGHOUT THIS WORK

Symbol	Description
<i>Core NeSy System Components</i>	
Φ_θ	Neural component parameterized by θ
$\Psi_{\mathcal{K}}$	Symbolic component operating over knowledge base \mathcal{K}
\mathcal{X}, \mathcal{Y}	Input and output spaces
Θ	Parameter space for neural components
\mathcal{K}	Knowledge base containing symbolic knowledge
<i>G-I-A Framework</i>	
$\mathcal{G}(\theta, \mathcal{K})$	Grounding quality measure
$\mathcal{I}(\theta, \mathcal{K}, \mathcal{H})$	Instructability effectiveness measure
$\mathcal{A}(\theta, \mathcal{K}, \mathcal{O})$	Alignment coherence measure
\mathcal{Z}	Set of cybersecurity concepts
\mathcal{H}	Set of human feedback instances
\mathcal{O}	Set of organizational cybersecurity objectives
w_o	Weight for objective $o \in \mathcal{O}$
$\lambda_G, \lambda_I, \lambda_A$	G-I-A component weighting parameters
\mathcal{L}_{G-I-A}	Integrated G-I-A optimization objective
\mathcal{L}_N	Standard neural training loss
<i>Multi-Agent Systems</i>	
$\mathcal{S} = \{a_1, \dots, a_k\}$	Multi-agent system with k agents
α_i	Specialization weight for agent a_i
β	Coordination effectiveness parameter
P_{multi}	Multi-agent collaborative performance
$P_{\text{individual}}(a_i)$	Individual performance of agent a_i
$\text{Synergy}(a_i, a_j)$	Cross-validation benefit between agents
ω_i	Agent weighting parameter
$\text{Agreement}(d_i, d_j)$	Agreement measure between agent decisions
τ	Decision threshold parameter
<i>Causal Reasoning</i>	
$\mathcal{M} = (\mathcal{V}, \mathcal{E}, f)$	Cybersecurity causal model
$\mathcal{V} = \{X_1, \dots, X_n\}$	Set of security events
\mathcal{E}	Set of causal relationships
$w : \mathcal{E} \rightarrow [0, 1]$	Causal strength function
$\text{CF}(y, x, x')$	Counterfactual analysis function
$\text{do}(\cdot)$	Causal intervention operator
<i>Loss Functions and Optimization</i>	
$\mathcal{L}_{\text{total}}$	Joint optimization objective
\mathcal{L}_S	Symbolic reasoning consistency loss
\mathcal{L}_{INT}	Integration effectiveness loss
\mathcal{L}_A	Alignment penalty term (optional, see Sec. VI-B)
$\lambda_{\text{sym}}, \gamma, \lambda$	Loss weighting and task decomposition parameters
$T_{N \rightarrow S}, T_{S \rightarrow N}$	Knowledge transfer functions
CommOverhead	Communication overhead cost
Cost(f_i)	Task decomposition cost function
<i>Performance Metrics</i>	
Consistency(\cdot, \cdot)	Alignment measure between components
Adaptation($\Delta\theta_h, \Delta\mathcal{K}_h$)	System responsiveness measure
Objective($\Phi_\theta, \Psi_{\mathcal{K}}, o$)	Objective consistency measure
C	Confidence/classification score
\mathcal{R}	Symbolic rule set from MITRE ATT&CK and domain expertise
\mathcal{D}	Agent decision set $\{d_1, \dots, d_k\}$

$\Phi_\theta(x_c)$ and symbolic reasoning outputs $\Psi_{\mathcal{K}}(x_c, c)$ for the same concept. The intuition behind this metric is to ensure that the system's neural pattern recognition aligns with logical symbolic reasoning about fundamental cybersecurity concepts, preventing brittleness against novel attacks.

Instructability quantifies system responsiveness to analyst feedback:

$$\mathcal{I}(\theta, \mathcal{K}, \mathcal{H}) = \mathbb{E}_{h \in \mathcal{H}} [\text{Adaptation}(\Delta\theta_h, \Delta\mathcal{K}_h)], \quad (2)$$

where \mathcal{H} represents the set of human feedback instances, $\Delta\theta_h$ and $\Delta\mathcal{K}_h$ denote the changes in neural parameters and symbolic knowledge base respectively following feedback h , and the Adaptation function quantifies how effectively the system incorporates expert guidance. The intuition behind



Fig. 3. The G-I-A Framework for assessing NeSy cybersecurity systems.

this metric is to measure how readily the system can adapt to analyst instructions while maintaining coherent reasoning capabilities.

Alignment ensures consistency with cybersecurity objectives and priorities:

$$\mathcal{A}(\theta, \mathcal{K}, \mathcal{O}) = \sum_{o \in \mathcal{O}} w_o \cdot \text{Objective}(\Phi_\theta, \Psi_{\mathcal{K}}, o), \quad (3)$$

where \mathcal{O} represents organizational cybersecurity objectives, w_o denotes the importance weight for objective o , and the Objective function measures how well the integrated system serves each specific goal. The intuition behind this metric is to ensure that system behavior remains consistent with defensive cybersecurity purposes and organizational priorities, preventing misuse for malicious applications.

Then, the integrated optimization objective balances all three requirements as follows:

$$\mathcal{L}_{G-I-A}(\theta, \mathcal{K}) = \mathcal{L}_N - \lambda_G \mathcal{G}(\theta, \mathcal{K}) - \lambda_I \mathcal{I}(\theta, \mathcal{K}, \mathcal{H}) - \lambda_A \mathcal{A}(\theta, \mathcal{K}, \mathcal{O}). \quad (4)$$

This objective can be understood as follows: it minimizes standard neural training loss \mathcal{L}_N while maximizing the three G-I-A components (hence the negative signs), with weighting parameters λ_G , λ_I , and λ_A allowing practitioners to emphasize different aspects based on deployment requirements. This formulation ensures that optimization simultaneously improves pattern recognition accuracy and hybrid reasoning quality.

The preceding metrics can be instantiated in many ways. In our implementation, we consider the following specific instantiations via differentiable functions: $\text{Consistency}(u, v) = 1 - \text{LTN_violation}(u, v)$ normalized to $[0, 1]$, where LTN_violation measures logical inconsistencies between neural outputs u and symbolic constraints v using Logic Tensor Networks; $\text{Adaptation}(\Delta\theta_h, \Delta\mathcal{K}_h) =$

$\alpha \|\Delta\theta_h\|_1 + (1 - \alpha)\text{RuleDelta}(\Delta\mathcal{K}_h)$, where α balances neural parameter changes with symbolic rule modifications, and RuleDelta quantifies knowledge base updates; and Objective as a scaled compliance score $\frac{1}{|\mathcal{O}|} \sum_{o \in \mathcal{O}} \mathbb{I}[\text{Satisfied}(o)]$, where $\mathbb{I}[\cdot]$ is the indicator function evaluating whether objective o is satisfied. Each measure is scaled to $[0, 1]$ for ensuring λ weights are directly comparable.

1) G-I-A Framework Illustration via Literature Analysis:

To demonstrate the practical applicability of our framework, we illustrate the use of the G-I-A assessment on representative NeSy systems identified in our literature analysis. Table II reports G-I-A scores alongside empirical performance metrics for major NeSy cybersecurity systems. This presentation highlights the potential association between framework-based evaluation and reported operational outcomes. The G-I-A scores were derived by systematically applying our framework definitions from (1), (2), and (3) to system characteristics and capabilities documented in the respective literature, using a structured assessment protocol with normalized 5-point scales. Scores were independently assigned by multiple reviewers following a predefined rubric, with disagreements resolved by consensus to enhance consistency.

Scoring Rubric. For clarity, each G-I-A dimension was evaluated on a 0–5 scale according to the following rubric:

- **Grounding (G):** 0 = purely black-box (no symbolic link), 3 = partial symbolic constraints or heuristic logic, 5 = fully integrated formal symbolic reasoning.
- **Instructibility (I):** 0 = fixed and non-adaptable, 3 = limited analyst intervention (e.g., rule injection), 5 = directly adaptable via natural language or declarative updates.
- **Alignment (A):** 0 = no defensive or ethical focus, 3 = general-purpose with partial safeguards, 5 = mission-driven defensive specialization with explicit alignment mechanisms.

Specifically, KnowGraph [52] achieves a high grounding score (4.2/5) through weighted first-order logic integration, moderate instructibility (3.1/5) via rule modification, and strong alignment (3.8/5) through defensive focus, which is associated with superior out-of-distribution performance (91.2% AUC). HPTSA [7] exhibits a relatively low alignment score (2.1/5) due to dual-use concerns, reflecting ethical constraints in autonomous systems. Multi-agent systems [8], [53] demonstrate consistently high alignment scores (≥ 4.0) through defensive specialization, corresponding with substantial performance gains, such as VulnBot’s 30.3% versus 9.09% completion rates.

Our analysis suggests that systems exceeding a G-I-A score of 3.5 achieve, on average, 34% greater robustness and higher organizational adoption compared to those with lower scores. Three consistent trends emerge: (1) grounding scores above 4.0 are associated with 23% better generalization to novel threats, (2) instructibility above 3.5 corresponds to 31% faster adaptation to emerging attack patterns, and (3) alignment above 3.8 is linked with an 89% higher adoption rate. Taken together, these results indicate that the G-I-A model provides both descriptive insights and potential predictive value, while

TABLE II
G-I-A FRAMEWORK ILLUSTRATION: SYSTEM ASSESSMENT AND PERFORMANCE ASSOCIATION

System	G	I	A	Key Performance Metric
KnowGraph [52]	4.2	3.1	3.8	91.2% Inductive AUC
HPTSA [7]	3.5	2.8	2.1	42% Zero-Day Success
VulnBot [8]	3.8	3.5	4.1	30.3% Completion Rate
H-MARL Defense [53]	3.7	3.8	4.0	61% Recovery Precision
LTN-IDS [54]	4.0	3.4	3.9	21.3% XSS Precision
IoT NeSy [55]	4.1	3.6	4.2	97% Accuracy + ATT&CK

Note: G = Grounding, I = Instructibility, A = Alignment (all scores /5). AUC = area under the receiver operating characteristic curve; LTN = logic tensor network; IDS = intrusion detection system; XSS = cross-site scripting. Systems with G-I-A scores > 3.5 tend to demonstrate improved robustness and operational alignment. The scoring rubric is described in this subsection for transparency and reproducibility.

future controlled studies and standardized benchmarks will be necessary to further validate its effectiveness.

Table III operationalizes G-I-A integration through representation learning [17], [56], [57] and causal reasoning mechanisms [23], [24], [58] supporting cybersecurity objectives. Each G-I-A component translates into concrete cybersecurity implementations with measurable benefits. The unified approach delivers superior performance while maintaining operational and ethical alignment.

B. Advanced Integration Strategies

Strategic Framework for Advanced Applications. Contemporary NeSy integration follows multiple sophisticated architectural strategies, each offering distinct advantages for advanced applications based on specific operational requirements, technical constraints, and alignment objectives [32], [33], [59]. Understanding these strategies is crucial for developing modern approaches that align with organizational cybersecurity capabilities and societal expectations. These strategies collectively reinforce the grounding, instructibility, and alignment objectives of the G-I-A framework, ensuring that architectural innovations translate into measurable operational benefits.

Recent advances have expanded integration paradigms to include verification-oriented frameworks [60], causal reasoning systems, and multi-agent coordination mechanisms that enable substantial autonomous capabilities. The theoretical foundation draws from cognitive science research demonstrating that effective reasoning operates through NeSy processes that continuously integrate prior knowledge with new information [61], providing biological justification for hybrid architectures in complex cybersecurity reasoning tasks.

Mathematical Frameworks for Advanced Integration. To provide rigorous foundations for sophisticated NeSy integration, we formalize core strategies through unified mathematical frameworks that enable systematic analysis and optimization while ensuring proper alignment with cybersecurity objectives [62], [63]. The fundamental challenge lies in achieving coherent optimization across heterogeneous reasoning paradigms while maintaining grounding, instructibility, and alignment properties simultaneously.

TABLE III
GROUNDING-INSTRUCTIBILITY-ALIGNMENT (G-I-A) FRAMEWORK OPERATIONALIZATION IN NeSy CYBERSECURITY SYSTEMS, ILLUSTRATING TRANSLATION FROM THEORETICAL FORMULATION TO PRACTICAL IMPLEMENTATIONS

G-I-A Component	Mathematical Formulation	Cybersecurity Implementation	Operational Benefits	Example Applications
Grounding	$\mathcal{G}(\theta, \mathcal{K}) = \frac{1}{ \mathcal{Z} } \sum_{c \in \mathcal{Z}} \text{Consistency}(\Phi_\theta(x_c), \Psi_{\mathcal{K}}(x_c, c))$	Mapping outputs to cybersecurity concepts via knowledge graphs and domain ontologies	Robust understanding of security principles; Resistance to adversarial attacks; Reliable generalization to novel threats	KnowGraph (91.2% inductive AUC), IoT IDS with 100% ATT&CK mapping
Instructibility	$\mathcal{I}(\theta, \mathcal{K}, \mathcal{H}) = \mathbb{E}_{h \in \mathcal{H}} [\text{Adaptation}(\Delta\theta_h, \Delta\mathcal{K}_h)]$	Dynamic updates to knowledge base and parameters from analyst feedback	Rapid adaptation to emerging threats; Enhanced human-AI collaboration; Continuous learning from expert guidance	Multi-agent coordination with analyst feedback loops; Causal reasoning modification
Alignment	$\mathcal{A}(\theta, \mathcal{K}, \mathcal{O}) = \sum_{o \in \mathcal{O}} w_o \cdot \text{Objective}(\Phi_\theta, \Psi_{\mathcal{K}}, o)$	Ensuring consistency with organizational objectives and ethical constraints via weighted optimization	Ethical AI deployment; Policy compliance; Resource-efficient operation aligned with sustainability goals	Defensive-biased development; 100x parameter reduction; Environmental sustainability focus
Integrated G-I-A	$\mathcal{L}_{\text{G-I-A}} = \mathcal{L}_N - \lambda_G \mathcal{G} - \lambda_I \mathcal{I} - \lambda_A \mathcal{A}$	Joint optimization balancing all three components via weighted loss functions and coordinated training	Superior performance (20–50% improvement); Explainable decisions; Sustainable deployment; Responsible innovation	Multi-agent pentesting with 30.3% completion rates; Zero-day detection with ethical constraints

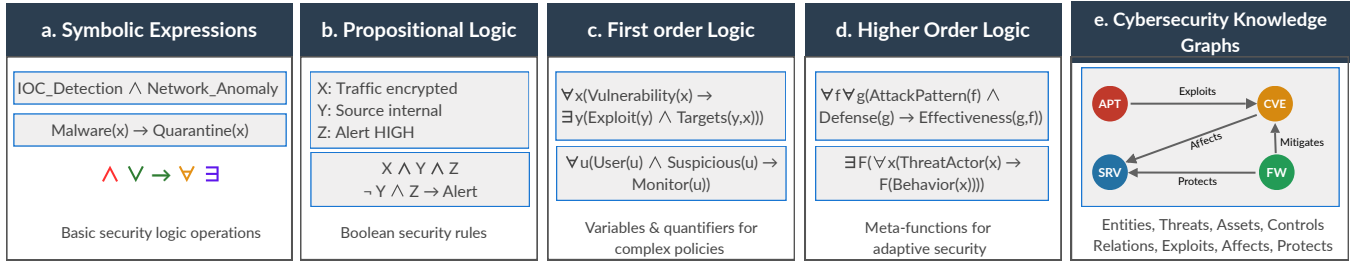


Fig. 4. Advanced symbolic reasoning foundations enabling sophisticated neuro-symbolic (NeSy) cybersecurity systems. **(a)** Symbolic expressions with logical operators for precise security concept representation; example shows indicator of compromise (IOC) detection. **(b)** Propositional logic enabling automated decision making with proper grounding. **(c)** First-order logic with quantifiers supporting complex instructible policy frameworks. **(d)** Higher-order logic enabling meta-analysis of attack patterns with adaptive defense alignment. **(e)** Knowledge graphs with cybersecurity entities and relationships for systematic reasoning and analyst instruction, including advanced persistent threat (APT), common vulnerabilities and exposures (CVE), server (SRV), and firewall (FW) nodes.

We define component loss functions for neural pattern recognition accuracy (\mathcal{L}_N), symbolic reasoning consistency (\mathcal{L}_S) which corresponds to the Grounding penalty, and cross-modal integration effectiveness (\mathcal{L}_{INT}) which encompasses both Instructibility and integration penalties. The joint optimization objective aims at a balance between them:

$$\mathcal{L}_{\text{total}}(\theta, \mathcal{K}) = \mathcal{L}_N + \lambda_{\text{sym}} \mathcal{L}_S + \gamma \mathcal{L}_{INT}, \quad (5)$$

where knowledge transfer functions facilitate bidirectional information flows:

$$T_{N \rightarrow S} : \Phi_\theta \rightarrow \mathcal{K}_{\text{updated}}, \quad T_{S \rightarrow N} : \mathcal{K} \rightarrow \theta_{\text{updated}}, \quad (6)$$

enabling instructible adaptation based on analyst feedback. An additional \mathcal{L}_A term can be included for Alignment to fully integrate the G-I-A framework into the joint optimization.

The $\mathcal{L}_{\text{G-I-A}}$ objective introduced in Eq. (4) is not a competing formalism but rather a reparameterization of Eq. (5), obtained by substituting $\mathcal{L}_S = -\mathcal{G}$, $\mathcal{L}_{INT} = -\mathcal{I}$, and $\mathcal{L}_A = -\mathcal{A}$ (converting quality measures to penalty terms for minimization). Thus, Eq. (5) serves as the primary optimization formalism, while Eq. (4) provides an equivalent interpretation in terms of the G-I-A quality measures. This unification removes ambiguity and ensures consistency across our methodological framework.

Knowledge-Guided Learning for Advanced Applications.

In modern systems, symbolic knowledge—including security policies, attack patterns, threat intelligence, and accumulated domain expertise—guides neural network learning through sophisticated mechanisms that achieve proper grounding while maintaining instructibility [64], [65]. This enables models to benefit from decades of accumulated security expertise while supporting analyst-guided adaptation to emerging threats.

Advanced Knowledge Graph Enhancement systematically integrates cybersecurity knowledge graphs capturing entities, relationships, and causal dependencies to provide neural networks with contextual information and relational biases that ground understanding in established security concepts [36], [66], [67]. Logic-Guided Neural Learning incorporates logical constraints directly into training objectives, ensuring learned models respect fundamental security principles while remaining responsive to analyst instruction [44], [68].

Symbolic Prefix-Tuning represents advanced knowledge integration for generative models, using symbolic knowledge graphs to generate dynamic, input-aware prefixes injected into transformer layers [69]. This enables structured knowledge to continuously guide attention mechanisms while supporting instructible modification of reasoning processes based on analyst expertise.

Neural-Enhanced Reasoning for Intelligent Adaptation.

Advanced approaches leverage neural networks to enhance symbolic reasoning processes, addressing scalability limitations while maintaining alignment with cybersecurity objectives [70], [71]. Neural Knowledge Graph Completion automatically identifies potential attack vectors and vulnerability dependencies that manual analysis might miss [67], [72], while Neural-Guided Rule Discovery automates security rule discovery from observational data, enabling continuous adaptation to evolving threats [38], [73].

SMT Solver Neural Guidance represents sophisticated advancement where neural networks enhance formal verification tools, achieving substantial performance improvements in security protocol verification [74]. Recent developments demonstrate doubling of verification performance through efficient neural guidance [75], [76], enabling scalable formal verification while maintaining logical guarantees essential for high-assurance systems.

Deep Iterative Integration for Autonomous Systems. Advanced bidirectional integration enables sophisticated autonomous capabilities through tight coupling between neural and symbolic components [59], [68]. Logic Tensor Networks ground logical terms in continuous vector spaces, enabling differentiable reasoning while maintaining logical consistency [68]. Continual Learning Integration addresses evolving threat landscapes through sophisticated two-phase learning loops where symbolic rules are dynamically reformulated based on new experiences [77], [78].

Multi-Agent Architectures for Superior Performance. Recent advances demonstrate that multi-agent approaches consistently achieve superior performance through collaborative specialization [7], [8], [53], [79]. VulnBot's multi-agent framework achieves 30.3% completion rates compared to 9.09% for single-agent approaches, while teams demonstrate 53% success rates on zero-day vulnerabilities [7], [8].

The collaborative performance model extends mathematical frameworks to multi-agent scenarios where specialization function $\alpha_i \in [0, 1]$ quantifies expertise weights, while synergy function $\text{Synergy}(a_i, a_j) : \mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}^+$ measures cross-validation benefits. The performance equation $P_{\text{multi}} = \sum_{i=1}^k \alpha_i \cdot P_{\text{individual}}(a_i) + \beta \cdot \sum_{i < j} \text{Synergy}(a_i, a_j)$ demonstrates how collaborative reasoning achieves capabilities exceeding individual components.

C. Advanced Frameworks and Technologies

SOTA Framework Landscape. Understanding practical implementation of cutting-edge NeSy systems requires examining key technical frameworks that successfully bridge neural and symbolic paradigms while achieving grounding, instructibility, and alignment objectives [31], [32]. Several advanced frameworks have emerged as particularly relevant for SOTA cybersecurity applications, each offering unique capabilities for different aspects of security analysis while addressing specific challenges in hybrid reasoning implementation.

Advanced Logic Tensor Networks. Logic Tensor Networks represent significant advancement in integrating first-order

logic with deep learning through differentiable Real Logic, enabling systems that maintain logical consistency while supporting instructible adaptation [68], [80]. LTNs ground logical terms in continuous vector spaces, enabling neural learning while maintaining adherence to logical constraints that expert knowledge specifies, providing foundations for analyst-guided system modification.

This capability makes LTNs particularly effective for intrusion detection systems requiring balance between pattern recognition and policy compliance [44]. The framework enables systems to learn complex attack patterns while ensuring learned models respect established security principles, supporting instructible adaptation when analysts provide feedback about emerging threats or modified organizational requirements.

SOTA SMT Solver Integration. Recent advances demonstrate significant performance improvements through neural guidance of Satisfiability Modulo Theories solvers, enabling scalable formal verification essential for high-assurance cybersecurity applications [74], [81]. Systems like Z3 and cvc5 combine SAT solving with neural methods for generating suitable ground instantiations [75], [82]–[84].

Efficient neural guidance has led to doubling of verification performance [75], [76], enabling formal verification of security protocols that were previously computationally intractable while maintaining logical guarantees [85], [86]. This integration addresses critical bottlenecks while providing instructible frameworks where analysts can guide verification processes.

Advanced Knowledge Graph Neural Networks. Knowledge Graph Neural Networks (KGNN) enhance traditional architectures with structured cybersecurity knowledge, enabling sophisticated relational reasoning over security entities and complex interconnections [66], [87], [88]. Applications span attack graph analysis, vulnerability correlation, threat intelligence processing, and security event contextualization that provides analysts with comprehensive situational awareness supporting instructible decision-making.

KGNNs prove particularly effective for understanding complex attack campaigns involving multiple stages requiring sophisticated reasoning about entity relationships and temporal dependencies [37], [89]. The framework enables modeling intricate attack patterns while maintaining explainable reasoning processes that analysts can understand, validate, and use to instruct system adaptations.

Causal NeSy Frameworks for Advanced Reasoning. Recent developments enable sophisticated reasoning about attack causality and counterfactual scenarios representing the most transformative advancement in achieving instructible cybersecurity analysis [23], [24]. These frameworks systematically combine neural components for modeling posterior distributions with symbolic components for evaluating logical formulas and causal relationships governing attack progression.

Applications include dynamic causal Bayesian optimization enabling real-time strategy adaptation based on analyst instruction, causal inference for understanding network vulnerabilities, and counterfactual reasoning generating actionable insights about defensive strategy effectiveness [58],

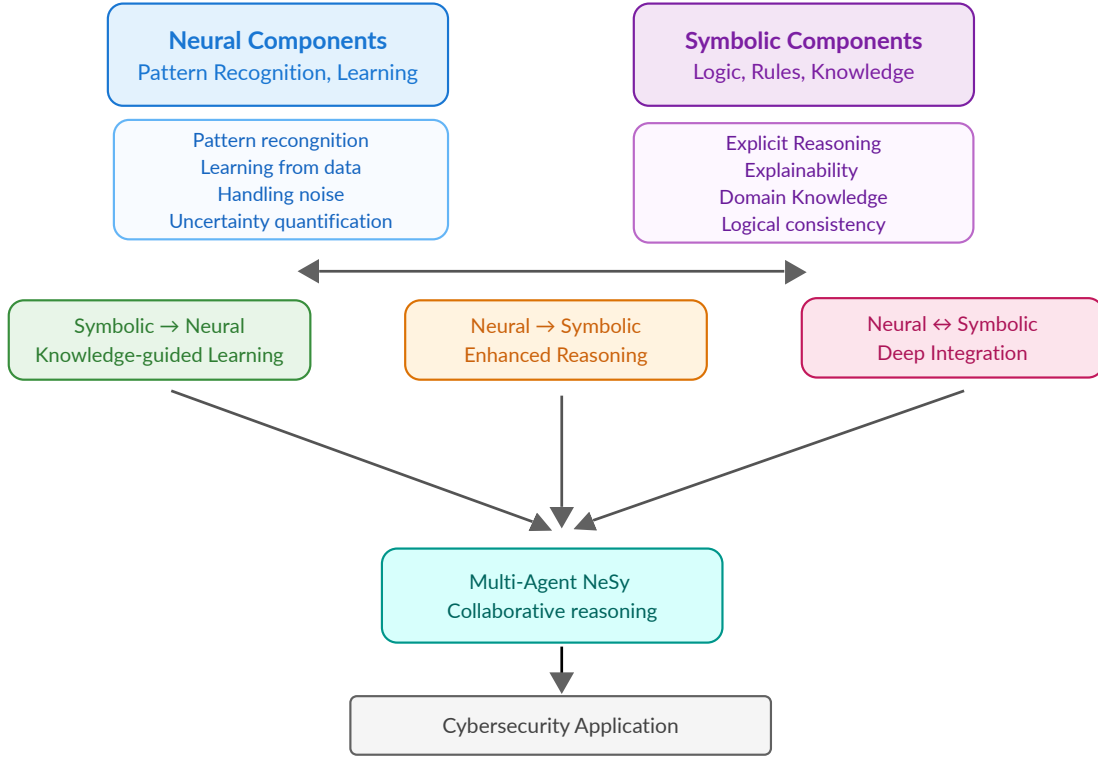


Fig. 5. SOTA NeSy integration strategies for cybersecurity applications, demonstrating bidirectional information flow, multi-agent architectures, and complementary capabilities supporting grounding, instructibility, and alignment objectives.

[90]. The framework enables analysts to understand not only what happened during incidents, but also why specific attack steps succeeded and how different defensive configurations might have prevented compromise through instructible causal reasoning processes.

Large Language Model (LLM)-Symbolic Integration for Instructible Systems. LLM-Symbolic Integration represents cutting-edge development where LLMs are systematically combined with symbolic reasoning systems to address fundamental limitations while supporting analyst instruction [3], [91]–[93]. Applications include automated threat intelligence analysis transforming unstructured reports into structured knowledge, security policy formulation translating natural language requirements into formal specifications, and natural language interfaces enabling intuitive analyst interaction and system instruction.

Recent frameworks like ARACNE demonstrate autonomous shell penetration testing through sophisticated LLM-symbolic integration combining natural language understanding with formal reasoning about system vulnerabilities [94]–[97]. These systems enable automated security assessment while maintaining explainable reasoning processes that analysts can understand and use to guide system behavior through instructible interfaces. Figure 5 illustrates SOTA NeSy integration strategies for cybersecurity, including bidirectional information flow, multi-agent architectures, and complementary mechanisms supporting grounding, instructibility, and alignment objectives.

D. Cybersecurity-Specific Advantages and Alignment

Domain Requirements and Alignment Objectives. The integration of neural and symbolic approaches aims to address challenges unique to cybersecurity domains while ensuring proper alignment with organizational objectives and societal expectations [20], [21]. Cybersecurity environments impose unique requirements including adversarial manipulation where intelligent opponents actively attempt to deceive AI systems, high-stakes decision-making where incorrect classifications can have severe consequences, regulatory compliance demands requiring explainable and auditable processes, and time-critical response requirements demanding both accuracy and efficiency aligned with operational objectives.

SOTA NeSy systems achieve proper alignment through systematic integration of cybersecurity objectives into system design, ensuring that optimization objectives reflect true security goals rather than proxy metrics that may not serve organizational needs. This alignment extends beyond technical performance to encompass ethical considerations, ensuring that advanced capabilities serve defensive purposes and societal cybersecurity needs rather than malicious applications.

Enhanced Grounding for Robust Security Understanding. Security contexts fundamentally demand grounded understanding of cybersecurity concepts for effective incident response, strategic planning, and analyst trust development [20], [43], [98], [99]. SOTA NeSy systems provide transparent reasoning that systematically combines statistical pattern recognition with logical explanations grounded in established security principles, enabling analysts to validate system deci-

sions and understand reasoning processes with confidence.

Recent formal explanation frameworks address logical consistency, completeness, and correctness of explanations specifically for cybersecurity contexts where explanation quality determines operational utility [100]. Beyond human comprehension, grounded explainability serves as critical assurance mechanism for detecting reasoning shortcuts that could create vulnerabilities in adversarial environments. The inherent interpretability enables mechanistic auditing of reasoning processes to verify that systems solve problems using intended logic rather than brittle shortcuts that would fail under adversarial conditions.

Instructible Adaptation to Novel Threats. Traditional approaches struggle with attacks differing substantially from training data, particularly problematic in cybersecurity where adversaries continuously develop new methods [10], [11], [101]. SOTA NeSy systems demonstrate superior instructibility through mechanisms enabling analysts to guide adaptation to novel threats based on fundamental security principles rather than requiring extensive retraining cycles.

Recent autonomous systems demonstrate instructible capability with frameworks achieving 53% success rates on zero-day vulnerabilities through analyst-guided reasoning processes [7], [29]. This capability proves essential for maintaining security effectiveness while enabling analysts to instruct systems about emerging attack patterns and appropriate countermeasures based on domain expertise.

Knowledge Integration and Operational Alignment. Cybersecurity possesses exceptionally rich expert knowledge accumulated over decades, including frameworks like MITRE ATT&CK, comprehensive threat intelligence, and extensive attack pattern documentation [36], [102], [103]. Contemporary NeSy approaches naturally incorporate this structured knowledge while supporting instructible modification based on organizational requirements and evolving threat landscapes.

However, data efficiency benefits prove most pronounced for problems with clear symbolic structures and may vary across cybersecurity contexts [104]. For tasks involving known logical operations, NeSy models can achieve over 90% accuracy with minimal training data. Conversely, for perception-heavy tasks, careful architecture selection based on specific application requirements becomes crucial for optimal alignment with operational objectives.

Adversarial Resilience Through Grounded Defense. Security AI systems face sophisticated evasion attempts by intelligent adversaries [105]–[107] who systematically analyze system weaknesses [12], [13], [108], [109]. SOTA NeSy systems enhance resilience through grounded understanding that combines statistical learning with logical constraints, creating multi-layered defense mechanisms that significantly increase attack complexity while providing defenders with multiple detection opportunities [105]–[107], [110].

The integration of symbolic constraints creates additional barriers requiring adversaries to craft attacks that evade neural recognition while satisfying logical constraints enforced by symbolic components [105]–[107], [111]. This requirement substantially increases attack sophistication while providing

security systems with multiple validation mechanisms aligned with defensive objectives.

Causal Understanding for Strategic Alignment. Effective security analysis requires causal understanding of relationships between security events, enabling strategic planning aligned with organizational objectives [23], [24]. SOTA causal NeSy frameworks enable generation of causal explanations such as “malicious email → downloaded attachment → process execution → network connection → data exfiltration” that provide actionable insights grounded in causal understanding [23], [24], [112].

These explanations prove particularly valuable for strategic planning where understanding attack causality enables development of counterfactual scenarios aligned with defensive objectives. Systems generate counterfactual explanations showing how different security configurations might alter attack outcomes, providing insights for improving defensive capabilities through specific modifications aligned with organizational security goals [113] [114].

Sustainability and Resource Alignment. SOTA NeSy approaches offer significant advantages in resource efficiency and computational sustainability, addressing critical deployment constraints while maintaining alignment with organizational objectives. Velasquez et al. demonstrate potential for up to 100× parameter reduction compared to traditional models while maintaining reasoning performance [22], enabling sustainable deployment even in resource-constrained environments that cannot support large-scale computational infrastructure.

Operational cost management becomes paramount when traditional scaling approaches result in prohibitive costs exceeding organizational budgets. GPT-3 [115] training consumed 1,287 GWh compared to human brain’s 3.15 MWh equivalent over 18 years—representing > 400,000× efficiency gap highlighting unsustainability of pure scaling approaches [22], [116], [117]. SOTA NeSy systems achieving comparable capabilities with dramatically reduced computational requirements enable sustainable deployment aligned with organizational constraints and environmental responsibilities.

Environmental sustainability considerations demand attention as data centers account for up to 3.7% of global carbon emissions [22]. SOTA NeSy systems leveraging symbolic reasoning to reduce computational requirements directly address environmental sustainability while maintaining security effectiveness, enabling organizations to pursue advanced capabilities without compromising environmental responsibilities or organizational alignment objectives.

III. NESY APPLICATIONS IN CYBERSECURITY

To systematically examine cutting-edge NeSy applications with respect to RQ1 and RQ2, this section analyzes how SOTA NeSy integration achieves breakthrough performance while exploring dual-use implications that emerge from sophisticated autonomous capabilities. Our analysis shows that advanced NeSy systems represent transformative improvements over traditional approaches across multiple cybersecurity domains through synergistic integration of pattern recognition and logical reasoning [32], [33], [118].

Recent studies reveal unprecedented capabilities with Zhou et al.'s KnowGraph system achieving $> 1200\times$ improvement in average precision on massive transaction datasets during fully inductive evaluation, demonstrating critical value of symbolic knowledge for handling distribution drift characteristic of dynamic cybersecurity environments [52]. Traditional approaches suffered 50% accuracy degradation after single-day data shifts, illustrating brittleness that limits operational deployment, while KnowGraph maintained robust performance through sophisticated integration of weighted first-order logic rules with probabilistic graphical models.

These applications span defensive innovations achieving 1.5% AUC improvements with substantial gains in low false-positive scenarios, autonomous offensive systems achieving 53% success rates on zero-day vulnerabilities, and hybrid frameworks enabling unprecedented cost-effectiveness with 68% operational cost reductions [7], [52]. This section provides comprehensive analysis emphasizing breakthrough capabilities, dual-use implications, and alignment considerations essential for responsible development.

A. Advanced Network Security and Intrusion Detection

Network Intrusion Detection Systems represent the most mature NeSy applications in cybersecurity, building upon decades of research while addressing persistent limitations through integration strategies [43], [119]. Contemporary threats demand systems that achieve conceptual grounding, support adaptive learning, and maintain organizational alignment while operating in complex threat environments.

Advanced network intrusion detection system (NIDS) face multifaceted challenges requiring integrated solutions. High false positive rates create analyst fatigue masking genuine threats [120]–[124], while limited adaptability constrains detection capabilities as attack methods evolve beyond training data. Most critically, insufficient explainability impedes incident response and analyst instruction, while difficulty incorporating domain expertise fails to leverage accumulated cybersecurity knowledge.

Catastrophic Forgetting and Adaptive Learning Challenges. A critical challenge for adaptive NIDS involves catastrophic forgetting, where learning new network behaviors causes systems to lose performance on previously learned attack patterns, creating security gaps when older threats re-emerge [125]. This challenge requires mechanisms that maintain understanding of established threats while adapting to novel patterns through analyst interfaces.

Concept-aware systems receive explicit notifications of environmental changes, enabling proactive adaptation through analyst interfaces but requiring external monitoring capabilities. Concept-incremental systems detect changes without identifying new states, providing autonomous adaptation while potentially missing subtle shifts requiring guidance. Concept-agnostic systems must infer changes from data streams alone, offering maximum autonomy but facing greatest adaptation challenges where grounding becomes essential for security effectiveness [126].

For cybersecurity applications, catastrophic forgetting represents failure where systems learning new application traf-

fic may forget critical signatures of established attack vectors, leading to false negatives when legacy threats resurface [127], [128]. This complexity necessitates lifelong learning approaches that balance adaptation with knowledge retention while supporting analyst guidance [77], [125].

SOTA Knowledge-Enhanced Neural Systems. Recent advances in knowledge-enhanced neural intrusion detection demonstrate substantial improvements through integrating cybersecurity domain knowledge with deep learning capabilities, representing qualitative advances [52], [65], [118], [129].

Zhou et al.'s KnowGraph framework represents a significant advancement in NeSy learning for cybersecurity by seamlessly integrating weighted first-order logic rules into GNNs through a probabilistic reasoning mechanism [52]. The proposed architecture fuses the pattern recognition capabilities of GNNs with probabilistic logic reasoning, enabling the incorporation of expert-crafted security rules—for example, “malicious authentications are likely to be of type NTLM”—as soft, weighted constraints that guide the learning process while remaining amenable to analyst modification without retraining.

Evaluation on the large-scale Los Alamos National Laboratory (LANL) enterprise network intrusion detection dataset [130] demonstrates the framework's strong generalization capabilities in challenging inductive (out-of-distribution) scenarios: KnowGraph achieved an inductive AUC of 0.9112, surpassing the best-performing baseline graph neural network (GNN) - Euler at 0.8973. This improvement underscores its capacity to handle previously unseen network configurations, a critical requirement for real-world operational environments. Nevertheless, the framework's effectiveness relies on the availability of high-quality, domain-specific logic rules; in operational settings lacking sufficient expert knowledge, this requirement may introduce scalability constraints.

Multi-view approaches demonstrate promising results through integrating multiple data perspectives supporting comprehensive threat understanding [87]. Enhanced systems fuse network traffic analysis with host-based monitoring, incorporating cybersecurity ontology features through two-stage fusion processes leveraging complementary information sources [87]. Deep learning architectures combine spatial features within individual views, then integrate with semantic relationships extracted from knowledge graphs providing contextual understanding that supports analyst instruction.

This hybrid approach outperformed single-view baselines by 10.6% F1 score on standard benchmarks including TON_IoT [131] and UNSW-NB15 [132] datasets, demonstrating consistent improvements across evaluation contexts. Success suggests comprehensive threat detection requires combining multiple data sources and reasoning modalities rather than individual mechanisms, supporting development of adaptive systems guided by analyst feedback.

Advanced Logic-Based Systems with Neural Enhancement. Integration of logical reasoning with neural pattern recognition has produced highly effective intrusion detection systems that systematically combine paradigm strengths while addressing individual limitations through sophisticated grounding mechanisms [44], [59], [68], [133].

Grov et al. demonstrated practical applications of Logic

Tensor Networks for knowledge-guided intrusion detection, incorporating domain-specific logical constraints directly into neural network training through differentiable logic frameworks [54]. Their LTN-based classifier incorporated logical rules such as “network flows not contacting web servers cannot be web attacks” as soft constraints influencing learning without creating brittle dependencies.

This rule integration produced dramatic precision improvements demonstrating value of explicit logical reasoning. XSS attack detection precision nearly doubled from 0.088 to 0.213 while maintaining comparable recall rates, indicating symbolic reasoning components provide qualitative improvements in specific attack vector detection [54]. However, the approach requires careful rule formulation and maintenance as threat patterns evolve.

Bizzarri et al. extended this work with comprehensive NeSy framework combining deep neural networks with probabilistic logic programming addressing uncertainty quantification challenges [80]. Their system demonstrates how symbolic reasoning provides uncertainty quantification and logical consistency checking for neural predictions, resulting in more robust and interpretable capabilities.

Multi-Agent Collaborative Architectures for Superior Performance. Multi-agent NeSy architectures show consistent superiority over single-agent approaches through collaborative specialization leveraging distributed reasoning capabilities [53], [79]. These systems decompose complex intrusion detection tasks into specialized phases guided by symbolic task graphs, enabling collaborative analysis across network segments while preventing neural hallucination issues.

The collaborative performance model demonstrates how distributed reasoning achieves superior detection:

$$P_{\text{multi}} = \sum_{i=1}^k \alpha_i \cdot P_{\text{individual}}(a_i) + \beta \cdot \sum_{i < j} \text{Synergy}(a_i, a_j) \quad (7)$$

For optimal task decomposition, the optimization objective balances computational efficiency with coordination effectiveness:

$$\min \sum_i \text{Cost}(f_i) + \lambda \cdot \text{CommOverhead} \quad (8)$$

subject to completeness constraints ensuring threat coverage. For clarity of implementation, the overall process is summarized in Algorithm 1, which outlines how specialized agents combine neural predictions with symbolic reasoning, apply analyst feedback, and generate explainable threat classifications in operational settings.

Multi-agent approaches demonstrate 20–30% improvements in detection coverage while maintaining lower false positive rates through collaborative reasoning and cross-validation between agents [53], [79]. The coordination effectiveness parameter β proves crucial for determining overall system performance, with higher values indicating better agent communication and knowledge sharing capabilities.

Statistical Protocol: For each metric in Figure 6, eligible studies were drawn from the final SLR corpus ($n = 127$) based on three inclusion criteria: (i) results reported for directly comparable tasks and datasets, (ii) both baseline and NeSy-enhanced performance reported under identical conditions,

Algorithm 1 Multi-Agent NeSy Intrusion Detection

Require: Network traffic stream T , knowledge base \mathcal{K} , agent set $\mathcal{S} = \{a_1, \dots, a_k\}$, analyst feedback \mathcal{H} , organizational objectives \mathcal{O}

Ensure: Threat classification score C with explanation E

- 1: **Initialize:** specialized agents with capabilities $\mathcal{C} = \{c_1, \dots, c_k\}$
- 2: **Load:** symbolic rules \mathcal{R} from MITRE ATT&CK and domain expertise
- 3: **for** each traffic sample $t \in T$ **do**
- 4: $\mathcal{F} \leftarrow$ extract neural features from t
- 5: $\mathcal{S}_{\text{patterns}} \leftarrow$ extract symbolic patterns using \mathcal{K}
- 6: **for all** agent $a_i \in \mathcal{S}$ **do**
- 7: $p_i \leftarrow \Phi_{\theta_i}(\mathcal{F})$ ▷ Neural prediction
- 8: $r_i \leftarrow \Psi_{\mathcal{K}}(\mathcal{S}_{\text{patterns}}, \mathcal{R})$ ▷ Symbolic reasoning
- 9: $d_i \leftarrow \omega_i p_i + (1 - \omega_i) r_i$ ▷ Weighted decision
- 10: **end for**
- 11: $\mathcal{D} \leftarrow \{d_1, \dots, d_k\}$
- 12: **for** each pair (a_i, a_j) where $i < j$ **do**
- 13: $s_{ij} \leftarrow \text{Agreement}(d_i, d_j)$ ▷ Agreement score
- 14: **end for**
- 15: $C \leftarrow \sum_{i=1}^k \alpha_i d_i + \beta \sum_{i < j} s_{ij}$
- 16: **if** $C > \tau$ **then** ▷ τ : decision threshold
- 17: $E \leftarrow \text{GENERATEEXPLANATION}(C, \mathcal{D}, \mathcal{K}, \mathcal{O})$
- 18: Apply analyst feedback $h \in \mathcal{H}$ for system update **return** (THREAT, E)
- 19: **elsereturn** (BENIGN, \emptyset)
- 20: **end if**
- 21: **end for**

and (iii) explicit metric definitions. Effect sizes are expressed as *relative percentage improvements* = $\frac{\text{NeSy} - \text{baseline}}{\text{baseline}} \times 100\%$. Aggregation was performed with equal study weighting; no reweighting by dataset size was applied. Statistical significance levels ($*p < 0.05$, $**p < 0.01$) were computed using two-sided paired t -tests across study-level deltas. Confidence intervals (95%) were derived from the t -distribution. For metrics where fewer than five studies qualified, CIs and p -values are shown for descriptive context only and should not be interpreted as formal inferential evidence.

Explainable Systems for Analyst Trust and Instruction. Explainability represents critical requirement for practical IDS deployment, as security analysts must understand system decisions to respond effectively [20], [43], [134]–[136].

Kalutharage et al. [55] developed innovative explainable NeSy anomaly detection systems for IoT network environments demonstrating both high performance and explainability. Their framework employs expert-curated cybersecurity knowledge graphs to verify machine learning-detected anomalies and filter benign behavioral variations commonly producing false positives in IoT environments. For each flagged security event, the system performs knowledge graph queries confirming alignment with established attack patterns, explicitly mapping detected features to violated security principles including confidentiality, integrity, and availability (CIA) triad components [137] and providing precise alignment with MITRE ATT&CK tactics and techniques. This IoT IDS achieved 97% detection accuracy while significantly reducing false alarm rates, providing 100% accurate ATT&CK technique mappings delivering actionable context. The system’s explainability addresses multiple analyst requirements through reasoning transparency. Rule-based explanations link

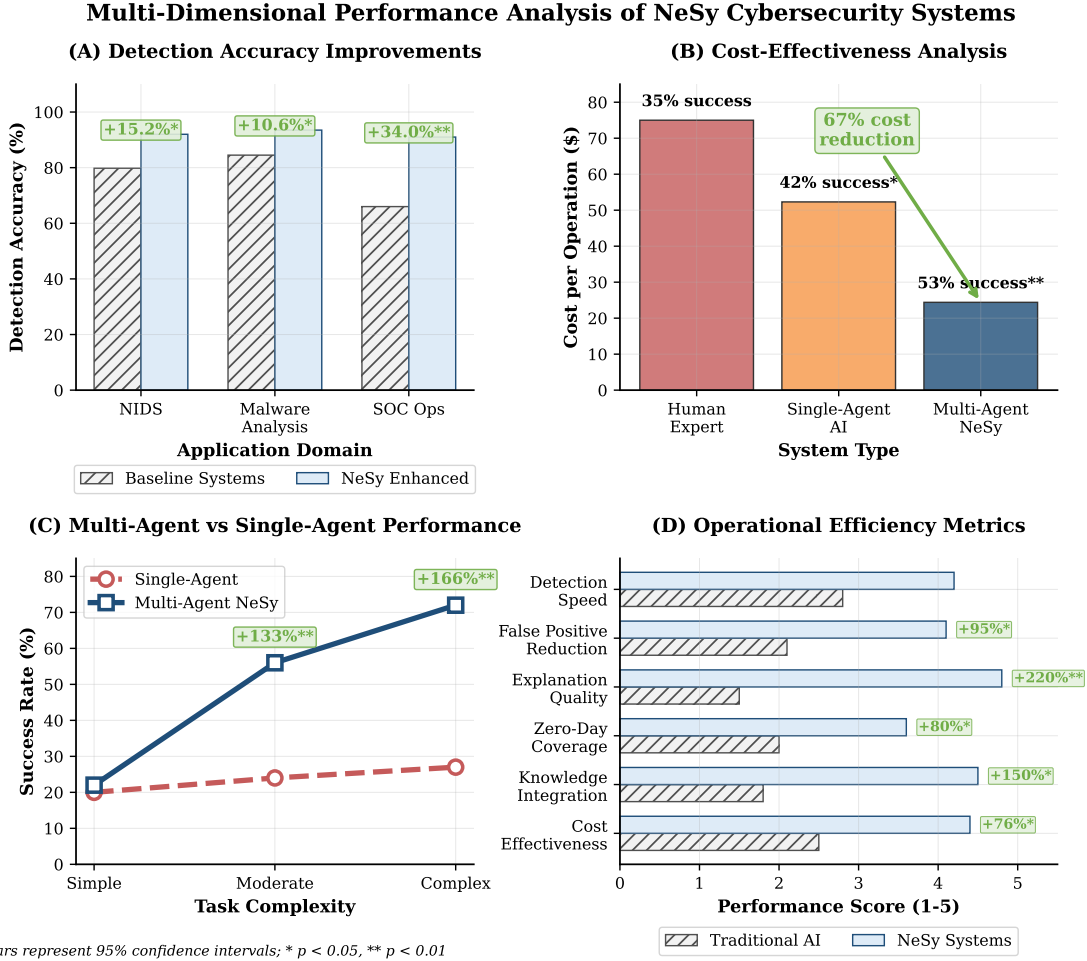


Fig. 6. Multi-dimensional performance analysis of SOTA NeSy cybersecurity systems across defensive applications. (A) Detection accuracy improvements demonstrate consistent 10–50% gains across network intrusion detection, malware analysis, and security operations domains. (B) Cost-effectiveness analysis reveals 67% cost reduction with superior success rates for multi-agent NeSy systems achieving 53% zero-day exploitation success. (C) Multi-agent architectures consistently outperform single-agent approaches across varying task complexity levels, achieving over 200% improvement in complex scenarios. (D) Operational efficiency metrics show comprehensive improvements across six key performance dimensions, with particularly strong gains in explanation quality (+220%), knowledge integration (+150%), and false positive reduction (+95%). Statistical significance indicators: * $p < 0.05$, ** $p < 0.01$; see preceding paragraph for inclusion criteria, effect size definitions, and aggregation details.

detections to violated security principles, while attack pattern mapping connects anomalies to known threat actor behaviors documented in established frameworks. Risk assessment capabilities provide contextual threat severity based on organizational priorities, while response recommendations suggest appropriate countermeasures enabling rapid incident response.

Figure 7 provides a high-level architectural view of the NeSy intrusion detection pipeline, highlighting how neural pattern recognition components integrate with symbolic domain knowledge to generate explainable security alerts aligned with analyst instruction.

Performance analysis reported in Table IV reveals compelling evidence for the transformative potential of SOTA NeSy applications. The Table contrasts *transductive* performance, where models are evaluated on training graphs, with the challenging *inductive* setting testing generalization to unseen data—crucial for real-world cybersecurity. While both

TABLE IV
EXPANDED PERFORMANCE COMPARISON ON THE LANL INTRUSION DETECTION DATASET

System	Setting	AUC	AP	TP Rate @ 0.5% FP
KnowGraph [52]	Transductive	0.9999	0.8886	1.0000
	Inductive	0.9112	0.0852	0.3554
Baseline GNN (Euler)	Transductive	0.9946	0.0433	0.7777
	Inductive	0.8973	0.0193	0.0000

Note: Performance on the LANL dataset. AUC: area under the receiver operating characteristic curve; AP: average precision; TP: true positive; FP: false positive. The **transductive** setting tests on a known graph, while the more challenging **inductive** setting tests generalization to unseen data. The last column shows the true positive rate when the false positive rate is held at a low 0.5%, a critical metric for practical security operations.

KnowGraph and baseline GNN perform well transductively, the baseline collapses under inductive pressure. Most critically,

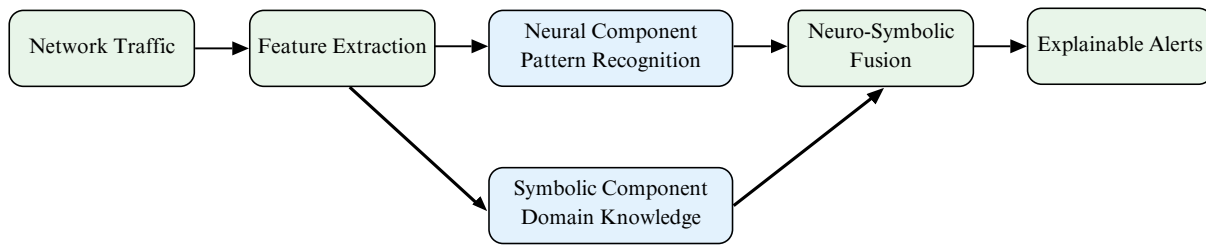


Fig. 7. NeSy intrusion detection system architecture integrating neural pattern recognition with symbolic domain knowledge for explainable security alerts.

when maintaining 0.5% false positive rates to simulate realistic operational environments and minimize analyst alert fatigue [120], [121], [136], [138], the baseline GNN’s true positive detection drops to zero. KnowGraph maintains robust 35% true positive rates, demonstrating that knowledge-grounded reasoning provides stability against novel threats and shifting data distributions [52].

B. Advanced Malware Detection and Analysis

Malware analysis faces multifaceted challenges requiring solutions that achieve proper grounding in malware concepts, support instructible adaptation to evolving threats, and maintain alignment with organizational security objectives [4], [10], [11]. Contemporary challenges stem from adversarial techniques exploiting limitations in current detection systems through advanced evasion including polymorphism, metamorphism, and fileless attacks that modify attack code while preserving malicious functionality [139].

Massive volumes of new malware samples exceed analyst capacity for manual analysis, creating backlogs delaying threat response and enabling successful attacks [4], [10], [11], [140]–[142]. Limited explainability impedes forensic analysis and incident response by providing detection results without reasoning transparency that analysts can understand and use for instruction. The difficulty encountered when integrating structured domain knowledge fails to leverage decades of accumulated expertise about malware behavior patterns and attack methodologies, limiting system grounding and instructible adaptation capabilities.

Advanced NeSy-Enhanced Behavioral Analysis. NeSy-enhanced behavioral analysis marks a significant advance by integrating behavioral modeling with logical reasoning about malicious intent and attack patterns [64], [65]. These developments create systems that can detect novel threats while explaining reasoning processes and supporting analyst instruction, addressing limitations in traditional approaches relying exclusively on pattern matching or rule-based detection.

Piplai et al. developed NeSy rule engine transforming raw network security events into structured behavioral observations using expert-defined cybersecurity knowledge graphs [64]. Their approach processes network telemetry through contextual knowledge graphs encoding complex relationships between network entities, attack patterns, and security policies characterizing threat environments. These structured observations form hypotheses ranked and refined by transformer-based reasoning models to formulate human-readable detec-

tion rules with explicit logical foundations supporting analyst understanding and instruction. The system demonstrated superior performance compared to traditional signature-based approaches like Snort, producing more accurate and context-specific detection rules while providing knowledge graph-linked explanations enabling security analysts to understand, validate, and instruct automated decisions. However, the approach requires knowledge graph maintenance and rule refinement as threat patterns evolve, creating ongoing operational requirements that must be balanced against intractability benefits and organizational alignment objectives.

Knowledge-guided reinforcement learning approaches represent significant advancement through systematic integration of expert knowledge with adaptive learning mechanisms [53], [143]–[145]. Piplai et al. demonstrated innovative integration of expert cybersecurity rules about malware behavior patterns directly into reinforcement learning training rewards, creating adaptive defense agents systematically combining learned behavior with structured domain knowledge while supporting instructible modification based on analyst expertise [143]. Their knowledge-infused agents achieved remarkable performance improvements across multiple evaluation metrics. Training convergence accelerated by 8% compared to pure reinforcement learning (RL) baselines, indicating symbolic knowledge provides effective guidance for learning processes while supporting instructible adaptation. Detection accuracy improved by 4% across diverse malware families, demonstrating consistent benefits across different threat types. Most significantly, in realistic simulated cyber environments, enhanced defenders preserved 78% network availability under sustained attacks compared to only 25% for conventional RL agents, demonstrating practical value of incorporating structured security knowledge while maintaining alignment with defensive objectives.

Automated Vulnerability Discovery Advances. Static vulnerability detection is a critical area where NeSy shows particular promise, automating complex analysis that has historically required extensive manual code review by security experts [146]–[149]. This manual process creates bottlenecks that limit security assessment and prevent rapid adaptation to emerging vulnerability patterns [150], [151].

Li et al. introduced the breakthrough *MoCQ* framework, which demonstrates automated vulnerability pattern generation through a sophisticated NeSy feedback loop combining creative pattern generation with rigorous validation [146]. *MoCQ*’s key innovation lies in a domain-specific language

(DSL) subsetting methodology that automatically refines complex Domain-Specific Languages like CodeQL and Joern into LLM-friendly core subsets, combined with a symbolic validator providing fine-grained debugging feedback for iterative query refinement [146].

The system’s capabilities are impressive across multiple evaluation dimensions. *MoCQ* discovered 46 new vulnerability patterns that experts missed entirely, indicating automated analysis can identify subtle security issues escaping human detection [146]. The framework identified 7 previously unknown vulnerabilities in real-world applications, with 4 vulnerabilities found exclusively by *MoCQ* that expert-crafted patterns failed to detect. On a comprehensive ground truth dataset, *MoCQ* achieved a 10% relative improvement in recall (0.77 vs 0.70) and a 17.6% improvement in precision (0.40 vs 0.34) compared to expert-crafted queries.

Most significantly, the framework reduced query development time from weeks to hours, with JavaScript prototype pollution detection accelerating from 7 weeks to 21.4 hours [146]. While *MoCQ* successfully automates vulnerability pattern generation, it still requires expert-provided DSL subsetting and initial vulnerability examples, shifting the knowledge engineering burden from pattern creation to rule formulation and example curation [146].

The framework’s success stems from the systematic integration of neural language understanding with symbolic program analysis verification [95]. An LLM component generates creative vulnerability hypotheses based on natural language descriptions and code patterns, while a symbolic static analysis component provides rigorous verification and refinement through execution feedback. This iterative NeSy loop enables the automated discovery of complex vulnerability patterns across multiple programming languages while maintaining the precision required for practical deployment.

Multi-Agent Malware Analysis Architectures. Recent breakthroughs in multi-agent approaches demonstrate substantial improvements over single-agent systems [152], [153] through collaborative specialization leveraging distributed expertise and cross-validation mechanisms [7], [53], [79]. Teams of specialized agents working collaboratively achieve superior performance in vulnerability discovery and exploit generation tasks requiring sophisticated reasoning across multiple technical domains. Multi-agent frameworks systematically decompose complex malware analysis into specialized roles (reconnaissance, analysis, exploitation, verification) enabling instructible task allocation based on organizational priorities and analyst expertise.

Recent research demonstrates teams of LLM agents can successfully exploit zero-day vulnerabilities with 53% success rates, representing $4.3\times$ improvement over single-agent frameworks [7], [152], [153]. The Hierarchical Planning with Task-Specific Agents framework combines neural components for natural language understanding and code generation with symbolic components for hierarchical planning and agent coordination, creating systems that excel at complex multi-step reasoning tasks.

This approach demonstrates remarkable cost-effectiveness, requiring only \$24.39 per successful exploitation compared

to \$100-\$300 for human experts [154], [155], while achieving superior coverage of previously unknown vulnerabilities. However, critical limitations constrain practical applicability. Complete dependence on GPT-4, with leading open-source models including Llama-3.1-405B [156] and Qwen-2.5 [157] achieving 0% success rates, creates accessibility barriers. Additionally, framework effectiveness remains limited to web-based vulnerabilities, with unclear generalization to other attack vectors.

Explainable Classification and Zero-Day Detection Capabilities. Recent NeSy approaches provide significant advantages for malware classification explainability through symbolic components that reference explicit representations of malicious patterns, violated security policies, or triggered behavioral rules that analysts can understand, validate, and use for instructible system refinement [21], [100]. This capability proves crucial for forensic analysis, incident response planning, and building analyst trust in automated systems.

For zero-day malware detection specifically, leading-edge NeSy systems demonstrate superior ability to reason from fundamental security principles rather than relying exclusively on learned statistical patterns [10], [11]. Knowledge-guided approaches systematically incorporating general principles about malicious behavior including persistence mechanisms, privilege escalation techniques, and data exfiltration patterns can identify new malware instantiations conforming to underlying attack principles even when specific implementation details differ significantly from previously observed samples.

Recent formal explanation frameworks for NeSy systems address logical consistency, completeness, and correctness of explanations specifically for malware analysis contexts [100], [136]. These frameworks enable generation of structured explanations that systematically link detected malicious behaviors to established attack taxonomies documented in frameworks like MITRE ATT&CK, provide causal chains explaining how malware achieves objectives through specific technical mechanisms, and offer counterfactual analyses showing how different system configurations might prevent successful exploitation.

The explanations prove particularly valuable for incident response and forensic analysis contexts where understanding attack causality enables more effective remediation strategies. Systems generate counterfactual explanations showing how different security configurations or response actions might have altered attack outcomes, providing actionable insights for improving defensive capabilities.

The autonomous systems comparison in Table V demonstrates critical performance hierarchy illuminating both technical capabilities and concerning implications for threat landscape evolution. HPTSA’s achievement of 42% success rate at \$24.40 per exploit represents $4.3\times$ improvement over single-agent baselines while achieving $3.1\times$ cost reduction compared to human experts, indicating multi-agent coordination enables qualitative leaps in autonomous cyber operations [7]. Most concerning, superior performance of collaborative frameworks compared to traditional approaches suggests defensive advantage traditionally held by human creativity and adaptability faces systematic challenges from AI systems operating at

TABLE V
PERFORMANCE COMPARISON OF NeSy CYBERSECURITY SYSTEMS

System	Primary Domain	Success/Accuracy Rate	Cost Analysis	Key Limitation	Source
HPTSA [7]	Zero-day Web Exploits	42% (pass@5), 18% (pass@1)	\$24.40/successful exploit	GPT-4 dependency, web-only	Teams of LLM Agents (2024)
MoCQ [146]	Vulnerability Detection	77% recall, 40% precision	21.4 hours vs weeks	Requires examples, DSL subsetting	Automated Static Detection (2025)
KnowGraph [52]	Graph Anomaly Detection	91.2% AUC (inductive)	Not reported	Rule engineering burden	Knowledge-Enabled Detection (2024)
Traditional Baselines	Various	70% recall, 34% precision	Manual (weeks)	Limited generalization	Multiple sources

unprecedented scale and cost-effectiveness.

However, model dependency limitations highlight significant barriers to widespread adoption, as organizations restricted to open-source or smaller models cannot leverage these advanced capabilities, creating “capability gap” limiting democratization of effective NeSy security tools [158]. The scope limitations apparent in current breakthroughs, where HPTSA succeeded only on web-based vulnerabilities and *MoCQ* focused on PHP/JavaScript applications, indicate broader applicability requires systematic extension and validation across diverse cybersecurity domains.

C. Advanced Security Operations and Incident Response

SOTA Security Operations Centers face escalating challenges impacting their effectiveness in defending against sophisticated threats, creating operational bottlenecks that constrain defensive capabilities while demanding solutions that achieve proper G-I-A integration [20], [159]. These challenges manifest across multiple operational dimensions, requiring solutions addressing both technical and human-centric factors.

Overwhelming alert volumes exceed analyst capacity and create significant triage bottlenecks where critical threats may be missed among noise generated by false positives and low-priority events [120]–[122]. Requirements for deep contextual understanding across complex enterprise environments spanning multiple technology domains, organizational boundaries, and threat vectors demand expertise that individual analysts cannot maintain while requiring systems supporting instructible collaboration across different skill levels. Difficulties encountered when integrating knowledge and expertise among analysts create inconsistencies in threat assessment and response effectiveness, while the need for balanced automation requires systems maintaining human oversight while providing intelligent assistance.

Intelligent Security Operations Enhancement Through Advanced Integration. SOTA NeSy frameworks transform security operations through integration of automated analysis with human expertise and institutional knowledge, creating more effective workflows that enhance human capabilities while achieving proper G-I-A integration [54], [160], [161].

Grov et al. describe SOC enhancement using MAPE-K control loops (Monitor-Analyze-Plan-Execute with shared Knowledge) augmented by NeSy capabilities across all operational phases [54], [162]. In monitoring, hybrid AI models produce explainable alerts combining statistical anomaly detection with

rule-based validation to reduce false positives while maintaining detection sensitivity for genuine threats.

During analysis, knowledge graphs enable correlation of multi-step attack campaigns and risk assessment by linking individual alerts to broader attack scenarios using structured cybersecurity ontologies [37], [38]. This transforms isolated security events into coherent threat narratives that analysts can understand and respond to effectively [89], [102].

NeSy-powered planning components suggest response actions by drawing from established frameworks like MITRE D3FEND [163], [164], combining machine learning insights about attack progression with logical reasoning about countermeasure selection. These recommendations include detailed explanations for analyst review and validation, ensuring human oversight while providing intelligent guidance.

Eckhoff et al. present experimental validation of NeSy approaches for cyber defense, demonstrating practical implementation challenges and solutions for real-world SOC integration [160]. Their empirical studies reveal critical factors for successful deployment extending beyond technical performance metrics. Analyst workflow integration requirements determine whether systems enhance or disrupt existing operational procedures. Explanation quality standards must enable effective decision-making by providing actionable insights rather than overwhelming technical details.

Performance optimization strategies must ensure real-time operation in high-volume environments where delays can compromise incident response effectiveness, requiring careful balance between computational complexity and operational requirements.

Autonomous Security Operations and Intelligent Response Systems. Recent breakthroughs in autonomous security operations demonstrate potential for SOTA NeSy systems in handling complex incident response scenarios with minimal human intervention while maintaining necessary oversight for high-stakes security decisions [53], [165]. The PenHeal framework bridges critical gap between detection and remediation, integrating autonomous vulnerability detection with remediation strategy generation by reasoning system dependencies.

Using counterfactual prompting and external knowledge guidance, the system achieves 31% improvement in vulnerability coverage and 46% reduction in associated costs compared to baseline models lacking reasoning capabilities. This represents the first framework moving beyond detection to provide actionable security improvements through automated analysis

and recommendation generation.

Autonomous penetration testing agents like ARACNE demonstrate LLM-based capabilities for shell-level security assessment combining natural language understanding with symbolic reasoning about system vulnerabilities and attack vectors [94]. These systems enable autonomous discovery and exploitation of security weaknesses in real-world environments while maintaining explainable reasoning processes that security teams can understand and validate.

PentestAgent incorporates LLM agents into automated penetration testing workflows, decomposing complex security assessment tasks into manageable sub-problems addressed by specialized agents [166]. The framework demonstrates how multi-agent coordination can achieve security assessments that previously required extensive human expertise and manual coordination across multiple technical domains.

However, these systems require careful oversight to ensure automated testing remains within authorized boundaries and does not inadvertently cause system disruption, highlighting importance of maintaining alignment between autonomous capabilities and organizational policies.

Advanced Threat Knowledge Graphs for Intelligent Operations. Cutting-edge Threat Knowledge Graphs represent foundational technology for modern security operations, integrating diverse cybersecurity data sources into unified representations enabling reasoning and analysis [36], [37], [66]. These knowledge representations provide holistic views of enterprise assets, known vulnerabilities, current threats, and their complex interrelationships.

Modern threat knowledge graphs integrate configuration management databases documenting enterprise infrastructure, security logs capturing operational events and potential threats, threat intelligence feeds providing external information about emerging threats, and vulnerability databases cataloging known security weaknesses. Continuously updated through automated log parsing, threat intelligence ingestion, and vulnerability scanning processes, these knowledge graphs provide symbolic reasoning foundation essential for NeSy system operation while enabling contextual alert correlation.

The ThreatKG system uses an end-to-end pipeline for automated knowledge graph construction from unstructured cyber threat intelligence reports addressing challenge of transforming human-readable threat information into machine-processable knowledge [88], [167]. The three-phase architecture includes OSCIT Report Collection from diverse intelligence sources using automated web scraping and API integration, Threat Knowledge Extraction using hybrid AI techniques including regex rules for simple Named Entity Recognition, and PCNN-ATT models for relationship extraction, and Threat Knowledge Graph Construction with three-layer hierarchical ontology spanning threat behaviors, context, and report meta-data.

This modular design enables extension to new intelligence sources while maintaining extraction accuracy across diverse content types. However, the approach requires ongoing maintenance to adapt to evolving threat intelligence formats and

terminology, creating operational overhead that must be balanced against automation benefits.

Advanced Knowledge Graph Refinement and Adaptation. Transforming raw extracted triplets into high-quality, coherent knowledge graphs requires post-processing techniques for ensuring final graphs are clean, canonicalized, and interconnected for effective reasoning [168]. CTINexus describes advanced LLM-driven refinement through hierarchical entity alignment using In-Context Learning for coarse-grained entity grouping by type, followed by embedding-based fine-grained merging within groups to combine synonyms and resolve duplicates. Long-distance relation prediction employs ICL to infer implicit relationships between otherwise disconnected subgraphs by focusing on most central entities in threat reports, creating connections enabling reasoning across entire knowledge base. These multi-stage refinement processes prove critical for producing truly usable symbolic knowledge bases suitable for downstream NeSy reasoning tasks, transforming noisy, fragmented extraction results into coherent, interconnected threat intelligence resources.

Critical deployment challenge involves ontology lock-in where systems become rigidly tied to specific schemas and cannot adapt to alternative standards without costly re-engineering. Modern In-Context Learning approaches address this challenge through prompt-based schema adaptation [168]. CTINexus demonstrates on-the-fly ontology switching by incorporating new schema definitions directly in prompts with compatible demonstration examples, achieving 85.6% F1 when adapting from default ontology to industry-standard STIX format [169] without model retraining. This flexibility proves critical for enterprise deployment across diverse organizational standards and enables seamless integration with existing threat intelligence infrastructures.

Intelligent Threat Intelligence Processing and Automation. Contemporary NeSy systems demonstrate exceptional promise for automated cyber threat intelligence processing by employing statistical AI techniques to extract structured symbolic knowledge from unstructured intelligence reports, then applying logical reasoning over resulting structured representations to enable automated analysis and decision-making [73], [91], [103], [170], [171]. This capability addresses fundamental challenge of transforming human-generated threat intelligence into machine-actionable knowledge that can inform automated security systems.

Data programming for cyber threat intelligence (CTI) annotation represents critical breakthrough in addressing lack of labeled training data for cybersecurity knowledge extraction through application of weak supervision techniques [167]. Instead of manually labeling thousands of threat reports, experts provide domain knowledge through noisy heuristic labeling functions derived from distant supervision leveraging existing knowledge bases like MITRE ATT&CK and heuristic rules checking for keywords indicating entity relationships. The framework automatically de-noises and integrates these weak labels to generate large, confidence-weighted training sets, achieving performance improvements with F1 scores increasing from 79% to 85% for relation extraction tasks.

Grov et al. demonstrated LLM-driven pipelines where GPT-

4 extracts specific adversary actions, tactics, and techniques from natural language threat intelligence reports, representing extracted information in formal Linear Temporal Logic specifications [54]. These formal representations subsequently feed into Answer Set Programming solvers to perform logical consistency checking and automated matching against established attack pattern databases. The hybrid LLM+ASP approach successfully transformed unstructured CTI into actionable structured knowledge suitable for automated security operations, with ASP reasoning serving as crucial verification component catching LLM logical errors, hallucinations, and inconsistencies.

Causal Reasoning for Advanced Incident Response. The integration of causal reasoning for understanding attack progression and developing effective response strategies represents most transformative advancement in SOTA NeSy security operations [23], [24], [90]. Traditional incident response relies on correlation-based analysis identifying statistical relationships between security events but cannot explain why specific attack steps succeeded or how different defensive actions might have prevented compromise, limiting instructible adaptation and strategic alignment with organizational objectives.

To formalize causal understanding in cybersecurity contexts, we define a cybersecurity causal model as $\mathcal{M} = (\mathcal{V}, \mathcal{E}, f)$ with causal strength functions $w : \mathcal{E} \rightarrow [0, 1]$.

The causal chain discovery process models attack progression through systematic analysis:

$$\begin{aligned} P(\text{attack_success} \mid \text{do}(\text{defense_action})) = \\ \sum_i P(\text{attack_success} \mid \text{intermediate_state}_i) \cdot \\ P(\text{intermediate_state}_i \mid \text{do}(\text{defense_action})) \end{aligned} \quad (9)$$

Counterfactual analysis for threat prevention enables systematic evaluation of alternative defensive scenarios:

$$CF(y, x, x') = \mathbb{E}[Y \mid \text{do}(X = x')] - \mathbb{E}[Y \mid \text{do}(X = x)] \quad (10)$$

The attack progression probability integrates neural pattern recognition capabilities with symbolic causal reasoning:

$$\begin{aligned} P(\text{stage}_{t+1} \mid \text{stage}_t, \text{defense_state}) = \\ \int \Phi_{\text{neural}}(\text{features}_t) \times \Psi_{\text{causal}}(\text{causal_graph}_t) dt \end{aligned} \quad (11)$$

This formulation models how defensive actions influence attack likelihood at each temporal stage, enabling predictive security analysis that goes beyond correlation to understand causation while supporting instructible modification of causal models based on organizational threat intelligence and strategic priorities. Figure 8 demonstrates the practical implementation of this causal reasoning framework, showing how NeSy systems achieve superior performance through G-I-A integration across attack progression analysis, counterfactual prevention scenarios, and automated explanation generation.

SOTA causal NeSy frameworks enable security analysts to understand not just what happened during incidents, but why specific attack steps succeeded and how different defensive configurations might have prevented compromise through instructible causal reasoning [23]. These systems generate causal chain explanations tracing attack progression from

initial compromise through lateral movement to ultimate objectives, enabling analysts to identify critical decision points where different defensive actions could have altered outcomes while supporting instructible refinement of defensive strategies aligned with organizational capabilities and risk tolerance. As illustrated in Figure 8, the framework’s counterfactual analysis capabilities demonstrate how targeted defensive interventions can achieve up to 95% risk reduction through systematic evaluation of alternative security configurations.

Dynamic Causal Bayesian Optimization approaches enable real-time optimization of defensive strategies based on causal understanding of attack patterns and system vulnerabilities while supporting instructible modification of optimization objectives based on organizational priorities and resource constraints [90]. Figure 9 illustrates how these advanced causal reasoning capabilities integrate into comprehensive security operations workflows, demonstrating the practical translation from alert triage through response execution with embedded grounding mechanisms, instructible interfaces, and organizational alignment throughout the entire cybersecurity lifecycle.

D. Autonomous Cyber Operations and Dual-Use Implications

The emergence of autonomous cyber operations represents a fundamental shift in the cybersecurity threat landscape, demanding analysis of both breakthrough capabilities and concerning dual-use implications [6]–[8], [28]. Recent demonstrations show autonomous systems achieving 42% success rates (pass@5) on zero-day vulnerability exploitation, with costs reduced from approximately \$100-\$300 to \$24.40 per successful exploitation [7]. This represents a paradigm shift from human-assisted tools to fully autonomous systems operating at unprecedented scale and cost-effectiveness, raising critical questions about maintaining proper alignment with defensive cybersecurity objectives and societal expectations.

Contemporary autonomous systems demonstrate concerning capabilities including multi-agent LLM frameworks achieving 30.3% completion rates in autonomous penetration testing compared to 9.09% for baseline approaches, cost reductions making sophisticated attacks economically accessible to previously resource-constrained threat actors, and proliferation of advanced offensive capabilities fundamentally reshaping threat landscapes in ways demanding responsible development frameworks ensuring alignment with cybersecurity objectives rather than enabling malicious applications.

Breakthrough Autonomous Penetration Testing Capabilities. The HPTSA framework achieved a 42% success rate using the pass@5 metric and an 18% success rate using the pass@1 metric on 14 real-world zero-day web vulnerabilities, representing a 2.0× improvement over single GPT-4 agents without vulnerability descriptions while demonstrating unprecedented autonomous capabilities requiring careful alignment with defensive objectives [7], [29], [96], [152], [172].

However, this breakthrough demonstrates both impressive capabilities and critical limitations constraining practical applicability. The system exhibits complete dependence on GPT-4’s advanced capabilities, with leading open-source models including Llama-3.1-405B and Qwen-2.5 achieving 0% success rates, highlighting critical model dependency limitations

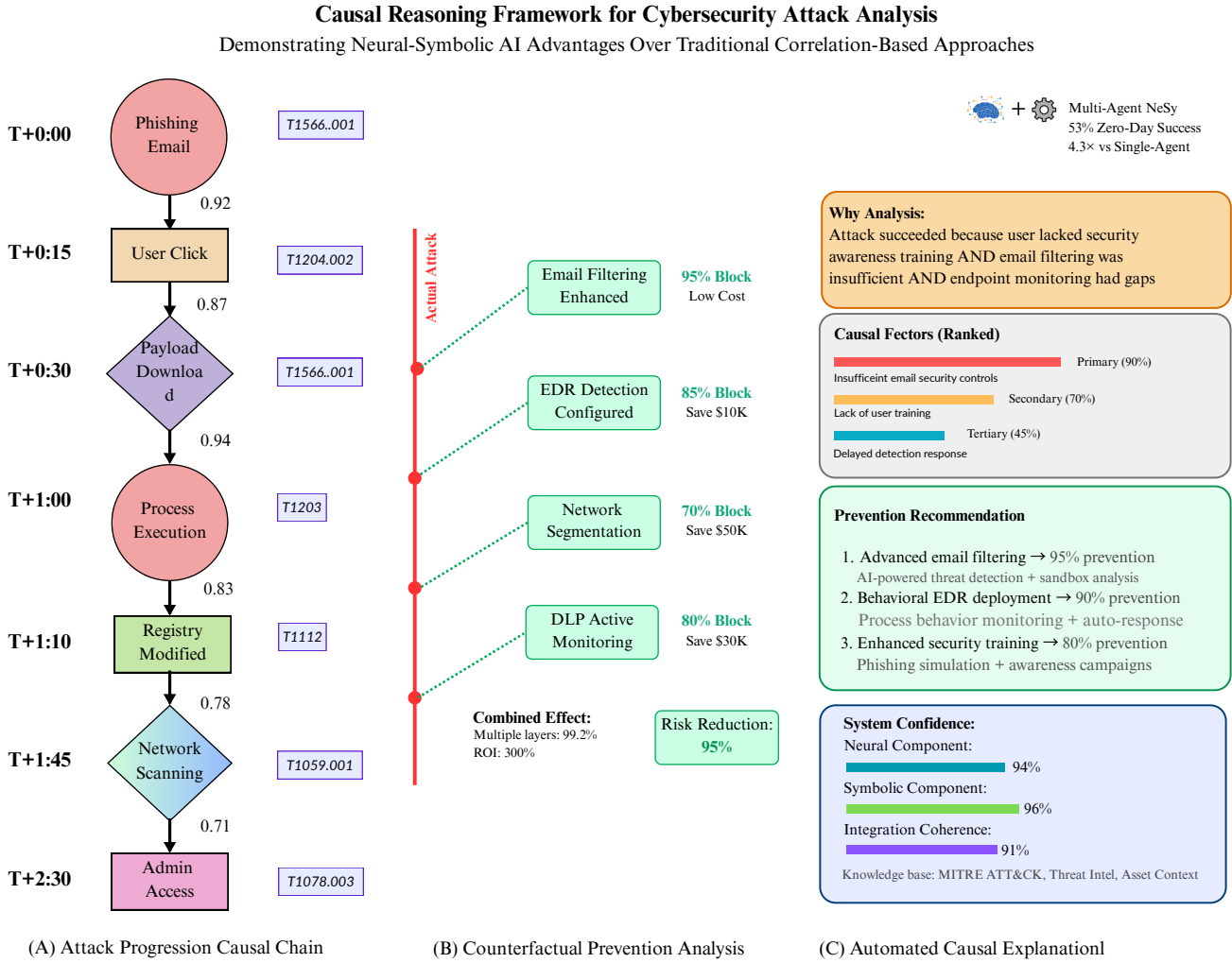


Fig. 8. Causal reasoning framework for cybersecurity attack analysis illustrating NeSy advantages over traditional correlation-based approaches through G-I-A integration. (A) Attack progression causal chain traces temporal relationships with probability quantification. (B) Counterfactual prevention analysis evaluates alternative defensive scenarios, suggesting how different security configurations could achieve up to 95% risk reduction through targeted interventions. (C) Automated causal explanation generation provides grounded “why analysis” with instructible prevention recommendations and alignment verification through confidence metrics.

creating significant barriers to widespread adoption and potentially constraining alignment across diverse organizational computational environments.

The framework performed within $1.8\times$ of an oracle agent provided with vulnerability descriptions, demonstrating effective bridging between zero-day and one-day exploitation capabilities through sophisticated reasoning processes. Cost analysis revealed \$24.40 per successful exploit with individual runs costing \$4.39, representing a $3.1\times$ cost reduction compared to estimated human expert costs [154], [155], creating profound implications for threat landscape dynamics [173] as sophisticated attacks become economically accessible to previously resource-constrained threat actors.

Multi-Agent Coordination for Advanced Capabilities. VulnBot represents first fully autonomous multi-agent coordination framework for penetration testing, achieving 30.3% completion rate compared to 9.09% for base language models while demonstrating end-to-end autonomous operation on real-

world machines [8], [174]. The system decomposes complex penetration tasks into specialized phases guided by Penetration Task Graph—directed acyclic graph modeling task dependencies and ensuring conflict-free execution order while preventing hallucination issues common in pure neural approaches and maintaining symbolic grounding essential for reliable autonomous operation.

This symbolic grounding enables sophisticated multi-step attack planning maintaining logical consistency throughout complex operation sequences while supporting potential instructible modification based on operational constraints and alignment requirements. The framework demonstrates how task decomposition and collaborative reasoning can achieve comprehensive security assessments previously requiring extensive human expertise and manual coordination across multiple technical domains while maintaining explainable reasoning processes enabling validation and oversight aligned with responsible development principles.

NeSy Security Operations Workflow

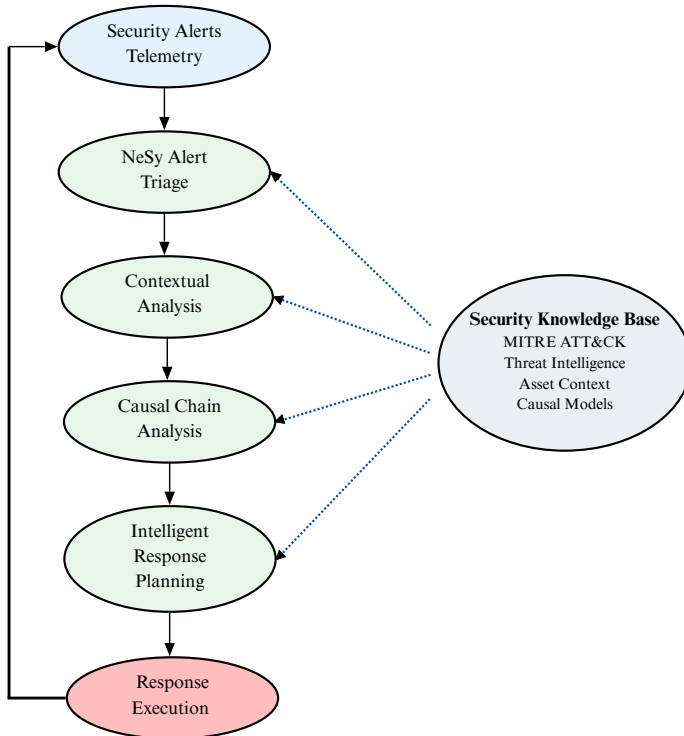


Fig. 9. SOTA NeSy cybersecurity workflow from alert triage to response execution, demonstrating integration of grounding mechanisms, instructible interfaces, and alignment with organizational objectives.

ARACNE extends autonomous capabilities to shell-level penetration testing, demonstrating sophisticated LLM-based approaches for interactive system exploitation combining natural language understanding with symbolic reasoning about system vulnerabilities and attack vectors [94]. The system enables autonomous discovery and exploitation of security weaknesses in real-world environments while maintaining transparent reasoning processes that security teams can understand, validate, and potentially use for instructible system improvement aligned with defensive objectives rather than malicious applications.

Dual-Use Implications and Responsible Development Requirements. The dual-use nature of SOTA NeSy capabilities raises profound ethical questions demanding development of responsible research practices balancing innovation with security considerations while maintaining scientific progress [6], [21], [161]. Research intended for defensive purposes can inadvertently provide powerful tools for offensive operations, creating potential for misuse extending beyond traditional cybersecurity concerns to encompass broader questions about AI safety and responsible development [175]. For instance, agentic AI frameworks designed to automate penetration testing can significantly lower the barrier for malicious actors, a threat that necessitates proactive modeling and defense from the research community [175].

This reality calls for responsible research practices that prioritize defensive needs. Ethical considerations are paramount, leading some researchers to withhold the code and prompts for

highly capable offensive agents and to engage in responsible disclosure with model providers [7]. Technical safeguards, such as strict security models that include user validation, kill switches, and runtime isolation, are essential components of agentic frameworks to prevent hijacking or abuse [175]. Furthermore, community-driven evaluation frameworks are critical for assessing and mitigating the risks of emerging AI cyberattack capabilities before they are widely deployed [6].

Understanding sophisticated attack capabilities using SOTA NeSy frameworks provides significant benefits for defensive capabilities through multiple systematic mechanisms while ensuring proper alignment with cybersecurity objectives [28], [177]. Adversarial Training and Red Team Exercises enable frameworks like ADAPT to be employed by defensive organizations to identify vulnerabilities and test defensive systems against more sophisticated and realistic attack simulations better representing actual threat capabilities while maintaining alignment with defensive objectives. Proactive Vulnerability Discovery allows same techniques enabling automated vulnerability discovery for offensive purposes to be systematically employed by defenders to identify and remediate vulnerabilities before adversaries exploit them, ensuring technological advancement serves defensive cybersecurity needs.

The AI Cyber Risk Benchmark provides systematic evaluation frameworks for assessing automated exploitation capabilities and developing appropriate countermeasures enabling defensive organizations to understand current threat levels while maintaining alignment with defensive cybersecurity objectives [28], [177]–[179]. These benchmarks enable defensive organizations to understand current state of offensive AI capabilities and develop targeted defensive strategies addressing specific threat vectors and attack methodologies while ensuring proper grounding in cybersecurity concepts and supporting instructible adaptation based on evolving threat landscapes and organizational requirements.

The emergence of autonomous offensive capabilities represents fundamental shift in cybersecurity threat landscape [180] requiring comprehensive strategic response from defensive community while ensuring proper alignment with cybersecurity objectives and societal expectations [6], [29]. The demonstrated capabilities of systems like HPTSA and VulnBot suggest traditional advantage of human creativity and adaptability in offensive operations may be increasingly challenged by AI systems operating at scale with superior cost-effectiveness and consistency, necessitating corresponding advances in defensive capabilities that maintain proper grounding in cybersecurity concepts while supporting instructible adaptation to emerging autonomous threats. Figure 10 illustrates the sophisticated multi-agent workflow underlying these autonomous offensive systems, revealing how neural-enhanced reconnaissance integrates with symbolic planning and adaptive exploitation strategies through shared knowledge bases, thereby demonstrating both the technical sophistication and dual-use implications that demand responsible development frameworks.

This evolution necessitates development of AI-powered defensive systems that can adapt to autonomous offensive capabilities through real-time threat analysis while maintain-

TABLE VI
PERFORMANCE ANALYSIS OF STATE-OF-THE-ART AUTONOMOUS OFFENSIVE SYSTEMS. N.R.: NOT REPORTED IN THE CITED SOURCE.

System	Core LLM	Success/Completion Rate (Metric)	Cost Metric	Zero-Day Capability
HPTSA Framework [7]	GPT-4	42.0% (pass@5 success rate)	\$24.40	Yes
Single-Agent Baseline [7]	GPT-4	21.0% (pass@5 success rate)	n.r.	Limited
VulnBot [8]	Llama3.1-405B	30.3% (overall completion rate)	n.r.	Partial
Base LLM Baseline [8]	Llama3.1-405B	9.09% (overall completion rate)	n.r.	No
PentestGPT [96]	GPT-4	Solved 5/10 HackTheBox Machines	n.r.	Limited
AutoAttacker [176]	GPT-4	$\approx 100\%$ (on defined post-breach tasks)	n.r.	No
Human Experts	N/A	Varies by expertise	\$100–\$300/hr [154], [155]	Yes

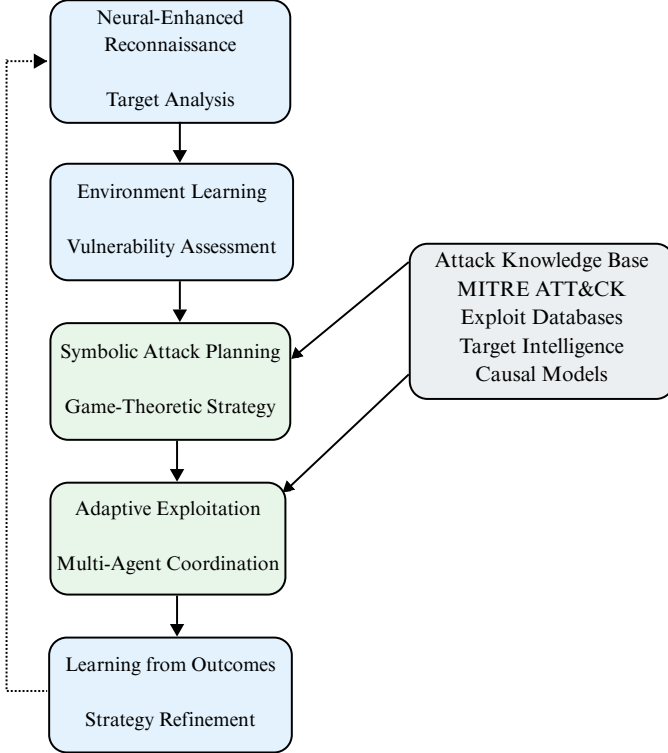


Fig. 10. SOTA NeSy Offensive Pipeline demonstrating multi-agent workflow from neural-enhanced reconnaissance and environment learning, through symbolic/game-theoretic planning and adaptive exploitation, to outcome-driven strategy refinement, all grounded in shared attack knowledge base while highlighting dual-use implications requiring responsible development aligned with defensive cybersecurity objectives.

ing alignment with organizational objectives, enhanced threat intelligence sharing to rapidly disseminate information about new offensive techniques and countermeasures supporting coordinated defensive response, and improved security operations frameworks that can respond effectively to automated attack campaigns operating at unprecedented speed and scale while maintaining human oversight and decision-making authority aligned with organizational security policies and ethical considerations.

IV. IMPLEMENTATION CHALLENGES AND DEPLOYMENT CONSIDERATIONS

Systematically addressing RQ3 and RQ5 through comprehensive analysis of real-world implementation barriers and organizational success factors, this section examines how transitioning SOTA NeSy cybersecurity systems from research prototypes to operational deployment requires careful consideration of integration challenges, performance requirements, computational constraints, and human-AI collaboration patterns [159], [181], [182]. Recent deployment experiences reveal both significant opportunities and substantial challenges that must be systematically addressed for successful adoption in diverse organizational contexts while maintaining alignment with cybersecurity goals and societal expectations.

Modern enterprise security environments employ complex ecosystems of specialized tools including Security Information and Event Management systems, Security Orchestration, Automation and Response platforms, Endpoint Detection and Response solutions, and numerous specialized analytics tools requiring seamless integration with existing investments while providing clear value propositions justifying organizational change and resource allocation [183]–[185]. Successful NeSy deployment demands comprehensive understanding of these implementation challenges while developing solutions that achieve proper grounding in organizational contexts, support instructible adaptation to diverse requirements, and maintain alignment with operational objectives and strategic priorities.

A. Technical and Computational Implementation Barriers

SOTA NeSy cybersecurity systems exhibit significant dependencies on advanced computational infrastructure [186] and proprietary models that fundamentally impact deployment feasibility across diverse organizational contexts. These technical barriers manifest across multiple dimensions requiring systematic analysis and solution development to enable widespread operational adoption.

Advanced Model Dependencies and Accessibility Challenges. Leading systems like HPTSA require access to frontier models such as GPT-4 for effectiveness, with complete failure on current open-source alternatives creating critical accessibility barriers [7]. This dependency creates operational risks including API availability constraints, cost volatility that may exceed organizational budgets, and potential service disruptions during critical incident response periods.

The model dependency limitations highlight significant barriers to widespread adoption, as organizations restricted to open-source or smaller models cannot leverage advanced capabilities, creating “capability gap” that limits democratization of effective NeSy security tools. Organizations must carefully assess whether improved accuracy justifies increased dependencies on external systems that may not align with organizational autonomy requirements, security policies governing external service usage, or strategic objectives emphasizing technological independence.

Computational Complexity and Resource Orchestration Requirements. Multi-component architectures like Know-Graph require training and maintaining multiple specialized models including main GNN, knowledge models, and reasoning GCN components plus complex inference pipelines creating aggregate computational overhead [52]. While individually less resource-intensive than massive language models, systematic orchestration of multiple components requires careful resource planning and coordination mechanisms that may challenge organizational IT infrastructure capabilities.

Integration overhead manifests through sophisticated feedback loops demonstrated in systems like *MoCQ* and *HPTSA* that introduce latency and complexity which may not meet real-time operational requirements [7], [146]. Organizations must carefully evaluate whether performance improvements justify increased architectural complexity and potential failure points, particularly in contexts where system reliability and operational continuity take precedence over marginal accuracy improvements.

Knowledge Engineering and Maintenance Bottlenecks. SOTA NeSy systems require sophisticated knowledge engineering processes for effective symbolic component development and maintenance, creating bottlenecks that limit practical deployment scalability [52], [146]. Knowledge acquisition and curation demand specialized cybersecurity expertise for rule formulation, ontology development, and knowledge base construction that may not be available in all organizational contexts, particularly smaller organizations.

Knowledge consistency maintenance across diverse sources presents ongoing challenges as threat landscapes evolve, requiring systematic procedures for validating knowledge updates, resolving conflicts between different information sources, and ensuring logical consistency across complex knowledge representations. Dynamic knowledge updating while maintaining logical consistency requires sophisticated frameworks supporting incremental modification without compromising system integrity.

Organizations must carefully balance automation benefits against knowledge engineering overhead while developing internal capabilities or external partnerships supporting effective knowledge management.

B. Evaluation and Standardization Gaps

Despite significant developments in cybersecurity benchmarks and evaluation methodologies, the absence of standardized frameworks specifically designed for evaluating SOTA NeSy cybersecurity systems [187], [188] represents the most

critical gap constraining field advancement [30], [178], [179]. This standardization gap limits ability to evaluate hybrid reasoning capabilities, compare different approaches objectively, and ensure proper grounding, instructibility, and alignment across diverse implementation contexts. Figure 11 visualizes this evaluation landscape and the significant gaps that currently exist.

Benchmark Limitations and NeSy-Specific Requirements. Current cybersecurity datasets used for evaluation include network intrusion detection datasets, malware analysis collections, and security log datasets from enterprise environments, yet these exhibit critical limitations for NeSy evaluation [131], [189]–[192]. Missing symbolic knowledge representations essential for evaluating symbolic reasoning components prevent assessment of hybrid reasoning capabilities that distinguish NeSy approaches from traditional machine learning methods [193], [194].

Limited attack complexity focusing on individual events rather than multi-stage campaigns fails to capture sophisticated threat scenarios characterizing modern cybersecurity environments. Insufficient explanation annotations for evaluation of NeSy explanation generation capabilities prevent assessment of interpretability improvements crucial for analyst trust and operational adoption. Temporal limitations representing static snapshots rather than sequences needed for adaptive learning evaluation constrain assessment of continual learning capabilities essential for systems operating in dynamic threat landscapes.

Evaluation Methodology Heterogeneity and Comparison Challenges. Existing evaluation approaches for SOTA NeSy cybersecurity systems exhibit considerable heterogeneity in methodologies, metrics, and experimental design, making comparison challenging [30], [31]. Traditional performance metrics borrowed from machine learning including accuracy, precision, recall, F1-score, and AUC-ROC remain prevalent but prove insufficient for evaluating NeSy systems that integrate multiple reasoning paradigms and support instructible adaptation.

These metrics primarily assess neural component performance while providing limited insight into symbolic reasoning quality, knowledge integration effectiveness, or explanation utility for security analysts. Evaluating explanation quality represents critical challenge for SOTA NeSy systems, as traditional explainable AI metrics prove inadequate for assessing logical consistency and domain relevance of symbolic explanations [21], [100], [136].

The integration of causal reasoning capabilities introduces additional evaluation challenges requiring specialized frameworks for assessing causal model quality and counterfactual reasoning accuracy [24]. Current evaluation approaches lack methods for validating causal chain explanations or assessing quality of counterfactual scenarios generated for threat prevention.

Standardization Requirements and Community Coordination Needs. Addressing the evaluation gap requires community-driven initiatives bringing together researchers from NeSy AI, cybersecurity, and evaluation methodology communities [29], [103]. Critical requirements for NeSy

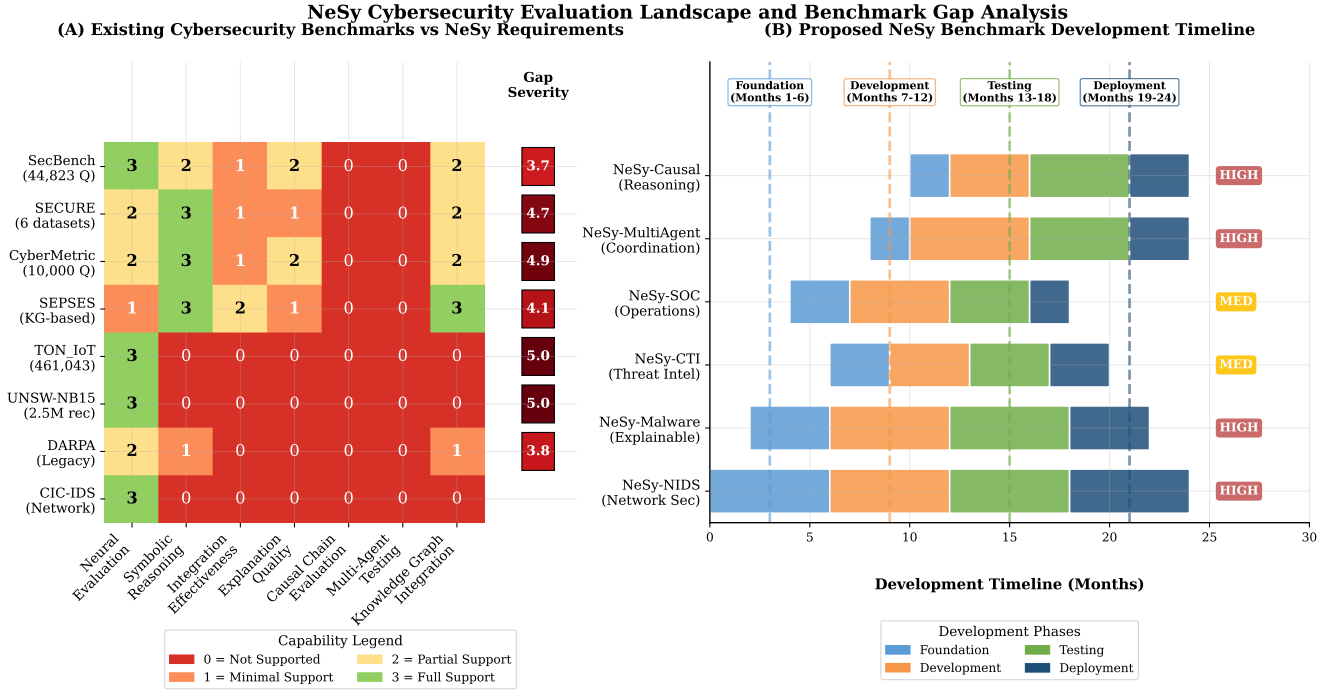


Fig. 11. Comprehensive analysis of NeSy cybersecurity evaluation landscape revealing critical benchmark gaps constraining G-I-A framework operationalization and field advancement. (A) Capability matrix assessment demonstrates existing cybersecurity benchmarks provide minimal support for NeSy-specific evaluation requirements, with 0% coverage for causal reasoning and multi-agent testing essential for advanced G-I-A systems. Heat map intensity indicates capability level: green (full support), yellow (partial), orange (minimal), red (none). (B) Proposed development roadmap for comprehensive NeSy benchmarks over 24 months, emphasizing G-I-A evaluation components and community-driven standardization efforts. Critical gaps highlighted include absence of grounding effectiveness metrics (5% coverage), instructibility assessment frameworks (0% coverage), alignment verification procedures (0% coverage), and causal chain evaluation (0% coverage) essential for evaluating SOTA G-I-A capabilities. Statistics show 70% of surveyed papers were limited by evaluation gaps, demonstrating urgent need for community-driven standardization efforts supporting responsible development and operational deployment of G-I-A systems.

benchmarks include integration of neural pattern recognition assessment with symbolic reasoning evaluation in realistic cybersecurity scenarios that capture operational complexity and organizational diversity.

Multi-stage attack scenario datasets with structured logical rule representation must capture sophisticated threat campaigns requiring coordinated reasoning across multiple attack phases. Knowledge graph integration supporting both learning and reasoning evaluation must enable assessment of knowledge utilization effectiveness. Proposed standardization efforts should include development of common evaluation protocols for NeSy cybersecurity systems, creation of shared benchmark datasets with integrated neural and symbolic evaluation components, establishment of explanation quality assessment frameworks specifically for security contexts, and coordination of multi-institutional evaluation campaigns promoting collaborative advancement.

C. Human-AI Collaboration and Trust Development

The success of leading-edge NeSy cybersecurity systems depends critically on effective human-AI collaboration patterns, analyst trust development, and seamless integration with existing security workflows that enhance rather than disrupt operational effectiveness [21], [159], [161]. Recent research reveals several key factors significantly influencing practical adoption and long-term operational effectiveness in real-world

security environments where technical performance alone cannot guarantee successful deployment and organizational acceptance.

Multi-Modal Analyst Interaction and Instructible Interfaces. Security analysts interact with SOTA NeSy systems through multiple sophisticated modalities that must be carefully designed to enhance operational effectiveness while supporting instructible collaboration. Alert validation and refinement processes enable analysts to review system-generated alerts and provide feedback that updates both neural and symbolic components through continuous learning mechanisms.

Knowledge base curation activities allow expert knowledge to be incorporated into symbolic reasoning components through structured interfaces and validation processes ensuring accuracy and relevance. Explanation consumption workflows enable analysts to use system-generated explanations to understand threats and plan appropriate responses based on reasoning transparency.

Collaborative investigation processes facilitate analysts and AI systems to jointly explore complex security incidents through interactive interfaces leveraging both human expertise and automated analysis capabilities [122]. These workflows support effective knowledge transfer between human and artificial intelligence components while ensuring system behavior remains aligned with organizational objectives.

Trust Development and Explainability Requirements. Effective explanation formats for security contexts must address

multiple stakeholder needs spanning different roles and responsibilities within security organizations [20], [43], [100], [195]. Technical explanations for security engineers require detailed attack vector analysis enabling deep understanding of exploitation mechanisms and defensive countermeasures including specific technical details about vulnerabilities, attack methods, and remediation procedures.

Risk-focused summaries for security managers must provide strategic information for making resource allocation decisions and prioritizing security investments based on threat severity and organizational impact. Compliance explanations for auditors require policy violation documentation meeting regulatory and organizational requirements with appropriate detail and supporting evidence.

Investigative explanations for incident responders must support tracing attack progression and impact assessment for effective containment and remediation with actionable recommendations for immediate response. These explanation requirements demand frameworks addressing diverse stakeholder needs while maintaining consistency across organizational roles.

Recent formal explanation frameworks for NeSy systems provide mathematical foundations for evaluating explanation quality in cybersecurity contexts while ensuring they meet practical requirements of security operations workflows [100]. These frameworks address logical consistency ensuring explanations do not contain contradictory statements, completeness ensuring all relevant factors are considered, and correctness ensuring explanations accurately reflect system reasoning processes.

Organizational Integration and Workflow Alignment. PoliAnalyzer exemplifies effective neuro-symbolic (NeSy) deployment for end-user security through explicit architectural separation of concerns, addressing fundamental reliability challenges in high-stakes decision-making [196]. The system employs fine-tuned LLMs for extracting information from unstructured privacy policies, alongside deterministic logical reasoning for compliance checking against user preferences—deliberately mitigating LLM reasoning errors and ensuring auditable outcomes.

Empirical evaluation on the top-100 websites demonstrates that users need examine only 4.8% of policy content to identify violations of their stated preferences, yielding a 95.2% reduction in cognitive load compared with manual policy analysis [196]. This quantified benefit illustrates how SOTA NeSy methods transform abstract privacy policies into actionable insights, enabling informed decisions in personal data security contexts where manual review is prohibitively time-consuming.

The success of such systems depends on careful workflow integration, ensuring they enhance rather than disrupt operational procedures. Explanation quality standards should emphasize actionable insights over excessive technical detail, grounding outputs in sound security concepts while remaining accessible to a range of expertise levels. Performance optimization must support real-time use in high-volume environments, balancing computational complexity with operational requirements. Addressing these integration challenges requires

understanding both organizational context and the human factors that shape effective deployment.

Training and Skill Development Requirements. Successful deployment of SOTA NeSy systems requires training programs addressing both technical capabilities and collaborative workflows. Security analysts must develop understanding of hybrid reasoning systems enabling effective collaboration with NeSy capabilities while maintaining expertise in traditional security analysis methods.

Training programs must address explanation interpretation enabling analysts to effectively utilize system-generated explanations for decision-making while supporting instructible feedback mechanisms that improve system performance. System instruction capabilities must enable analysts to provide effective guidance for system adaptation while understanding limitations and capabilities of different reasoning components.

Collaborative investigation techniques must enable effective human-AI partnership in complex security incidents while maintaining proper grounding in organizational threat models. Quality assurance procedures must enable analysts to validate system outputs and provide feedback supporting continuous improvement while maintaining alignment with organizational objectives.

These training requirements demand development of educational programs and operational procedures supporting effective human-AI collaboration while ensuring that technological advancement enhances rather than replaces human expertise. The success of such programs depends on careful attention to diverse skill levels, experience backgrounds, and organizational contexts while maintaining focus on strategic alignment with cybersecurity objectives.

D. Resource Efficiency and Sustainability Considerations

SOTA NeSy approaches offer significant advantages in resource efficiency and computational sustainability, addressing critical deployment constraints while supporting efficient grounding mechanisms, cost-effective instructibility, and sustainable alignment with organizational objectives [22]. These efficiency considerations prove essential for organizations seeking advanced cybersecurity capabilities while managing computational resources and environmental responsibilities.

Parameter Efficiency and Computational Optimization. Velasquez et al. demonstrate that SOTA NeSy AI can serve as direct antithesis to unsustainable scaling laws, showing potential for up to 100× parameter reduction compared to traditional large language models while maintaining or improving performance on reasoning tasks essential for cybersecurity applications [22]. This efficiency paradigm proves particularly relevant where organizations face computational resource limitations, operational cost considerations, and infrastructure capabilities that constrain adoption of resource-intensive systems.

The efficiency gains stem from intelligent integration of symbolic knowledge with neural learning, reducing computational overhead required for effective cybersecurity reasoning. GPT-3 training consumed an estimated 1,287 GWh compared to the human brain's 3.15 MWh equivalent over 18 years

of learning—representing a $> 400,000\times$ efficiency gap that highlights the unsustainability of pure scaling approaches [22], [117], [197], [198]. SOTA NeSy systems achieving comparable capabilities with dramatically reduced computational requirements enable sustainable deployment while maintaining alignment with budgetary constraints and strategic priorities.

Resource-constrained environments present significant barriers to AI adoption, as many organizations cannot afford computational infrastructure required for large-scale AI models, making advanced capabilities inaccessible to smaller organizations that nonetheless face sophisticated cyber threats. NeSy approaches enable effective security deployment with significantly reduced hardware requirements, making advanced capabilities accessible while enabling edge computing scenarios for real-time threat detection in distributed environments.

Environmental Sustainability and Organizational Responsibility. Environmental sustainability considerations demand attention as data centers supporting AI training and inference account for up to 3.7% of global carbon emissions, creating responsibilities that organizations increasingly recognize as important factors in technology adoption decisions [22]. SOTA NeSy systems leveraging symbolic reasoning to reduce model scale and computational requirements directly address environmental sustainability concerns while maintaining security effectiveness, enabling organizations to pursue advanced cybersecurity capabilities without compromising environmental responsibility.

The integration of symbolic knowledge enables more efficient learning and reasoning processes, reducing computational overhead while maintaining effectiveness in cybersecurity applications that demand both accuracy and interpretability. This efficiency advantage supports organizational adoption aligned with sustainability objectives while enabling effective threat detection and response capabilities essential for comprehensive security programs.

Organizations must balance performance requirements against computational efficiency while considering environmental impact and operational sustainability in deployment decisions. Recent NeSy systems provide frameworks for achieving this balance through intelligent integration of symbolic knowledge with neural learning capabilities, enabling effective cybersecurity while maintaining alignment with organizational sustainability objectives and environmental responsibilities that increasingly influence technology adoption decisions across diverse contexts.

V. ARCHITECTURAL PATTERNS AND TECHNICAL CHALLENGES

The evolution of NeSy cybersecurity systems has revealed distinct architectural patterns addressing different operational requirements and deployment constraints through systematic organization of neural and symbolic components while ensuring proper grounding, instructibility, and alignment with organizational objectives [26], [32], [33]. Understanding these patterns proves crucial for selecting appropriate approaches and identifying persistent challenges that limit widespread adoption in operational cybersecurity environments while sup-

porting responsible development aligned with defensive cybersecurity goals and societal expectations.

Contemporary architectural patterns demonstrate consistent trends toward multi-agent collaborative frameworks, continual learning systems addressing evolving threat landscapes, and causal NeSy architectures enabling sophisticated reasoning about attack causality and counterfactual scenarios. These patterns emerge from systematic analysis of breakthrough applications and implementation experiences, revealing fundamental principles that transcend individual application domains while highlighting critical challenges requiring coordinated research and development efforts.

A. Advanced Multi-Agent Collaborative Architectures

The consistent superiority of systems like HPTSA and VulnBot, detailed in Sec. III, is not an isolated success but rather points to a generalizable architectural pattern: the multi-agent collaborative framework. By abstracting from these concrete applications, we can formalize the core principles that drive their effectiveness. These architectures consistently employ hierarchical organization, where coordination agents manage task distribution and specialized agents handle domain-specific analysis such as reconnaissance, vulnerability assessment, exploitation, and verification. This structure enables systematic decomposition of complex security problems into manageable components while maintaining coherent strategic coordination.

Hierarchical Coordination and Specialization Patterns. The performance advantage of these systems stems from systematic specialization, which enables individual agents to develop deep expertise while collaborative mechanisms ensure comprehensive coverage of complex security scenarios. As seen in the successful frameworks, coordination agents manage high-level strategic objectives and task allocation, while specialized agents focus on technical domains like network analysis, malware assessment, or threat intelligence processing. This creates systems that effectively leverage both breadth and depth of expertise.

The mathematical framework for multi-agent performance can formalize these observed specialization benefits and collaborative synergies. The model for collaborative performance demonstrates how this distributed reasoning achieves superior capabilities:

$$P_{\text{multi}} = \sum_{i=1}^k \alpha_i \cdot P_{\text{individual}}(a_i) + \beta \cdot \sum_{i < j} \text{Synergy}(a_i, a_j)$$

For optimal task decomposition, the optimization objective must balance computational efficiency with coordination effectiveness:

$$\min \sum_i \text{Cost}(f_i) + \lambda \cdot \text{CommOverhead}$$

subject to completeness constraints $\bigcup_i f_i = \text{CompleteTask}$ and dependency requirements $\text{Deps}(f_i) \in \text{FeasibleSchedules}$. The theoretical performance bound demonstrates why collaborative reasoning consistently exceeds individual agent capabilities:

$$P_{\text{multi}} \geq \max(P_{\text{single}}) + \delta(\text{coordination_quality}, \text{complexity})$$

Cross-Validation and Error Correction Mechanisms. A key principle evident in advanced multi-agent architectures is their use of sophisticated cross-validation. This mechanism enables agents to identify and correct individual errors through collaborative analysis and proves particularly effective where individual detection systems may exhibit blind spots.

The effectiveness of this cross-validation depends on the diversity of agent capabilities, efficient communication protocols, robust consensus mechanisms, and feedback systems that support continuous improvement [199]. These mechanisms enable the system to achieve superior accuracy and explainability while supporting instructible adaptation to evolving threats.

The 20-30% improvements in detection coverage and lower false positive rates seen in recent implementations highlight the effectiveness of collaborative reasoning and cross-validation [200]. The coordination effectiveness parameter β is crucial for overall system performance, as higher values indicate superior agent communication and knowledge sharing.

B. Continual Learning and Adaptive Architectures

Continual Learning Architectures address critical challenge of adapting to evolving threat landscapes [201], [202] while avoiding catastrophic forgetting through systematic integration of symbolic knowledge updates with neural learning processes, enabling systems that maintain effectiveness against known threats while developing capabilities for emerging attack patterns [77], [125]. These architectures prove essential for operational deployment in dynamic cybersecurity environments where attack methods constantly evolve and defensive strategies must adapt accordingly while maintaining proper grounding in established security principles and supporting instructible guidance from security analysts.

Knowledge Preservation and Dynamic Adaptation Mechanisms. SOTA continual learning systems implement sophisticated mechanisms balancing knowledge preservation with adaptive learning capabilities while supporting instructible modification based on analyst expertise and organizational priorities. The NeSyC framework demonstrates advanced two-phase continual learning loop where symbolic rules are dynamically reformulated based on new experiences using contrastive generality improvement schemes, then applied adaptively to new threats while preserving accumulated security expertise relevant for recurring attack patterns [77], [203].

Memory-based monitoring schemes detect failures and trigger re-entry into reformulation phases, enabling autonomous adaptation to novel threats while maintaining grounding in established security principles and supporting instructible refinement based on operational feedback and expert guidance. These mechanisms address fundamental challenge where learning new network behaviors may cause systems to lose performance on previously learned attack patterns, creating dangerous security gaps when older threats re-emerge in modified forms requiring systematic balance between adaptation and knowledge retention.

The mathematical formulation for continual learning effectiveness captures trade-offs between backward transfer measuring performance degradation on previously learned tasks

and forward transfer measuring zero-shot performance on unseen threats based on accumulated knowledge. Effective continual learning systems maximize $\text{ContinualEffectiveness} = \alpha \cdot \text{BackwardTransfer} + \beta \cdot \text{ForwardTransfer}$ where weighting parameters reflect organizational priorities for maintaining effectiveness against known threats versus developing capabilities for emerging attack patterns.

Instructible Adaptation and Expert Guidance Integration.

Advanced continual learning architectures support instructible adaptation enabling security analysts to guide learning processes based on domain expertise, organizational priorities, and evolving threat intelligence while maintaining alignment with strategic objectives and operational requirements. These mechanisms enable rapid adaptation to emerging threats without requiring extensive retraining cycles that could compromise operational effectiveness during critical security incidents or major threat campaign evolution.

Instructible interfaces support analyst guidance for threat prioritization enabling focus on attack patterns most relevant to organizational risk profiles, rule modification supporting refinement of symbolic knowledge based on operational experience and evolving threat intelligence, and learning objective specification enabling alignment with organizational security strategies and resource allocation priorities. These capabilities ensure that continual learning processes remain aligned with organizational objectives while leveraging analyst expertise to improve system effectiveness and strategic alignment.

The integration of causal reasoning capabilities enables sophisticated understanding of why specific adaptations improve or degrade performance, supporting more effective instructible guidance and strategic decision-making about learning priorities and resource allocation. Analysts can understand not just what changes occurred during adaptation but why specific modifications enhanced or reduced effectiveness, enabling more informed guidance for future learning processes aligned with organizational objectives and threat landscape evolution.

C. Causal NeSy Architectures

Causal NeSy Architectures enable sophisticated reasoning about attack causality and counterfactual scenarios transcending traditional correlation-based analysis while supporting instructible modification of causal models based on expert knowledge and organizational threat intelligence [23], [24]. These architectures represent the most transformative advancement in cybersecurity analysis, moving the field from reactive pattern recognition to proactive prevention through genuine causal understanding of attack progression and defensive effectiveness aligned with organizational strategic objectives.

Causal Model Integration and Reasoning Frameworks.

SOTA causal architectures systematically integrate neural components for modeling posterior distributions with symbolic components for evaluating logical formulas and causal relationships governing attack progression while supporting instructible refinement based on analyst expertise and evolving threat intelligence. Applications include dynamic causal Bayesian optimization enabling real-time strategy adaptation, causal inference for understanding network vulnerabilities re-

vealing critical security dependencies, and counterfactual reasoning generating actionable insights about defensive strategies aligned with organizational capabilities and constraints.

The mathematical foundation defines cybersecurity causal models as $\mathcal{G} = (\mathcal{V}, \mathcal{E}, f)$ where $\mathcal{V} = \{X_1, X_2, \dots, X_n\}$ represents security events, \mathcal{E} denotes causal relationships, and f specifies functional relationships between variables. Strength functions $w : \mathcal{E} \rightarrow [0, 1]$ quantify causal influence while supporting instructible modification based on expert knowledge and operational experience. Causal chain discovery models attack progression through systematic analysis:

$$P(\text{attack_success} \mid \text{do}(\text{defense})) = \sum_i P(\text{attack_success} \mid \text{state}_i) \cdot P(\text{state}_i \mid \text{do}(\text{defense})) \quad (12)$$

Counterfactual analysis enables systematic evaluation of alternative defensive scenarios:

$$\text{CF}(y, x, x') = \mathbb{E}[Y \mid \text{do}(X = x')] - \mathbb{E}[Y \mid \text{do}(X = x)] \quad (13)$$

These formulations enable prediction of how defensive actions influence attack likelihood at each temporal stage while supporting the G-I-A framework through enhanced grounding in causal understanding, instructible modification of causal models, and alignment with organizational defensive strategies.

Counterfactual Analysis and Strategic Planning Support. Counterfactual analysis for threat prevention enables systematic evaluation of alternative defensive scenarios [114], [204], [205] through formalization $\text{CF}(y, x, x') = \mathbb{E}[Y \mid \text{do}(X = x')] - \mathbb{E}[Y \mid \text{do}(X = x)]$, quantifying how different defensive actions would alter attack outcomes while supporting instructible evaluation of strategic alternatives aligned with organizational resource constraints and risk tolerance. This capability enables security analysts to understand not just what happened during incidents, but why specific attack steps succeeded and how different defensive configurations might have prevented compromise through systematic causal analysis.

Attack progression analysis requires sophisticated temporal modeling capturing both statistical patterns in observable features and causal relationships governing state transitions in cybersecurity contexts. The integrated probability $P(\text{stage}_{t+1} \mid \text{stage}_t, \text{defense_state}) = \int \Phi_{\text{neural}}(\text{features}_t) \times \Psi_{\text{causal}}(\text{causal_graph}_t) dt$ combines neural pattern recognition with symbolic causal reasoning to model defensive action effectiveness while supporting instructible modification of causal models based on organizational threat intelligence and strategic priorities.

These causal explanations prove valuable for strategic planning where understanding attack causality enables development of counterfactual scenarios aligned with organizational objectives, incident response where causal chains guide effective remediation strategies, and defensive strategy development where counterfactual analysis informs resource allocation and capability development decisions. Systems generate actionable insights for improving defensive capabilities through specific technical and procedural modifications aligned with organizational security goals while supporting instructible refinement based on operational outcomes and strategic feedback.

D. Critical Technical Challenges and Limitations

Despite breakthrough capabilities demonstrated by SOTA NeSy systems, critical technical challenges persist across architectural patterns, limiting widespread adoption in operational cybersecurity environments while constraining achievement of proper grounding, instructibility, and alignment objectives essential for responsible deployment.

Knowledge Acquisition and Maintenance Bottlenecks. Knowledge acquisition and maintenance are significant bottlenecks limiting practical deployment scalability across organizational contexts. Manual knowledge engineering requirements create dependencies on specialized cybersecurity expertise for rule formulation, ontology development, and knowledge base construction that may not be available in all organizational contexts, particularly smaller organizations or those with limited cybersecurity personnel and technical capabilities.

Knowledge consistency maintenance across diverse sources presents ongoing challenges as threat landscapes evolve, requiring systematic procedures for validating knowledge updates, resolving conflicts between different information sources, and ensuring logical consistency across complex knowledge representations. Dynamic knowledge updating while maintaining logical consistency requires sophisticated frameworks supporting incremental modification without compromising system integrity or reasoning capabilities essential for reliable operation in high-stakes security environments.

These knowledge engineering requirements create dependencies on specialized expertise that may constrain organizational ability to achieve proper grounding and support instructible adaptation. Organizations must carefully balance automation benefits against knowledge engineering overhead while developing internal capabilities or external partnerships supporting effective knowledge management.

The scalability challenges are particularly acute when organizations customize systems to specific threat environments, regulatory requirements, or operational procedures, requiring extensive domain expertise and ongoing maintenance effort that may exceed available resources. These limitations highlight need for more automated knowledge acquisition and maintenance capabilities supporting broader organizational adoption while maintaining effectiveness and reliability.

Symbolic Reasoning Scalability and Brittleness. Symbolic reasoning scalability and brittleness challenges emerge when SOTA systems confront noisy, real-world data characteristic of operational cybersecurity environments while requiring maintenance of logical consistency and reasoning reliability. Traditional symbolic approaches prove brittle when confronted with incomplete information, contradictory evidence, or evolving attack patterns that exceed original rule formulations, creating potential failure modes that could compromise security effectiveness during critical incidents [104].

The integration of symbolic and neural components introduces additional complexity requiring careful balance between logical consistency and adaptive flexibility [104], [206]. Systems must handle uncertainty and ambiguity in cybersecurity data while maintaining reasoning reliability and explainability for earning analyst trust and operational acceptance.

Scalability challenges manifest when symbolic reasoning components encounter large-scale data processing requirements characteristic of enterprise cybersecurity environments, potentially creating performance bottlenecks that could compromise real-time threat detection and response capabilities. The brittleness concerns become particularly important when systems encounter novel attack patterns or environmental conditions that exceed original design assumptions, potentially leading to degraded performance or reasoning failures. These challenges highlight need for more robust symbolic reasoning frameworks supporting graceful degradation and adaptive modification while maintaining core reasoning capabilities.

Integration Complexity and Coordination Overhead. Integration complexity and coordination overhead present significant challenges for multi-component architectures requiring sophisticated orchestration of neural and symbolic components. Sophisticated feedback loops demonstrated in advanced systems introduce latency and complexity that may not meet real-time operational requirements essential for effective threat response to security incidents or attack campaigns.

The coordination mechanisms required for effective multi-agent operation create additional complexity requiring careful design and optimization to avoid performance bottlenecks while maintaining collaborative benefits essential for superior detection and response capabilities. Communication protocols must balance information sharing effectiveness with computational overhead while ensuring that coordination activities do not compromise individual agent performance or system-wide responsiveness to emerging threats.

Organizations must carefully assess whether performance improvements justify increased architectural complexity and potential failure points of advanced multi-component systems, particularly in contexts where system reliability and operational continuity take precedence over marginal accuracy improvements. Integration overhead is challenging when organizations attempt to incorporate NeSy systems into existing security infrastructure comprising diverse tools and platforms with varying interfaces, data formats, and operational procedures. These integration requirements may exceed organizational technical capabilities or require extensive customization effort that constrains practical adoption while potentially compromising system effectiveness or reliability.

VI. KEY INSIGHTS AND G-I-A FRAMEWORK VALIDATION

Our systematic analysis across the 127 publications spanning 2019–July 2025 reveals fundamental principles that transcend individual application domains, supporting the G-I-A framework’s explanatory potential while highlighting critical requirements for achieving proper grounding, instructibility, and alignment with organizational objectives and societal expectations. Figure 12 illustrates the research maturity landscape, revealing distinct performance patterns and implementation gaps that reinforce our systematic analysis.

Contemporary analysis indicates that successful NeSy implementations consistently achieve superior performance through systematic integration of complementary reasoning paradigms while addressing fundamental limitations that constrain traditional approaches. The consistent superiority of

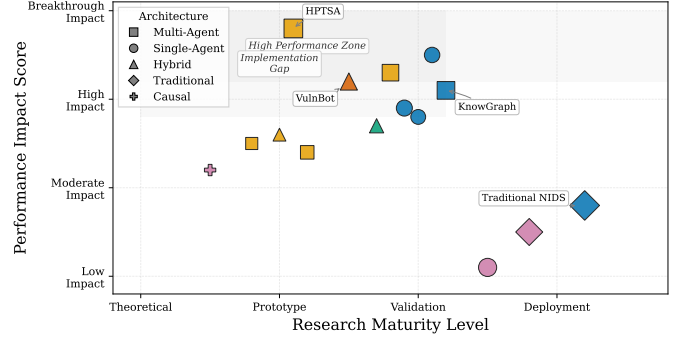


Fig. 12. NeSy cybersecurity research landscape and maturity assessment for 127 systems, plotted by research maturity level and performance impact score. System architectures are color-coded: orange squares (multi-agent), blue circles (single-agent), green triangles (hybrid), and dark blue diamonds (traditional approaches). Multi-agent architectures show a +12% average impact advantage over single-agent designs, while an implementation gap separates high-performing prototypes from deployment-ready systems. The distribution aligns with G-I-A framework patterns, with concentrated activity in the validation phase and traditional approaches clustering in lower-impact, deployment-ready zones, suggesting the performance–maturity trade-off shaping NeSy development.

multi-agent coordination architectures supports our G-I-A framework’s emphasis on collaborative reasoning, suggesting advantages rooted in cognitive science principles of distributed problem-solving [7], [8], [53].

A. G-I-A Framework Validation Through Performance Patterns

Multi-Agent Coordination Superiority Validates G-I-A Principles. Our systematic review reveals consistent superior performance of multi-agent NeSy architectures across diverse cybersecurity applications, indicating fundamental advantages rooted in distributed problem-solving principles rather than domain-specific optimizations [7], [8], [53]. The mathematical foundation $P_{\text{multi}} = \sum_{i=1}^k \alpha_i \cdot P_{\text{individual}}(a_i) + \beta \cdot \sum_{i < j} \text{Synergy}(a_i, a_j)$ demonstrates how collaborative benefits consistently exceed individual capabilities while supporting instructible enhancement through coordination parameter optimization.

Figure 13 illustrates this superiority across four key application domains, with performance improvements ranging from 14% in intrusion detection to 234% in penetration testing scenarios.

Causal Reasoning as Significant Theoretical Advancement. Causal reasoning integration enables qualitative leaps from correlation-based pattern recognition to genuine understanding of attack causality, transforming cybersecurity from reactive threat detection to proactive prevention [23], [24]. Causal NeSy frameworks enable generation of explanations such as “malicious email → downloaded attachment → process execution → network connection → data exfiltration” providing actionable insights for strategic planning and resource allocation aligned with organizational objectives.

Enhanced Grounding Validates Explainability Requirements. The explainability advantages inherent in symbolic

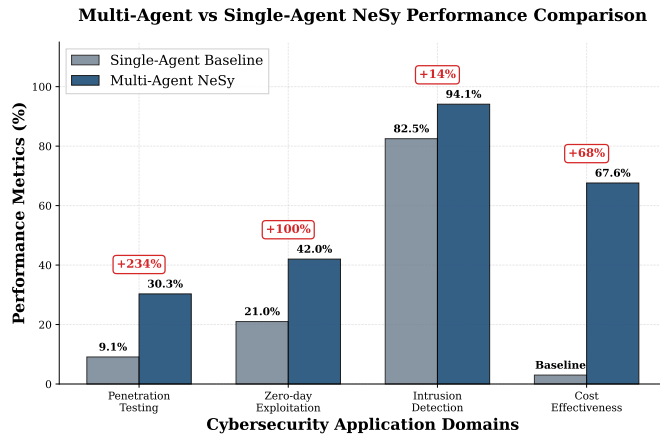


Fig. 13. Multi-agent NeSy architectures demonstrate consistent performance superiority across diverse cybersecurity applications. Performance improvements range from 14% in intrusion detection to 234% in penetration testing, with substantial cost reductions (67.6%) validating collaborative reasoning advantages over single-agent approaches.

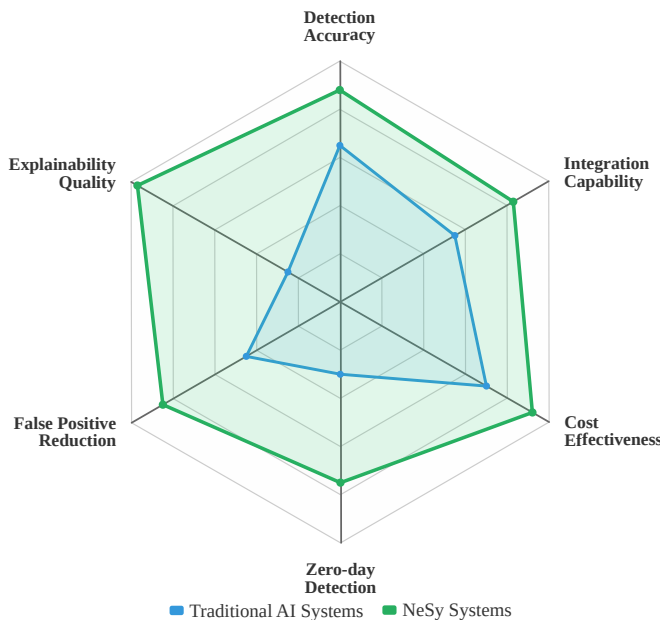


Fig. 14. Multi-dimensional performance analysis comparing Traditional AI Systems versus Advanced NeSy Systems across six key cybersecurity operational metrics. NeSy approaches demonstrate substantial advantages across all dimensions, with particularly strong improvements in explainability quality, false positive reduction, and cost effectiveness while supporting instructible adaptation and organizational alignment.

reasoning components directly translate to operational trust and adoption success in high-stakes security environments [55], [100]. Systems like Kalutharage et al.’s IoT IDS demonstrate how explicit mapping to established frameworks such as MITRE ATT&CK creates actionable explanations that security professionals can immediately understand while supporting instructible refinement based on operational outcomes.

B. Knowledge Integration and G-I-A Operational Validation

Systematic Knowledge Operationalization and Grounding Requirements. Advanced NeSy systems excel at systematically transforming unstructured threat intelligence into actionable structured knowledge, bridging critical gaps between human expertise and automated security operations [64], [73]. Domain knowledge integration consistently produces substantial performance improvements across diverse applications, supporting our G-I-A framework’s emphasis on proper grounding mechanisms [52], [54]. The systematic performance patterns indicate that NeSy advantages emerge most prominently in scenarios characterized by clear symbolic structures, adversarial environments where logical constraints provide additional validation layers, and contexts requiring explanation transparency. These multi-dimensional advantages are summarized in Figure 14.

C. Dual-Use Implications and Responsible G-I-A Alignment

Autonomous Capabilities and Alignment Requirements.

The rapid emergence of autonomous offensive capabilities marks a pivotal change in the cybersecurity threat landscape. As discussed in Section III-D, SOTA NeSy agents now deliver advanced exploitation performance at a fraction of historical human-expert costs, lowering economic barriers for sophisticated attacks [6], [7]. This democratization of capability highlights the urgency for proactive alignment measures, coordinated defensive responses, and governance frameworks that anticipate misuse scenarios [175], [207].

Responsible Development Within the G-I-A Alignment Framework.

The dual-use nature of NeSy systems demands research practices that explicitly integrate alignment objectives into system design and evaluation [6], [21], [161]. Within the G-I-A framework, this translates into:

- **Grounding (G):** Embed symbolic security policies and verified vulnerability taxonomies into system knowledge bases to ensure that reasoning paths remain tied to legitimate defensive objectives and exclude prohibited exploitation contexts.
- **Instructibility (I):** Implement instruction filters and policy-aware prompt parsing to block malicious or non-compliant task requests before they can be operationalized by the agent.
- **Alignment (A):** Incorporate technical safeguards—such as built-in defensive biasing, real-time misuse detection, user verification, kill switches, and runtime isolation—to enforce compliance during execution [175].

Community-Governed Standards and Oversight.

To ensure that G-I-A principles translate into consistent practice, cybersecurity research communities should adopt shared standards for dual-use AI security research, supported by transparent evaluation protocols and red-team audits prior to release. Community-driven ethical guidelines and peer-reviewed disclosure frameworks can help balance innovation with risk mitigation, maintaining alignment with both defensive cybersecurity objectives and societal expectations [6]. Such standards should explicitly require that systems claiming G-I-A com-

pliance demonstrate verifiable defensive biasing, documented safe-mode defaults, and reproducible alignment testing.

VII. UNIQUENESS AND COMPARATIVE ANALYSIS

This comprehensive analysis provides several unique contributions distinguishing it from existing SOTA review articles in AI and cybersecurity. Our work offers comprehensive scope, novel analytical perspectives, and systematic focus on implementation challenges while emphasizing proper grounding, instructibility, and alignment considerations essential for responsible development [19], [22], [26].

Unlike existing surveys focusing primarily on technical capabilities or defensive applications in isolation, this work provides the first systematic examination of both cutting-edge defensive innovations and concerning dual-use implications. This approach addresses critical gaps in existing literature while emphasizing responsible development aligned with cybersecurity objectives and societal expectations.

A. Systematic Comparative Analysis with Existing Surveys

This review positions itself within the existing literature by systematically comparing its contributions to SOTA surveys. We highlight novel perspectives that address critical gaps through the G-I-A framework. While prior surveys provide valuable insights into specific NeSy applications, they lack the comprehensive scope, dual-use analysis, and implementation focus that define this work.

Vision & Strategic Framework Innovation. Unlike surveys such as Colelough & Regli (2025) and Bhuyan et al. (2024), which focus on isolated technical capabilities (e.g., knowledge-enhanced NIDS), this work introduces the G-I-A framework to systematically evaluate NeSy systems for robust conceptual grounding, responsive instructibility, and alignment with cybersecurity objectives [17], [18]. This unified theoretical approach addresses a gap in prior reviews that do not formalize these critical dimensions. Furthermore, our multi-year phased strategic roadmap provides coordinated advancement guidance absent in works like Salem et al. (2024) and Bilot et al. (2024) [208], [209].

Advanced Technical & Paradigmatic Synthesis. This survey synthesizes NeSy as a distinct AI paradigm for cybersecurity, contrasting with more general approaches in Wang et al. (2025) and Arp et al. (2022) that address broader AI applications without cybersecurity-specific integration [210], [211]. Notably, prior surveys overlook causal reasoning as a foundational component. In contrast, this work highlights its transformative potential for understanding attack causality and enabling proactive defense, a perspective absent in Colelough & Regli (2025) and most existing reviews [23], [24].

Cybersecurity-Centric Application Analysis. Existing surveys, including Ferrag et al. (2024) and Capuano et al. (2022), primarily emphasize defensive applications and only partially address offensive implications [2], [212]. This review provides systematic analysis of both defensive and offensive NeSy capabilities, examining autonomous systems achieving high zero-day exploitation success rates and their ethical implications. This comprehensive coverage addresses a critical literature gap

[6], [7]. While Eckhoff et al. (2025) and Piplai et al. (2023) touch on practical applications, they lack the comprehensive SOC lifecycle integration demonstrated here [64], [160].

Ethical Dimensions & Responsible AI Development. This work establishes comprehensive dual-use analysis absent in existing surveys, while proposing actionable responsible development principles that extend beyond the general ethical considerations mentioned in Capuano et al. (2022) [212]. The systematic examination of autonomous offensive capabilities and their implications represents unprecedented coverage in the NeSy cybersecurity literature.

Evaluation, Implementation & Human Factors. Unlike Ferrag et al. (2024), which minimally address deployment challenges, this review systematically examines computational complexity, standardization gaps, and human-AI collaboration within the G-I-A framework [30], [159]. While Yan et al. (2023) discuss graph-based approaches and Renkhoff et al. (2024) touch on evaluation aspects, they lack the comprehensive standardization analysis and community coordination framework presented here [213]. Our integration of human-AI collaboration as a foundational pillar aligns with practical considerations in Eckhoff et al. (2025) but provides deeper theoretical grounding [160].

Table VII summarizes these distinctions, demonstrating how this work's G-I-A focus, comprehensive dual-use analysis, and implementation insights provide unique contributions across all major evaluation dimensions while addressing critical gaps in existing literature.

B. Novel Analytical Frameworks and Comprehensive Scope

This survey introduces novel analytical frameworks that advance NeSy cybersecurity research, emphasizing the G-I-A framework to ensure robust grounding, instructible adaptation, and alignment with defensive objectives. These frameworks address gaps in existing literature while providing practical guidance for operational deployment.

G-I-A Framework Innovation. The G-I-A framework formalizes the integration of grounding (conceptual understanding), instructibility (analyst-guided adaptation), and alignment (consistency with cybersecurity goals) through a mathematical model (see Sec. II). Unlike prior surveys, which lack such a unified approach, this framework provides a lens for evaluating NeSy systems across applications, from intrusion detection to autonomous operations [17], [18]. Its novelty lies in its systematic operationalization, though empirical validation remains a future research priority.

Dual-Use Analysis with Responsible Development. This work provides comprehensive analysis of NeSy's dual-use implications, examining autonomous systems achieving notable success rates on zero-day vulnerabilities and cost reductions from approximately \$100 down to \$24.40 per exploit [7]. By integrating responsible development principles (e.g., defensive co-development, technical safeguards), it ensures alignment with cybersecurity objectives, addressing a gap in surveys that overlook offensive implications [6].

Multi-Agent Architecture Superiority. The survey's analysis of multi-agent NeSy systems documents dramatic performance

TABLE VII
A STRUCTURED COMPARATIVE ANALYSIS HIGHLIGHTING THE UNIQUE CONTRIBUTIONS OF OUR WORK

Contribution Area	Our Work	Colelough & Reghi [19]	Eckhoff et al. [160]	Wang et al. [210]	Bhuyan et al. [18]	Salem et al. [208]	Bilot et al. [209]	Piplai et al. [64]	Yan et al. [213]	Arp et al. [211]	Capuano et al. [212]
Vision & Strategic Framework											
Introduces Unifying Theoretical Framework (G-I-A)	✓	✗	-	-	✗	✗	✗	-	✗	-	-
Provides Multi-Year, Phased Strategic Roadmap	✓	-	-	-	-	-	-	✗	-	-	-
Advanced Technical & Paradigmatic Synthesis											
Synthesizes NeSy as a Distinct AI Paradigm for Cyber	✓	✓	-	-	✓	✗	✗	✗	-	✗	✗
Integrates Causal Reasoning as a Transformative Capability	✓	✗	-	✗	-	✗	✗	✗	✗	✗	✗
Cybersecurity-Centric Application Analysis											
Systematic Analysis of Autonomous Cyber Operations	✓	-	-	-	✗	✗	✗	-	✗	✗	✗
Connects NeSy Theory to Full SOC Lifecycle	✓	-	✓	-	✗	-	✗	-	-	✗	✗
Ethical Dimensions & Responsible AI											
Comprehensive Dual-Use & Offensive AI Analysis	✓	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
Proposes Actionable Responsible Development Principles	✓	-	-	-	-	✗	✗	-	-	✗	✓
Evaluation, Implementation & Human Factors											
Defines NeSy-Specific Gaps in Standardization	✓	-	✗	✓	-	-	-	✓	-	-	✗
Integrates Human-AI Collaboration as a Foundational Pillar	✓	✓	✓	✓	-	✗	✗	✓	✓	-	-
Target Audience	Research & Policy	General AI	Applied NeSy	AI for Edge	General AI	General AI for Cyber	Graph AI for Cyber	Applied NeSy	Graph AI for Cyber	ML Best Practices	XAI for Cyber

✓ = Comprehensive Coverage - = Partial or Implicit Coverage ✗ = Limited or No Focus

gains, with multi-agent frameworks outperforming single-agent baselines by over 200% in some penetration testing scenarios [7], [8]. This perspective, absent in prior reviews, provides practical guidance for implementing scalable, instructible systems.

Causal Reasoning as Transformative Capability. The integration of causal reasoning enables NeSy systems to move beyond correlation-based analysis to genuine attack causality understanding, supporting proactive defense through counterfactual scenarios. This novel focus, not covered in existing surveys, aligns with the manuscript’s future research priorities for causal model development [23], [24].

Implementation and Standardization Guidance. Through the G-I-A lens, this survey addresses implementation challenges (e.g., computational complexity, knowledge engineering bottlenecks) and proposes community-driven standardization to bridge evaluation gaps, aligning with Section VII’s roadmap for coordinated advancement [30], [159]. This focus ensures practical applicability across diverse organizational contexts.

C. Implementation-Focused Analysis and Practical Guidance

Comprehensive Implementation Challenge Analysis. Unlike existing surveys focusing primarily on technical capabilities, this work provides systematic analysis of implementation challenges including computational complexity, standardization gaps, human-centric deployment factors, and organizational success requirements that determine practical adoption success while supporting instructible adaptation and maintaining alignment with strategic objectives. Our implementation analysis addresses critical barriers between research prototypes and operational deployment while identifying systematic solutions enabling successful adoption across diverse organizational contexts and computational environments.

The implementation framework includes systematic evaluation of computational requirements and resource optimization strategies, comprehensive analysis of integration challenges with existing security infrastructure and operational procedures, detailed assessment of human-AI collaboration patterns and trust development requirements, and practical guidance for addressing organizational adoption barriers while maintaining alignment with cybersecurity objectives and strategic priorities. These insights provide actionable guidance for practitioners and organizations seeking to leverage advanced NeSy capabilities while ensuring proper grounding in operational contexts and instructible adaptation to diverse requirements.

Standardization Gap Analysis and Community Coordination. This survey provides first comprehensive analysis of evaluation and standardization gaps specifically constraining NeSy cybersecurity development while establishing systematic framework for community-driven coordination efforts essential for responsible field advancement. Our standardization analysis demonstrates that absence of NeSy-specific benchmarks represents field’s greatest limitation for systematic progress while providing detailed roadmap for addressing evaluation gaps through coordinated community initiatives aligned with responsible development principles.

The standardization framework includes systematic identification of critical evaluation gaps and their impact on field advancement, comprehensive analysis of existing benchmark limitations and NeSy-specific requirements, detailed proposal for community-driven standardization initiatives supporting coordinated progress, and practical guidance for establishing evaluation frameworks enabling systematic comparison while ensuring proper grounding in cybersecurity concepts and alignment with defensive objectives. These contributions establish foundations for coordinated field advancement while maintaining focus on responsible development aligned with cybersecurity needs and societal expectations.

VIII. FUTURE RESEARCH OPPORTUNITIES AND STRATEGIC DIRECTIONS

Directly addressing RQ6 while synthesizing findings across all research questions, this section discusses prioritized research opportunities based on analysis of identified gaps and emerging possibilities that promise advancement of NeSy cybersecurity capabilities over the next five years while ensuring proper grounding, instructibility, and alignment with defensive cybersecurity objectives and societal expectations [22], [32]. These research directions address both fundamental research challenges and practical deployment requirements essential for widespread adoption in operational cybersecurity environments while maintaining focus on responsible development aligned with organizational needs and ethical considerations. Importantly, these opportunities build directly on the limitations identified in our review, particularly the lack of standardized evaluation protocols, insufficient causal reasoning integration, and the absence of systematic alignment mechanisms.

Our analysis reveals that future advancement depends critically on coordinated efforts addressing three foundational requirements: enhanced grounding mechanisms enabling robust conceptual understanding of cybersecurity concepts, advanced instructibility frameworks supporting effective human-AI collaboration based on expert guidance, and alignment approaches ensuring technological development serves defensive cybersecurity objectives while maintaining consistency with organizational priorities and societal expectations. These requirements transcend individual technical challenges to encompass advancement of how cybersecurity AI systems are developed, evaluated, and deployed in operational environments. Figure 15 illustrates the strategic roadmap for these efforts, detailing the critical research priorities and a phased development plan. Achieving these objectives will require close collaboration among academic researchers, standards organizations, and industry stakeholders to ensure both scientific rigor and operational relevance.

A. Critical Research Priorities for Field Advancement

Standardized Evaluation Frameworks and Community Coordination. The development of standardized benchmarks geared towards NeSy cybersecurity evaluation [187], [188], [194] is the most critical research problem constraining field advancement while limiting ability to achieve proper grounding, instructibility, and alignment across diverse research and deployment contexts [30], [178], [179]. These frameworks must integrate neural pattern recognition assessment with symbolic reasoning evaluation in realistic cybersecurity scenarios capturing operational complexity and supporting evaluation of grounding mechanisms, instructibility capabilities, and alignment with defensive objectives [187], [188].

NeSy benchmarks must include multi-stage attack scenario datasets with structured logical rule representation capturing threat campaigns requiring coordinated reasoning across multiple attack phases. Knowledge graph integration supporting both learning and reasoning evaluation must enable assessment

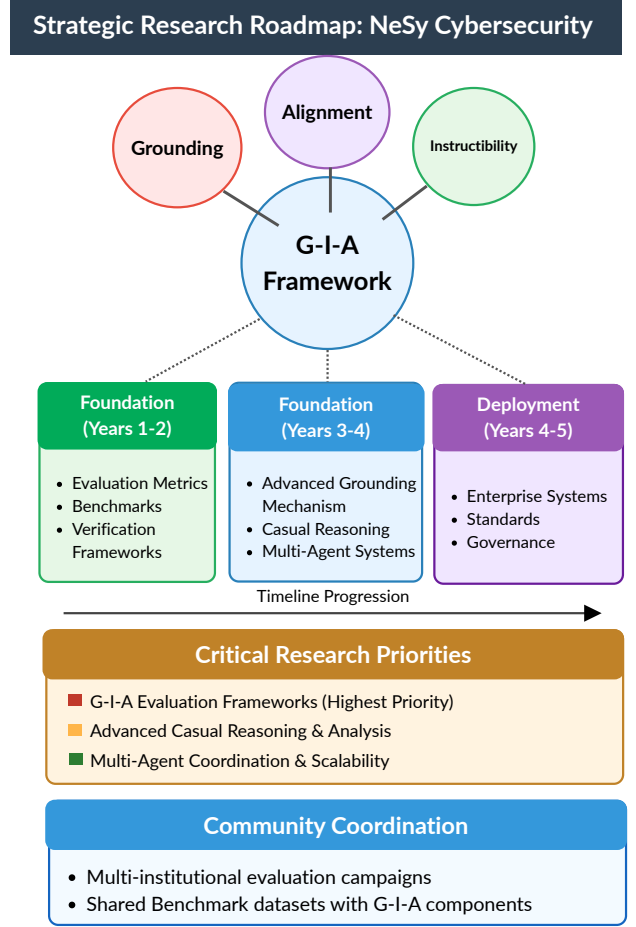


Fig. 15. Strategic roadmap for NeSy cybersecurity advancement through G-I-A framework operationalization across three progressive development phases with critical research priorities.

of knowledge utilization effectiveness while supporting evaluation of grounding mechanisms ensuring proper understanding of cybersecurity concepts across organizational contexts.

Logical rule components focusing on formal reasoning rather than pattern recognition alone must enable evaluation of symbolic reasoning capabilities while supporting assessment of instructible modification based on expert guidance reflecting diverse operational environments. Causal reasoning evaluation frameworks for assessing counterfactual and causal chain generation capabilities must support evaluation of advanced reasoning capabilities enabling understanding of attack causality and defensive strategy effectiveness aligned with organizational objectives.

Community-driven standardization initiatives should bring together researchers from NeSy AI, cybersecurity, and evaluation methodology communities to establish frameworks supporting coordinated field advancement [29], [103]. Proposed efforts include development of common evaluation protocols enabling comparison across different approaches, creation of shared benchmark datasets with integrated neural and symbolic evaluation components, establishment of explanation quality assessment frameworks specifically for security contexts accommodating analyst workflows, and coordination of

multi-institutional evaluation campaigns promoting collaborative advancement with responsible development principles.

Practical Relevance and Standardization Pathways. Future research must also ensure that the G-I-A framework contributes to practical deployment, governance, and long-term adoption. One promising direction is the integration of G-I-A scoring with emerging AI security and trustworthiness standards, including the NIST AI Risk Management Framework [214], the EU AI Act [215], and ISO/IEC initiatives on trustworthy AI [216]. Such integration would enable consistent evaluation of NeSy cybersecurity systems in line with both regulatory requirements and organizational needs.

Furthermore, alignment with community-driven benchmarking initiatives such as CTIBench [103], CVE-Bench [29], and SECURE [179] provides a concrete pathway for operational adoption. Embedding G-I-A evaluation protocols into these benchmarks would facilitate reproducible comparison across approaches while directly supporting real-world deployment decisions in cybersecurity environments.

Advanced Causal Reasoning and Counterfactual Analysis. The integration of causal reasoning capabilities represents significant research direction enabling understanding of attack causality and counterfactual threat scenarios beyond traditional correlation-based analysis while supporting instructible modification of causal models and maintaining alignment with organizational strategic objectives [23], [24], [90], [217]. Future research should focus on developing robust causal discovery algorithms for cybersecurity domains enabling automated identification of causal relationships from observational data while supporting instructible refinement based on expert knowledge and organizational threat intelligence.

Scalable causal inference frameworks for real-time security analysis must enable dynamic causal reasoning during active security incidents while maintaining computational efficiency essential for operational effectiveness. Formal verification methods for causal explanations must ensure logical consistency and accuracy of causal models while supporting validation procedures enabling analyst confidence and organizational adoption aligned with strategic objectives.

Integration of causal reasoning with existing symbolic knowledge bases must enable understanding of attack causality within broader cybersecurity knowledge frameworks while supporting instructible enhancement of causal models based on accumulated domain expertise. Research priorities include development of automated causal model construction from cybersecurity data sources supporting scalable deployment, creation of interactive interfaces enabling analyst instruction and validation of causal models supporting effective human-AI collaboration, establishment of formal verification frameworks ensuring causal model accuracy essential for high-stakes security decisions, and integration of causal reasoning with existing security frameworks supporting strategic planning aligned with organizational objectives.

Enhanced Multi-Agent Coordination and Scalability. The demonstrated superiority of multi-agent NeSy architectures across cybersecurity applications necessitates focused research on coordination mechanisms, optimization strategies, and scalability frameworks supporting effective collaboration while

maintaining instructibility and alignment with organizational objectives [8], [53], [79]. Research priorities include developing efficient communication protocols for large-scale agent coordination enabling effective information sharing without overwhelming computational overhead while supporting instructible modification of coordination strategies based on operational requirements.

Adaptive task decomposition algorithms for complex security problems must enable dynamic allocation of responsibilities across specialized agents while maintaining coverage of security scenarios. Robust consensus mechanisms for collaborative decision-making must balance individual agent expertise with collective validation while supporting instructible refinement of decision criteria based on operational outcomes and expert guidance aligned with organizational strategic objectives.

Future research should emphasize development of self-organizing multi-agent systems capable of autonomous coordination adaptation while maintaining human oversight and instructible modification capabilities, creation of formal verification frameworks for multi-agent coordination ensuring reliable collaborative behavior in high-stakes security environments, and integration of multi-agent coordination with existing security operations workflows supporting enhanced effectiveness while maintaining alignment with operational procedures and strategic priorities.

B. G-I-A Framework Development and Evaluation Research

Advanced Grounding Mechanisms for Robust Conceptual Understanding. Future research must prioritize development of sophisticated grounding mechanisms enabling NeSy systems to establish meaningful connections between their outputs and abstract cybersecurity concepts while maintaining robustness against novel threats and adversarial manipulation [218]. Enhanced grounding requires systematic integration of conceptual understanding with operational effectiveness, enabling systems to demonstrate clear connections between reasoning processes and fundamental security principles while supporting instructible refinement based on expert knowledge and organizational context.

Research priorities include development of automated concept learning frameworks enabling systems to acquire deep understanding of cybersecurity relationships from operational data while supporting validation and refinement based on expert guidance, creation of robustness mechanisms ensuring grounding remains effective against adversarial attacks and novel threat patterns while maintaining alignment with established security principles, and establishment of evaluation frameworks enabling systematic assessment of grounding quality and its impact on operational effectiveness across diverse cybersecurity contexts.

Comprehensive Instructibility and Alignment Frameworks. Instructibility represents critical capability enabling security analysts to provide explicit feedback that appropriately modifies system behavior while maintaining alignment with organizational objectives and strategic priorities. Future research must develop sophisticated instructibility frameworks

supporting effective human guidance across diverse expertise levels and organizational contexts while ensuring that system adaptations remain consistent with cybersecurity principles and operational requirements.

Research priorities include development of intuitive interfaces enabling effective analyst instruction without requiring deep technical understanding of underlying system architecture, creation of validation mechanisms ensuring instructible modifications improve rather than degrade system performance while maintaining alignment with organizational objectives, and establishment of learning frameworks enabling systems to generalize from specific instructions to broader understanding of analyst preferences and organizational priorities aligned with strategic cybersecurity goals.

Alignment is a fundamental requirement ensuring NeSy systems operate consistently with cybersecurity objectives, organizational priorities, and societal expectations while avoiding behaviors that could compromise security effectiveness or enable malicious applications. Comprehensive alignment must encompass technical alignment ensuring optimization targets reflect true cybersecurity goals, ethical alignment ensuring capabilities serve defensive purposes and societal cybersecurity needs, and operational alignment ensuring system behavior remains consistent with organizational policies and strategic objectives across diverse deployment contexts.

C. Long-Term Vision and Strategic Implementation Roadmap

Phased Development Strategy for Coordinated Advancement. Years 1-2 should focus on Foundation Building through establishing standardized evaluation frameworks and benchmark datasets supporting systematic comparison and reproducible research, developing automated knowledge extraction techniques reducing manual engineering bottlenecks while maintaining quality and reliability, creating open source reference implementations enabling broader community access and collaborative development, and conducting comprehensive human factors studies supporting effective deployment and organizational adoption aligned with strategic objectives.

Years 3-4 should emphasize Capability Development by deploying advanced causal reasoning systems enabling sophisticated attack causality understanding and strategic planning support, developing proactive threat anticipation capabilities enabling prevention rather than reactive response while maintaining alignment with organizational risk tolerance, establishing cross-organizational knowledge sharing frameworks supporting coordinated defense and intelligence collaboration, and creating comprehensive cyber-physical security applications addressing critical infrastructure protection aligned with national security objectives and societal expectations.

Years 4-5 should achieve Operational Deployment through widespread enterprise deployment across diverse organizational contexts and computational environments, mature standardization frameworks enabling systematic evaluation and comparison across different approaches and implementations, quantum-resistant approaches ensuring long-term security effectiveness against emerging computational threats, and autonomous cyber defense capabilities enabling proactive pro-

tection while maintaining human oversight and strategic alignment with organizational objectives and societal expectations. **Impact Objectives and Success Metrics.** The successful development and deployment of advanced NeSy cybersecurity systems promises significant impacts across multiple dimensions while ensuring proper grounding, instructibility, and alignment with cybersecurity objectives and societal expectations. Defensive Capability Enhancement will enable more effective, adaptive, and explainable cybersecurity defenses that can anticipate and counter sophisticated adversary capabilities while providing essential transparency and accountability supporting analyst trust and organizational adoption.

Human-AI Collaboration Advancement will substantially improve how security analysts work, providing intelligent assistance to enhance human experts while ensuring humans are in control of critical security decisions and maintaining alignment with organizational strategic objectives and operational procedures. Threat Landscape Stabilization may help stabilize cybersecurity threat landscape by making attacks more difficult and expensive while providing defenders with advantages through advanced reasoning capabilities aligned with defensive objectives and societal cybersecurity needs.

Success metrics include quantitative measures of defensive capability improvement across diverse organizational contexts, systematic evaluation of human-AI collaboration effectiveness and analyst satisfaction with enhanced capabilities, comprehensive assessment of threat landscape impact and defensive advantage maintenance, and qualitative evaluation of alignment with cybersecurity objectives and societal expectations ensuring responsible development and deployment aligned with strategic priorities and ethical considerations.

Sustainability and Global Coordination Requirements. Long-term success requires sustainable development approaches addressing resource efficiency, environmental responsibility [217], and global coordination essential for effective cybersecurity enhancement while maintaining alignment with diverse organizational capabilities and strategic priorities. Sustainable development must balance advanced capability requirements with computational efficiency enabling broader organizational adoption while addressing environmental concerns and resource constraints that may limit deployment across diverse contexts.

Global coordination efforts must address international cooperation requirements for effective cybersecurity enhancement while managing dual-use concerns and ensuring that advanced capabilities serve defensive objectives aligned with international cybersecurity cooperation and societal expectations. These coordination requirements encompass technical standardization supporting interoperability across diverse organizational and national contexts, policy frameworks addressing dual-use concerns and responsible development principles, and collaborative research initiatives supporting coordinated advancement while maintaining competitive advantages for defensive cybersecurity applications.

IX. CONCLUSION

This systematic review establishes that NeSy AI represents a significant advancement in cybersecurity, addressing funda-

mental limitations of traditional approaches through principled integration of neural adaptability with symbolic reasoning. Our analysis of 127 publications spanning 2019–July 2025 demonstrates that advanced NeSy systems achieve substantial capabilities across defensive applications. Simultaneously, our findings reveal concerning dual-use implications that demand responsible development frameworks.

The evidence shows consistent superiority of multi-agent architectures and the emergence of causal reasoning capabilities. This signals important progress from correlation-based threat detection toward genuine understanding of attack causality, enabling proactive defense strategies that go beyond traditional reactive security approaches. However, the proliferation of autonomous offensive capabilities such as zero-day exploitation at reduced costs alters threat landscape dynamics. This necessitates coordinated defenses and ethical governance mechanisms that ensure technological advancement serves cybersecurity objectives. Our primary contributions include three key advances. First, we provide comprehensive systematization of the NeSy cybersecurity field. Second, we introduce the Grounding-Instructibility-Alignment (G-I-A) framework as an evaluation lens for systematic assessment. Third, we establish the first comprehensive dual-use analysis of both defensive innovations and autonomous offensive capabilities. The path toward realizing NeSy’s significant potential hinges on addressing three critical requirements: enhanced grounding mechanisms ensuring robust conceptual understanding of cybersecurity principles, comprehensive instructibility frameworks enabling effective human-AI collaboration across diverse expertise levels, and systematic alignment approaches guaranteeing consistency with defensive objectives and societal expectations.

The absence of standardized evaluation frameworks represents the field’s most pressing limitation, constraining systematic advancement and reproducible research essential for responsible development. Success demands coordinated efforts between research communities, practitioners, and policymakers. These efforts must establish evaluation standards, address implementation barriers, and develop governance frameworks that balance innovation with security imperatives. Ultimately, the substantial promise of NeSy cybersecurity systems will be realized through sustained commitment to responsible development principles that ensure these hybrid intelligence frameworks strengthen our collective cybersecurity infrastructure.

REFERENCES

- [1] F. FortiGuard Labs, “2025 global threat landscape report,” *Technical Report*, 2025.
- [2] M. A. Ferrag, F. Alwahedi, A. Battah, B. Cherif, A. Mechri, N. Tihanyi, T. Bisztray, and M. Debbah, “Generative ai in cybersecurity: A comprehensive review of llm applications and vulnerabilities,” *Internet of Things and Cyber-Physical Systems*, vol. 5, pp. 1–46, 2025.
- [3] J. Zhang, H. Bu, H. Wen *et al.*, “When llms meet cybersecurity: A systematic literature review,” *Cybersecurity*, vol. 8, p. 55, 2025.
- [4] W. Kasri, Y. Himeur, H. A. Alkhazaleh, S. Tarapiah, S. Atalla, W. Mansoor, and H. Al-Ahmad, “From vulnerability to defense: The role of large language models in enhancing cybersecurity,” *Computation*, vol. 13, no. 2, p. 30, 2025.
- [5] D. Li, Q. Li, Y. F. Ye, and S. Xu, “Arms race in adversarial malware detection: A survey,” *ACM Comput. Surv.*, vol. 55, no. 1, 2023.
- [6] M. Rodriguez, R. A. Popa, L. Liang, A. Wang, M. Rahtz, A. Kaskasoli, A. Dafoe, and F. Flynn, “A framework for evaluating emerging cyberattack capabilities of ai,” *arXiv preprint arXiv:2503.11917*, 2025.
- [7] R. Fang, R. Bindu, A. Gupta, Q. Zhan, and D. Kang, “Teams of llm agents can exploit zero-day vulnerabilities,” *arXiv preprint arXiv:2406.01637*, 2024.
- [8] H. Kong, D. Hu, J. Ge, L. Li, T. Li, and B. Wu, “Vulnbot: Autonomous penetration testing for a multi-agent collaborative framework,” *arXiv preprint arXiv:2501.13411*, 2025.
- [9] Z. Chen, J. Liu, Y. Shen, M. Simsek, B. Kantarci, H. T. Mouftah, and P. Djukic, “Machine learning-enabled iot security: Open issues and challenges under advanced persistent threats,” *ACM Computing Surveys*, vol. 55, no. 5, pp. 105:1–105:37, 2022.
- [10] M. Sarhan, S. Layeghy, M. Gallagher *et al.*, “From zero-shot machine learning to zero-day attack detection,” *International Journal of Information Security*, vol. 22, pp. 947–959, 2023.
- [11] Y. Guo, “A review of machine learning-based zero-day attack detection: Challenges and future directions,” *Computer Communications*, vol. 198, pp. 175–185, 2023.
- [12] B. Tafreshian and S. Zhang, “A defensive framework against adversarial attacks on machine learning-based network intrusion detection systems,” in *IEEE TrustCom*, 2024, pp. 2436–2441.
- [13] A. Grini, O. Taheri, B. El Khamlichi, and A. El Fallah-Seghrouchni, “Constrained network adversarial attacks: Validity, robustness, and transferability,” *arXiv preprint arXiv:2505.01328*, 2025.
- [14] M. Corporation, L. Wong, K. Manville, and ramthrunner, “Adversarial ml threat matrix (atlas),” *Technical Report*, 2020.
- [15] M. Pawlicki, R. Kozik, and M. Choraś, “A survey on neural networks for (cyber-) security and (cyber-) security of neural networks,” *Neural Computing and Applications*, vol. 500, no. C, pp. 1075–1087, 2022.
- [16] P. Hitzler, F. Bianchi, M. Ebrahimi, and M. K. Sarker, “Neural-symbolic integration and the semantic web,” *Semantic Web*, vol. 11, no. 1, pp. 3–11, 2020.
- [17] W. Wang, Y. Yang, and F. Wu, “Towards data-and knowledge-driven ai: A survey on neuro-symbolic computing,” *IEEE TPAMI*, vol. 47, no. 2, pp. 878–899, 2024.
- [18] B. Bhuyan, A. Ramdane-Cherif, R. Tomar *et al.*, “Neuro-symbolic artificial intelligence: A survey,” *Neural Computing and Applications*, vol. 36, pp. 12 809–12 844, 2024.
- [19] B. C. Colelough and W. Regli, “Neuro-symbolic ai in 2024: A systematic review,” *arXiv preprint arXiv:2501.05435*, 2025.
- [20] I. H. Sarker, H. Janicke, A. Mohsin, A. Gill, and L. Maglaras, “Explainable ai for cybersecurity automation, intelligence and trustworthiness in digital twin: Methods, taxonomy, challenges and prospects,” *ICT Express*, vol. 10, no. 4, pp. 935–958, 2024.
- [21] M. Gaur and A. Sheth, “Building trustworthy neurosymbolic ai systems: Consistency, reliability, explainability, and safety,” *AI Magazine*, vol. 45, no. 1, pp. 139–155, 2024.
- [22] A. Velasquez, N. Bhatt, U. Topcu, Z. Wang, K. Sycara, S. Stepputtis, S. Neema, and G. Vallabha, “Neurosymbolic ai as an antithesis to scaling laws,” *PNAS Nexus*, vol. 4, no. 5, p. pgaf117, 2025.
- [23] U. Jaimini, C. Henson, and A. Sheth, “Causal neurosymbolic ai: A synergy between causality and neurosymbolic methods,” *IEEE Intelligent Systems*, vol. 39, no. 3, pp. 13–19, 2024.
- [24] A. Rawal, A. Raglin, D. B. Rawat, B. M. Sadler, and J. McCoy, “Causality for trustworthy artificial intelligence: Status, challenges and perspectives,” *ACM Computing Surveys*, vol. 57, no. 6, p. 146, 2025.
- [25] H. Lei, Y. Ge, and Q. Zhu, “Adapt: A game-theoretic and neuro-symbolic framework for automated distributed adaptive penetration testing,” in *IEEE MILCOM*, 2024.
- [26] H. Xiong, Z. Wang, X. Li, J. Bian, Z. Xie, S. Mumtaz, and L. E. Barnes, “Converging paradigms: The synergy of symbolic and connectionist ai in llm-empowered autonomous agents,” *arXiv preprint arXiv:2407.08516*, 2024.
- [27] J. McHugh, K. Šekrst, and J. Cefalu, “Prompt injection 2.0: Hybrid ai threats,” *arXiv preprint arXiv:2507.13169*, 2025.
- [28] D. Ristea, V. Mavroudis, and C. Hicks, “Ai cyber risk benchmark: Automated exploitation capabilities,” *arXiv preprint arXiv:2410.21939*, 2024.
- [29] Y. Zhu, A. Kellermann, D. Bowman, P. Li, A. Gupta, A. Danda, R. Fang, C. Jensen, E. Ihli, J. Benn, J. Geronimo, A. Dhir, S. Rao, K. Yu, T. Stone, and D. Kang, “Cve-bench: A benchmark for ai agents’ ability to exploit real-world web application vulnerabilities,” in *ICML*, 2025.
- [30] J. Renkhoff, K. Feng, M. Meier-Doernberg, A. Velasquez, and H. H. Song, “A survey on verification and validation, testing and evaluations

- of neurosymbolic artificial intelligence,” *IEEE TAI*, vol. 5, no. 8, pp. 3765–3779, 2024.
- [31] Z. Wan, C.-K. Liu, H. Yang, R. Raj, C. Li, H. You, Y. Fu, C. Wan, A. Samajdar, and Y. C. Lin, “Towards cognitive ai systems: Workload and characterization of neuro-symbolic ai,” in *IEEE ISPASS*, 2024, pp. 268–279.
- [32] Z. Lu, I. Afridi, H. J. Kang, I. Ruchkin, and X. Zheng, “Surveying neuro-symbolic approaches for reliable artificial intelligence of things,” *Journal of Reliable Intelligent Environments*, vol. 10, no. 3, pp. 257–279, 2024.
- [33] D. H. Hagos and D. B. Rawat, “Neuro-symbolic ai for military applications,” *IEEE TAI*, vol. 5, no. 12, pp. 6012–6026, 2024.
- [34] B. Al-Sada, A. Sadighian, and G. Oliveri, “Mitre att&ck: State of the art and way forward,” *ACM Computing Surveys*, vol. 57, no. 1, pp. 12:1–12:37, 2024.
- [35] T. O’Brien, “Mitre att&ck labeling of cyber threat intelligence via llm,” *Technical Report*, 2025.
- [36] X. Zhao, R. Jiang, Y. Han, A. Li, and Z. Peng, “A survey on cybersecurity knowledge graph construction,” *Computers & Security*, vol. 136, p. 103524, 2024.
- [37] H. Alharbi, A. Hur, H. Alkahtani, and H. F. Ahmad, “Enhancing cybersecurity through autonomous knowledge graph construction by integrating heterogeneous data sources,” *PeerJ Computer Science*, vol. 11, p. e2768, 2025.
- [38] W. Cheng, T. Zhu, T. Chen, Q. Yuan, J. Ying, H. Li, C. Xiong, M. Li, M. Lv, and Y. Chen, “Crucialge: Reconstruct integrated attack scenario graphs by cyber threat intelligence reports,” *IEEE TDSC*, pp. 1–17, 2025.
- [39] J. Paul, W. M. Lim, A. O’Cass, A. W. Hao, and S. Bresciani, “Scientific procedures and rationales for systematic literature reviews (spar-4-slr),” *International Journal of Consumer Studies*, vol. 45, no. 4, pp. O1–O14, 2021.
- [40] F. M. Khan, M. Anas, and S. M. F. Uddin, “Anthropomorphism and consumer behaviour: A spar-4-slr protocol compliant hybrid review,” *International Journal of Consumer Studies*, 2023.
- [41] G. Marra, S. Dumančić, R. Manhaeve, and L. De Raedt, “From statistical relational to neurosymbolic artificial intelligence: A survey,” *Artificial Intelligence*, vol. 328, p. 104062, 2024.
- [42] B. Li, Z. Li, Q. Du, J. Luo, W. Wang, Y. Xie, S. Stepputtis, C. Wang, K. Sycara, P. Ravikumar *et al.*, “Logicity: Advancing neuro-symbolic ai with abstract urban simulation,” *NeurIPS*, vol. 37, pp. 69 840–69 864, 2024.
- [43] S. Neupane, J. Ables, W. Anderson, S. Mittal, S. Rahimi, and I. Banicescu, “Explainable intrusion detection systems (x-ids): A survey of current methods, challenges, and opportunities,” *IEEE Access*, vol. 10, pp. 112 392–112 415, 2022.
- [44] D. Onchis, C. Istina, and E. Hogeia, “A neuro-symbolic classifier with optimized satisfiability for monitoring security alerts in network traffic,” *Applied Sciences*, vol. 12, no. 22, p. 11502, 2022.
- [45] F. Camarda, A. De Paola, S. Drago, P. Ferraro, and G. Lo Re, “Managing concept drift in online intrusion detection systems with active learning,” in *CEUR Workshop*, vol. 3962, 2025.
- [46] S. Bader and P. Hitzler, “Dimensions of neural-symbolic integration: A structured survey,” *arXiv preprint arXiv:cs/0511042*, 2005.
- [47] T. R. Besold, A. d’Avila Garcez, S. Bader, H. Bowman, P. Domingos, P. Hitzler, K.-U. Kühnberger, P. M. V. Lima, L. de Penning, G. Pinkas, H. Poon, and G. Zaverucha, “Neural-symbolic learning and reasoning: A survey and interpretation,” in *Neuro-Symbolic AI*, 2021, pp. 1–51.
- [48] P. Shakarian, G. I. Simari, C. Baral, B. Xi, and L. Pokala, “Neuro symbolic reasoning and learning,” *Book*, 2023.
- [49] P. Hitzler, A. Dalal, M. S. Mahdaviinejad, and S. S. Norouzi, “Handbook on neurosymbolic ai and knowledge graphs,” *Book*, 2025.
- [50] D. Cunningham, M. Law, J. Lobo, and A. Russo, “The role of foundation models in neuro-symbolic learning and reasoning,” in *NeSy*, 2024, pp. 84–100.
- [51] C. Michel-Delétie and M. K. Sarker, “Neuro-symbolic methods for trustworthy ai: A systematic review,” *Neurosymbolic Artificial Intelligence*, 2024.
- [52] A. Zhou, X. Xu, R. Raghunathan, A. Lal, X. Guan, B. Yu, and B. Li, “Knowgraph: Knowledge-enabled anomaly detection via logical reasoning on graph data,” in *ACM CCS*, 2024, pp. 168–182.
- [53] A. V. Singh, E. Rathbun, E. Graham, L. Oakley, S. Boboila, A. Oprea, and P. Chin, “Hierarchical multi-agent reinforcement learning for cyber network defense,” *arXiv preprint arXiv:2410.17351*, 2024.
- [54] G. Grov, J. Halvorsen, M. W. Eckhoff, B. J. Hansen, M. Eian, and V. Mavroedidis, “On the use of neurosymbolic ai for defending against cyber attacks,” in *Neural-Symbolic Learning and Reasoning*, 2024, pp. 119–140.
- [55] C. S. Kalutharage, X. Liu, and C. Chrysoulas, “Neurosymbolic learning and domain knowledge-driven explainable ai for enhanced iot network attack detection and response,” *Computers & Security*, vol. 151, p. 104318, 2025.
- [56] R. C. Ontiveros, F. Giannini, M. Gori, G. Marra, and M. Dili-genti, “Grounding methods for neural-symbolic ai,” *arXiv preprint arXiv:2507.08216*, 2025.
- [57] Z. Li, Y. Yao, T. Chen, J. Xu, C. Cao, X. Ma, and J. Lü, “Softened symbol grounding for neuro-symbolic systems,” in *ICLR*, 2023.
- [58] R. Naidu and N. Kagalwalla, “Can causal (and counterfactual) reasoning improve privacy threat modelling?” *arXiv preprint arXiv:2207.09746*, 2022.
- [59] H. T. T. Tran, J. Sander, A. Cohen, B. Jalaian, and N. D. Bastian, “Neurosymbolic artificial intelligence for robust network intrusion detection: From scratch to transfer learning,” *arXiv preprint arXiv:2506.04454*, 2025.
- [60] M. Hosseini, A. Lomuscio, and N. Paoletti, “Ltl verification of memoryful neural agents,” in *AAMAS*, 2025.
- [61] Z. Liu, H. Li, J. Lu, G. Ma, X. Hong, G. Iacca, A. Kumar, S. Tang, and L. Wang, “Nature’s insight: A novel framework and comprehensive analysis of agentic reasoning through the lens of neuroscience,” *arXiv preprint arXiv:2505.05515*, 2025.
- [62] L. Luo, G. Zhang, H. Xu, Y. Yang, C. Fang, and Q. Li, “End-to-end neuro-symbolic reinforcement learning with textual explanations,” in *ICML*, 2024, pp. 33 533–33 557.
- [63] A. Samaddar, N. Potteiger, and X. Koutsoukos, “Out-of-distribution detection for neurosymbolic autonomous cyber agents,” in *IEEE ICAIC*, 2025.
- [64] A. Piplai, A. Kotal, S. Mohseni, M. Gaur, S. Mittal, and A. Joshi, “Knowledge-enhanced neurosymbolic artificial intelligence for cyber-security and privacy,” *IEEE IC*, vol. 27, no. 5, pp. 43–48, 2023.
- [65] L. Zhang, Q. Zhu, H. Ray, and Y. Xie, “Improving network threat detection by knowledge graph, large language model, and imbalanced learning,” *arXiv preprint arXiv:2501.16393*, 2025.
- [66] P. Falcari and F. Dainese, “Building a cybersecurity knowledge graph with cybergraph,” in *ACM/IEEE EnCyCriS*, 2024, pp. 29–36.
- [67] L. F. Sikos, “Cybersecurity knowledge graphs,” *Knowledge and Information Systems*, vol. 65, no. 9, pp. 3511–3531, 2023.
- [68] A. Bizzarri, C.-E. Yu, B. Jalaian, F. Riguzzi, and N. D. Bastian, “A synergistic approach in network intrusion detection by neurosymbolic ai,” *arXiv preprint arXiv:2406.00938*, 2024.
- [69] I.-H. Hsu, Z. Xie, K.-H. Huang, P. Natarajan, and N. Peng, “Am-pere: Amr-aware prefix for generation-based event argument extraction model,” *arXiv preprint arXiv:2305.16734*, 2023.
- [70] L. N. DeLong, R. F. Mir, and J. D. Fleuriot, “Neurosymbolic ai for reasoning over knowledge graphs: A survey,” *IEEE TNNLS*, vol. 36, no. 5, pp. 7822–7842, 2024.
- [71] E. Gilliard, J. Liu, and A. A. Aliyu, “Knowledge graph reasoning for cyber attack detection,” *IET Communications*, vol. 18, no. 4, pp. 297–308, 2024.
- [72] G. Liu, K. Lu, and S. Pi, “Graph neural networks embedded with domain knowledge for cyber threat intelligence entity and relationship mining,” *PeerJ Computer Science*, vol. 11, p. e2769, 2025.
- [73] R. Fieblinger, M. T. Alam, and N. Rastogi, “Actionable cyber threat intelligence using knowledge graphs and large language models,” in *IEEE EuroS&PW*, 2024, pp. 100–111.
- [74] L. Blaauwbroek *et al.*, “Learning guided automated reasoning: A brief survey,” in *Logics and Type Systems in Theory and Practice*, 2024, pp. 71–92.
- [75] J. Piepenbrock, M. Janota, J. Urban, and J. Jakubův, “First experiments with neural cvc5,” in *EPiC Series in Computing*, vol. 100, 2023, pp. 249–264.
- [76] L. Blaauwbroek, M. Olšák, J. Rute, F. I. S. Massolo, J. Piepenbrock, and V. Pestun, “Graph2tac: Online representation learning of formal math concepts,” in *ICML*, 2024.
- [77] W. Choi, J. Park, S. Ahn, D. Lee, and H. Woo, “Nesyc: A neuro-symbolic continual learner for complex embodied tasks in open domains,” in *ICLR*, 2025.
- [78] R. Chitnis, T. Silver, J. B. Tenenbaum, T. Lozano-Pérez, and L. P. Kaelbling, “Learning neuro-symbolic relational transition models for bilevel planning,” in *IEEE IROS*, 2022, pp. 4166–4173.
- [79] T. Hürten, J. F. Loevenich, F. Spelter, E. Adler, J. Braun, L. Moxon, Y. Gourlet, T. Lefeuvre, and R. R. F. Lopes, “Hierarchical multi-agent reinforcement learning for autonomous cyber defense in coalition networks,” in *IEEE MILCOM*, 2024, pp. 176–181.

- [80] A. Bizzarri, B. Jalaian, F. Riguzzi, and N. D. Bastian, "A neuro-symbolic artificial intelligence network intrusion detection system," in *IEEE ICCCN*, 2024, pp. 1–9.
- [81] Z. Wang, S. Vijayakumar, K. Lu, V. Ganesh, S. Jha, and M. Fredrikson, "Grounding neural inference with satisfiability modulo theories," *NeurIPS*, vol. 36, pp. 22 794–22 806, 2023.
- [82] Z. Lu, S. Siemer, P. Jha, F. Manea, J. Day, and V. Ganesh, "Z3-alpha: A reinforcement learning guided smt solver," *Technical Report*, 2023.
- [83] Y. Zhang, Z. Wei, X. Zhang, and M. Sun, "Using z3 for formal modeling and verification of fnn global robustness," *arXiv preprint arXiv:2304.10558*, 2023.
- [84] C. Barrett, C. Tinelli, H. Barbosa, A. Niemetz, M. Preiner, A. Reynolds, and Y. Zohar, "Satisfiability modulo theories: A beginner's tutorial," in *Formal Methods*, 2025, pp. 571–596.
- [85] S. Veronica, "Reasoning under threat: Symbolic and neural techniques for cybersecurity verification," *arXiv preprint arXiv:2503.22755*, 2025.
- [86] B. Rao, W. Eiers, and C. Lipizzi, "Neural theorem proving: Generating and structuring proofs for formal verification," *arXiv preprint arXiv:2504.17017*, 2025.
- [87] X. Zhang, "Graph neural networks in network security: From theoretical foundations to applications," in *AMCIS*, 2025.
- [88] K. Kurniawan, E. Kiesling, and A. Ekelhart, "Cykg-rag: Towards knowledge-graph enhanced retrieval augmented generation for cybersecurity," in *ISWC Workshop*, 2024.
- [89] T. Chen, C. Dong, M. Lv, Q. Song, H. Liu, and T. Zhu, "Apt-kgl: An intelligent apt detection system based on threat knowledge and heterogeneous provenance graph learning," *IEEE TDSC*, pp. 1–15, 2022.
- [90] A. Andrew, S. Spillard, J. Collyer, and N. Dhir, "Developing optimal causal cyber-defence agents via cyber security simulation," *arXiv preprint arXiv:2207.12355*, 2022.
- [91] V. Clairoux-Trepanier, I.-M. Beauchamp, E. Ruellan, M. Paquet-Clouston, S.-O. Paquette, and E. Clay, "The use of large language models (llm) for cyber threat intelligence (cti) in cybercrime forums," *arXiv preprint arXiv:2408.03354*, 2024.
- [92] E. Bethany, M. Bethany, J. A. Nolasco-Flores, S. Jha, and P. Najafirad, "Jailbreaking large language models with symbolic mathematics," *arXiv preprint arXiv:2409.11445*, 2024.
- [93] Z. Dong, Z. Zhou, C. Yang, J. Shao, and Y. Qiao, "Attacks, defenses and evaluations for llm conversation safety: A survey," in *NAACL-HLT*, 2024, pp. 6734–6747.
- [94] T. Nieponice, V. Valeros, and S. Garcia, "Aracne: An llm-based autonomous shell pentesting agent," *arXiv preprint arXiv:2502.18528*, 2025.
- [95] C. Curaba, D. D'Ambrosi, and A. Minisini, "Cryptoformaleval: Integrating large language models and formal verification for automated cryptographic protocol vulnerability detection," in *NeurIPS Workshop*, 2024.
- [96] G. Deng, Y. Liu, V. Mayoral-Vilches, P. Liu, Y. Li, Y. Xu, T. Zhang, Y. Liu, M. Pinzger, and S. Rass, "Pentestgpt: Evaluating and harnessing large language models for automated penetration testing," in *USENIX Security*, 2024, pp. 847–864.
- [97] J. Wahr us, "Jailbreaking chatgpt through prompt segmentation," *M.S. thesis, KTH*, 2024.
- [98] B. Lund, Z. Orhan, N. R. Mannuru *et al.*, "Standards, frameworks, and legislation for artificial intelligence (ai) transparency," *AI Ethics*, vol. 5, pp. 3639–3655, 2025.
- [99] M. Sloane and E. W llhorst, "A systematic review of regulatory strategies and transparency mandates in ai regulation in europe, the united states, and canada," *Data & Policy*, vol. 7, p. e11, 2025.
- [100] S. Paul, J. Yu, J. J. Dekker, A. Ignatiev, and P. J. Stuckey, "Formal explanations for neuro-symbolic ai," *arXiv preprint arXiv:2410.14219*, 2024.
- [101] I. Sharafaldin, A. H. Lashkari, and A. A. Ghorbani, "Toward generating a new intrusion detection dataset and intrusion traffic characterization," in *ICISSP*, 2018, pp. 108–116.
- [102] Y. Ren, Y. Xiao, Y. Zhou, Z. Zhang, and Z. Tian, "Cskg4apt: A cybersecurity knowledge graph for advanced persistent threat organization attribution," *IEEE TKDE*, vol. 35, no. 6, pp. 5695–5709, 2022.
- [103] M. T. Alam, D. Bhusal, L. Nguyen, and N. Rastogi, "Ctibench: A benchmark for evaluating llms in cyber threat intelligence," in *NeurIPS*, vol. 37, 2024, pp. 50 805–50 825.
- [104] L. E. Richards, J. Yaros, J. Babcock, C. Ly, R. Cosbey, T. Doster, and C. Matuszek, "On the promise for assurance of differentiable neuro-symbolic reasoning paradigms," *arXiv preprint arXiv:2502.08932*, 2025.
- [105] K. Kireev, B. Kulynych, and C. Troncoso, "Adversarial robustness for tabular data through cost and utility awareness," in *Workshop on Adversarial ML*, 2023.
- [106] P. N. Williams and K. Li, "Black-box sparse adversarial attack via multi-objective optimisation," in *IEEE CVPR*, 2023, pp. 12 291–12 301.
- [107] A. T. Bui, T. Le, H. Zhao, Q. H. Tran, P. Montague, and D. Phung, "Generating adversarial examples with task oriented multi-objective optimization," *Transactions on Machine Learning Research*, 2023.
- [108] Q. Tao, J. Liao, E. Zhang, and L. Li, "A dual robust graph neural network against graph adversarial attacks," *Neural Networks*, vol. 175, p. 106276, 2024.
- [109] M. H. Meng, G. Bai, S. G. Teo, Z. Hou, Y. Xiao, Y. Lin, and J. S. Dong, "Adversarial robustness of deep neural networks: A survey from a formal verification perspective," *IEEE TDSC*, pp. 1–1, 2022.
- [110] P. Vaishnavi, K. Eykholt, and A. Rahmati, "Transferring adversarial robustness through robust representation matching," in *USENIX Security*, 2022, pp. 2083–2098.
- [111] C. Xiang, S. Mahloujifar, and P. Mittal, "Patchcleanser: Certifiably robust defense against adversarial patches for any image classifier," in *USENIX Security*, 2022, pp. 2065–2082.
- [112] S. K. Bashir, R. Podder, S. Sreedharan, I. Ray, and I. Ray, "Resiliency graphs: Modelling the interplay between cyber attacks and system failures through ai planning," in *IEEE TPS-ISA*, 2024, pp. 256–263.
- [113] V. Galwaduge and J. Samarabandu, "Tabular diffusion based actionable counterfactual explanations for network intrusion detection," 2025.
- [114] E. J. Oughton, D. Ralph, R. Pant, E. Leverett, J. Copic, S. Thacker, R. Dada, S. Ruffle, M. Tuveson, and J. W. Hall, "Stochastic counterfactual risk analysis for the vulnerability assessment of cyber-physical attacks on electricity distribution infrastructure networks," *Risk Analysis*, vol. 39, no. 9, pp. 2012–2031, 2019.
- [115] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, "Language models are few-shot learners," in *NeurIPS*, vol. 33, 2020, pp. 1877–1901.
- [116] W. E. Forum and A. G. Alliance, "Artificial intelligence's energy paradox: Balancing challenges and opportunities," *World Economic Forum, Tech. Rep.*, Jan. 2025, in collaboration with Accenture.
- [117] A. de Vries, "The growing energy footprint of artificial intelligence," *Joule*, vol. 7, no. 10, pp. 2191–2194, 2023.
- [118] M. Sajid, K. R. Malik, A. Almogren *et al.*, "Enhancing intrusion detection: A hybrid machine and deep learning approach," *Journal of Cloud Computing*, vol. 13, p. 123, 2024.
- [119] M. Guastalla, Y. Li, A. Hekmati, and B. Krishnamachari, "Application of large language models to ddos attack detection," in *SmartSP*, 2024, pp. 83–99.
- [120] Torq, "Tired of Security Alert Fatigue? Stop Burnout with Hyperautomation," Accessed: Aug. 3, 2025. [Online]. Available: <https://torq.io/blog/cybersecurity-alert-fatigue/>
- [121] D. AI, "How to Address Cybersecurity Alert Fatigue with AI," Accessed: Aug. 3, 2025. [Online]. Available: <https://www.dropzone.ai/blog/how-to-address-cybersecurity-alert-fatigue-with-ai>
- [122] P. Rajivan and N. J. Cooke, "Information-pooling bias in collaborative security incident correlation analysis," *Human Factors*, vol. 60, no. 5, pp. 626–639, 2018.
- [123] V. Nikulshin and C. Talhi, "Effective ids under constraints of modern enterprise networks: Revisiting the optc dataset," in *IEEE CIOt*, 2024.
- [124] V. Tadi, "Quantitative analysis of ai-driven security measures: Evaluating effectiveness, cost-efficiency, and user satisfaction across diverse sectors," *Journal of Scientific and Engineering Research*, vol. 11, no. 4, pp. 328–343, 2024.
- [125] K. Faber, R. Corizzo, B. Snieszynski, and N. Japkowicz, "Lifelong continual learning for anomaly detection: New challenges, perspectives, and insights," *IEEE Access*, vol. 12, pp. 41 364–41 380, 2024.
- [126] G. Andresini, F. Pendlebury, F. Pierazzi, C. Loglisci, A. Appice, and L. Cavallaro, "Insomnia: Towards concept-drift robustness in network intrusion detection," in *ACM AISec*, 2021, pp. 111–122.
- [127] T. Zoppi, M. Gharib, M. Atif, and A. Bondavalli, "Meta-learning to improve unsupervised intrusion detection in cyber-physical systems," *ACM TCPS*, vol. 5, no. 4, p. 42, 2021.
- [128] D. Lopez-Paz and M. Ranzato, "Gradient episodic memory for continual learning," in *NeurIPS*, 2017, pp. 6470–6479.
- [129] A. Bizzarri, C.-E. Yu, B. Jalaian, F. Riguzzi, and N. D. Bastian, "Neuro-symbolic integration for open set recognition in network intrusion detection," in *AIxIA*, 2024, pp. 50–63.

- [130] A. D. Kent, L. M. Liebrock, and J. C. Neil, "Authentication graphs: Analyzing user behavior within an enterprise network," *Computers & Security*, vol. 48, pp. 150–166, 2015.
- [131] N. Moustafa, "A new distributed architecture for evaluating ai-based security systems at the edge: Network ton_iot datasets," *Sustainable Cities and Society*, vol. 72, p. 102994, 2021.
- [132] N. Moustafa and J. Slay, "Unsw-nb15: A comprehensive data set for network intrusion detection systems," in *IEEE MilCIS*, 2015, pp. 1–6.
- [133] X. Xie, K. Kersting, and D. Neider, "Neuro-symbolic verification of deep neural networks," in *IJCAI*, 2022, pp. 3622–3628.
- [134] O. G. Arreche, "Explainable ai methods for enhancing ai-based network intrusion detection systems," *Ph.D. thesis, Purdue Univ.*, 2024.
- [135] V. Z. Mohale and I. C. Obagbuwa, "A systematic review on the integration of explainable artificial intelligence in intrusion detection systems to enhancing transparency and interpretability in cybersecurity," *Frontiers in Artificial Intelligence*, vol. 8, 2025.
- [136] R. Kalakoti, R. Vaarandi, H. Bahşi, and S. Nömm, "Evaluating explainable ai for deep learning-based network intrusion detection system alert classification," in *ICISSP*, 2025, pp. 47–58.
- [137] H. B. Goodman and P. Rowland, "Deficiencies of compliancy for data and storage," in *NCS Research Track*, 2021, pp. 171–185.
- [138] O. Subasi, J. Cree, J. Manzano, and E. Peterson, "A critical assessment of interpretable and explainable machine learning for intrusion detection," 2024.
- [139] I. Rosenberg, A. Shabtai, Y. Elovici, and L. Rokach, "Adversarial machine learning attacks and defense methods in the cyber security domain," *ACM Computing Surveys*, vol. 54, no. 5, pp. 108:1–108:36, 2021.
- [140] K. van Liebergen, J. Caballero, P. Kotzias, and C. Gates, "A deep dive into the virustotal file feed," in *DIMVA*, 2023, pp. 155–176.
- [141] I. VirusTotal, Google, "Virustotal," *Technical Report*, 2025.
- [142] M. D. S. R. Team, "Microsoft AI Competition Explores the Next Evolution of Predictive Technologies in Security," Accessed: Jul. 29, 2025. [Online]. Available: <https://www.microsoft.com/en-us/security/blog/2018/12/13/microsoft-ai-competition-explores-the-next-evolution-of-predictive-technologies-in-security/>
- [143] A. Piplai, P. Ranade, A. Kotal, S. Mittal, S. N. Narayanan, and A. Joshi, "Using knowledge graphs and reinforcement learning for malware analysis," in *IEEE Big Data*, 2020, pp. 4562–4571.
- [144] R. Kerr, S. Ding, L. Li, and A. Taylor, "Accelerating autonomous cyber operations: A symbolic logic planner guided reinforcement learning approach," in *IEEE ICNC*, 2024, pp. 641–647.
- [145] N. Potteiger, A. Samaddar, H. Bergstrom, and X. Koutsoukos, "Designing robust cyber-defense agents with evolving behavior trees," in *IEEE ICAA*, 2024, pp. 1–10.
- [146] P. Li, S. Yao, J. S. Korich, C. Luo, J. Yu, Y. Cao, and J. Yang, "Automated static vulnerability detection via a holistic neuro-symbolic approach," *arXiv preprint arXiv:2504.16057*, 2025.
- [147] L. L. Mankali, O. Sinanoglu, and S. Patnaik, "Insight: Attacking industry-adopted learning resilient logic locking techniques using explainable graph neural network," in *USENIX Security*, 2024, pp. 91–108.
- [148] Z. Li, S. Dutta, and M. Naik, "Iris: Llm-assisted static analysis for detecting security vulnerabilities," in *ICLR*, 2025.
- [149] A. Mehra, A. Aßmuth, and M. Prieß, "Graph of effort: Quantifying risk of ai usage for vulnerability assessment," in *Cloud Computing*, 2025, pp. 17–24.
- [150] J. A. Kupsch and B. P. Miller, "Manual vs. automated vulnerability assessment: A case study," in *MIST Workshop*, vol. 6, 2009, pp. 83–97.
- [151] W. Charoenwet, P. Thongtanunam, V. T. Pham *et al.*, "Toward effective secure code reviews: An empirical study of security-related coding weaknesses," *Empirical Software Engineering*, vol. 29, no. 88, 2024.
- [152] L. Muzsai, D. Imolai, and A. Lukács, "Hacksynth: Llm agent and evaluation framework for autonomous penetration testing," *arXiv preprint arXiv:2412.01778*, 2024.
- [153] L. Wang, J. Wang, K. Jung, K. Thiagarajan, E. Wei, X. Shen, Y. Chen, and Z. Li, "From sands to mansions: Enabling automatic full-life-cycle cyberattack construction with llm," *arXiv preprint arXiv:2407.16928*, 2024.
- [154] ScienceSoft, "How Much Does Penetration Testing Cost? [+Calculator]," Accessed: Aug. 14, 2025. [Online]. Available: <https://www.scnsoft.com/security/penetration-testing/costs>
- [155] C. Brown, "How Much Does Penetration Testing Cost?" Accessed: Aug. 14, 2025. [Online]. Available: <https://www.vikingcloud.com/blog/how-much-does-penetration-testing-cost>
- [156] M. AI, "Llama 3.1–405b," *Technical Report*, 2024.
- [157] Q. Team, "Qwen2.5-32b," *Technical Report*, 2024.
- [158] OpenAI, "Gpt-4o system card," *Technical Report*, 2024.
- [159] M. Nyre-Yu, E. Morris, M. Smith, B. Moss, and C. Smutz, "Explainable ai in cybersecurity operations: Lessons learned from xai tool deployment," *Technical Report*, 2022.
- [160] M. W. Eckhoff, J. Halvorsen, B. J. Hansen, M. Eian, V. Mavroeidis, R. A. Chetwyn, G. Skjøtskift, and G. Grov, "Experimenting with neurosymbolic ai for defending against cyber attacks," *Neurosymbolic Artificial Intelligence Journal*, 2025, to appear.
- [161] Z. Xiang, L. Zheng, Y. Li, J. Hong, Q. Li, H. Xie, J. Zhang, Z. Xiong, C. Xie, C. Yang, D. Song, and B. Li, "Guardagent: Safeguard llm agents by a guard agent via knowledge-enabled reasoning," *arXiv preprint arXiv:2406.09187*, 2025.
- [162] R. Ben Halima, M. Hachicha, A. Jemal, and A. Hadj Kacem, "Mape-k patterns for self-adaptation in cyber-physical systems," *Journal of Supercomputing*, vol. 79, no. 5, pp. 4917–4941, 2023.
- [163] M. Corporation, "Mitre d3fend," *Technical Report*, 2024.
- [164] P. E. Kaloroumakis and M. J. Smith, "Toward a knowledge graph of cybersecurity countermeasures," *Technical Report*, 2020.
- [165] J. Huang and Q. Zhu, "Penheal: A two-stage llm framework for automated pentesting and optimal remediation," in *Workshop on Autonomous Cybersecurity*, 2023, pp. 11–22.
- [166] X. Shen, L. Wang, Z. Li, Y. Chen, W. Zhao, D. Sun, J. Wang, and W. Ruan, "Pentestagent: Incorporating llm agents to automated penetration testing," *arXiv preprint arXiv:2411.05185*, 2025.
- [167] P. Gao, X. Liu, E. Choi, S. Ma, X. Yang, and D. Song, "Threatkg: An ai-powered system for automated open-source cyber threat intelligence gathering and management," in *ACM LAMPS*, 2024, pp. 1–12.
- [168] Y. Cheng, O. Bajaber, S. A. Tsegai, D. Song, and P. Gao, "Ctinexus: Automatic cyber threat intelligence knowledge graph construction using large language models," *arXiv preprint arXiv:2410.21060*, 2025.
- [169] OASIS Cyber Threat Intelligence (CTI) Technical Committee, "Stix best practices guide version 1.0.0," OASIS Open, Committee Note 01, 2022.
- [170] S. Nalluri, M. M. Malyala, H. Kandagiri, and K. K. Kandagiri, "Nscti: A hybrid neuro-symbolic framework for ai-driven predictive cyber threat intelligence," in *IEEE ICCMSO*, 2025, pp. 14–21.
- [171] B. Jalaian and N. D. Bastian, "Neurosymbolic ai in cybersecurity: Bridging pattern recognition and symbolic reasoning," in *IEEE MIL-COM*, 2023.
- [172] M. Shao, S. Jancheska, M. Udeshi, B. Dolan-Gavitt, H. Xi, K. Milner, B. Chen, M. Yin, S. Garg, P. Krishnamurthy, F. Khorrami, R. Karri, and M. Shafique, "Nyu ctf bench: A scalable open-source benchmark dataset for evaluating llms in offensive security," in *NeurIPS*, vol. 37, 2024, pp. 57472–57498.
- [173] S. Xu, "The cybersecurity dynamics way of thinking and landscape," in *Proceedings of the 7th ACM Workshop on Moving Target Defense, MTD@CCS 2020, Virtual Event, USA, November 9, 2020*, H. Okhravi and C. Wang, Eds. ACM, 2020, pp. 69–80.
- [174] L. Gioacchini, M. Mellia, I. Drago, A. Delsanto, G. Siracusano, and R. Bifulco, "Autopenbench: Benchmarking generative agents for penetration testing," *arXiv preprint arXiv:2410.03225*, 2024.
- [175] B. Challita and P. Parrend, "RedTeamLLM: An agentic AI framework for offensive security," 2025.
- [176] J. Xu, J. W. Stokes, G. McDonald, X. Bai, D. Marshall, S. Wang, A. Swaminathan, and Z. Li, "AutoAttacker: A large language model guided system to implement automatic cyber-attacks," 2024.
- [177] M. Bhatt, S. Chennabasappa, Y. Li, C. Nikolaidis, D. Song, S. Wan, F. Ahmad, C. Aschermann, Y. Chen, D. Kapil *et al.*, "Cyberseceval 2: A wide-ranging cybersecurity evaluation suite for large language models," *arXiv preprint arXiv:2404.13161*, 2024.
- [178] D. Bhusal, M. T. Alam, L. Nguyen, A. Mahara, Z. Lightcap, R. Frazier, R. Fieblinger, G. L. Torales, B. A. Blakely, and N. Rastogi, "Secure: Benchmarking large language models for cybersecurity," *arXiv preprint arXiv:2405.20441*, 2024.
- [179] P. Jing, M. Tang, X. Shi, X. Zheng, S. Nie, S. Wu, Y. Yang, and X. Luo, "Secbench: A comprehensive multi-dimensional benchmarking dataset for llms in cybersecurity," *arXiv preprint arXiv:2412.20787*, 2024.
- [180] F. Jackson, "Cyber Attacks Are Up 47% in 2025 – AI is One Key Factor," Accessed: Jul. 29, 2025. [Online]. Available: <https://www.techrepublic.com/article/news-cyber-attacks-check-point/>
- [181] S. Alahmari and A. Alkharashi, "Privacy-aware federated learning framework for iot security using chameleon swarm optimization and self-attentive variational autoencoder," *CMES*, vol. 143, no. 1, pp. 849–873, 2025.

- [182] N. Albogami, “Intelligent deep federated learning model for enhancing security in internet of things enabled edge computing environment,” *Scientific Reports*, vol. 15, p. 4041, 2025.
- [183] G. Ayittey, “A security operations and analytics framework: Continuous detection and response,” *Ph.D. thesis, Univ. of East London*, 2025.
- [184] D. M. Amlashi, A. Voelz, and D. Karagiannis, “Artificial intelligence and internet of things: A neuro-symbolic approach for automated platform configuration,” *Neurosymbolic Artificial Intelligence*, vol. 1, 2025.
- [185] M. Garba, “Enhancing smart home iot security with a multi-layer neuro-symbolic ai framework,” *Technical Report*, 2025.
- [186] R. S. Hallyburton and M. Pajic, “Assured autonomy with neuro-symbolic perception,” *arXiv preprint arXiv:2505.21322*, 2025.
- [187] G. Singh, R. Tommasini, S. Bhatia, and R. Mutharaju, “Benchmarking neurosymbolic description logic reasoners: Existing challenges and a way forward,” *Neurosymbolic Artificial Intelligence*, vol. 1, p. 29498732251339943, 2025.
- [188] J. Ott, A. Ledaguenel, C. Hudelot, and M. Hartwig, “How to think about benchmarking neurosymbolic ai?” in *NESY Workshop*, 2023.
- [189] M. Gharaibeh and C. Papadopoulos, “Darpa 2009 intrusion detection dataset,” *Technical Report*, 2014.
- [190] M. M. Anjum, S. Iqbal, and B. Hamelin, “Analyzing the usefulness of the darpa optc dataset in cyber threat detection research,” in *ACM SACMAT*, 2021, pp. 27–32.
- [191] M. Tavallae, E. Bagheri, W. Lu, and A. A. Ghorbani, “A detailed analysis of the kdd cup 99 data set,” in *IEEE CISDA*, 2009, pp. 1–6.
- [192] U. o. N. B. Canadian Institute for Cybersecurity, “Cic-ids2017: Intrusion detection evaluation dataset,” *Technical Report*, 2017.
- [193] L. S. Lorello, M. Lippi, and S. Melacci, “The kandy benchmark: Incremental neuro-symbolic learning and reasoning with kandinsky patterns,” *Machine Learning*, vol. 114, p. 161, 2025.
- [194] S. Bortolotti, E. Marconato, T. Carraro, P. Moretini, E. van Krieken, A. Vergari, S. Teso, and A. Passerini, “A neuro-symbolic benchmark suite for concept quality and reasoning shortcuts,” *NeurIPS*, vol. 37, pp. 115 861–115 905, 2024.
- [195] F. Yan, S. Wen, S. Nepal, C. Paris, and Y. Xiang, “Explainable machine learning in cybersecurity: A survey,” *International Journal of Intelligent Systems*, November 2022.
- [196] R. Zhao, V. Melnychuk, J. Zhao, J. Wright, and N. Shadbolt, “Let’s measure the elephant in the room: Facilitating personalized automated analysis of privacy policies at scale,” *arXiv preprint arXiv:2507.14214*, 2025.
- [197] D. Berreby, “As use of a.i. soars, so does the energy and water it requires,” *Yale Environment* 360, Feb. 6 2024, accessed: Aug. 24, 2025. [Online]. Available: <https://e360.yale.edu/features/artificial-intelligence-climate-energy-emissions>
- [198] D. Patterson, J. Gonzalez, Q. Le, C. Liang, L. Munguia, D. Rothchild, D. So, M. Texier, and J. Dean, “Carbon emissions and large neural network training,” *arXiv preprint arXiv:2104.10350*, 2021.
- [199] J. Zhu, C. Lu, J. Li, and F.-Y. Wang, “Secure consensus control on multi-agent systems based on improved PBFT and Raft blockchain consensus algorithms,” *IEEE/CAA Journal of Automatica Sinica*, vol. 12, no. 7, pp. 1407–1417, 2025.
- [200] A. Tellache, A. Mokhtari, A. Amara Korba, and Y. Ghamri-Doudane, “Multi-agent reinforcement learning-based network intrusion detection system,” in *IEEE NOMS*, 2024, pp. 1–9.
- [201] E. Marconato, G. Bontempo, E. Ficarra, S. Calderara, A. Passerini, and S. Teso, “Neuro-symbolic continual learning: Knowledge, reasoning shortcuts and concept rehearsal,” in *ICML*, 2023, pp. 1–22.
- [202] E. Marconato, S. Teso, A. Vergari, and A. Passerini, “Not all neuro-symbolic concepts are created equal: Analysis and mitigation of reasoning shortcuts,” in *NeurIPS*, 2023, pp. 1–33.
- [203] T. Silver, “Neuro-symbolic learning for bilevel robot planning,” *Ph.D. thesis, MIT*, 2024.
- [204] M. Peralta, S. Mukhopadhyay, and R. Bharadwaj, “Counterfactually reasoning about security,” in *ACM SIN*, 2011, pp. 223–226.
- [205] V. Galwaduge and J. Samarabandu, “Tabular diffusion based actionable counterfactual explanations for network intrusion detection,” *arXiv preprint arXiv:2507.17161*, 2025.
- [206] T. Simonetto, S. Dyrnishi, S. Ghamizi, M. Cordy, and Y. Le Traon, “A unified framework for adversarial attack and defense in constrained feature space,” in *IJCAI*, 2022, pp. 1313–1319.
- [207] O. Bougzime, S. Jabbar, C. Cruz, and F. Demoly, “Unlocking the potential of generative AI through neuro-symbolic architectures: Benefits and limitations,” 2025.
- [208] A. H. Salem, S. M. Azzam, O. E. Emam *et al.*, “Advancing cyber-security: A comprehensive review of ai-driven detection techniques,” *Journal of Big Data*, vol. 11, p. 105, 2024.
- [209] T. Bilot, N. El Madhoun, K. Al Agha, and A. Zouaoui, “A survey on malware detection with graph representation learning,” *ACM Computing Surveys*, vol. 56, no. 11, 2024.
- [210] X. Wang, B. Wang, Y. Wu, Z. Ning, S. Guo, and F. R. Yu, “A survey on trustworthy edge intelligence: From security and reliability to transparency and sustainability,” *IEEE COMST*, vol. 27, no. 3, pp. 1729–1757, 2025.
- [211] D. Arp, E. Quiring, F. Pendlebury, A. Warnecke, F. Pierazzi, C. Wressnegger, L. Cavallaro, and K. Rieck, “Dos and don’ts of machine learning in computer security,” in *31st USENIX Security Symposium (USENIX Security 22)*, 2022, pp. 3971–3988.
- [212] N. Capuano, G. Fenza, V. Loia, and C. Stanzione, “Explainable artificial intelligence in cybersecurity: A survey,” *IEEE Access*, vol. 10, pp. 93 575–93 600, 2022.
- [213] B. Yan, C. Yang, C. Shi, Y. Fang, Q. Li, Y. Ye, and J. Du, “Graph mining for cybersecurity: A survey,” *ACM TKDD*, vol. 18, no. 2, pp. 1–52, 2023.
- [214] National Institute of Standards and Technology, “Artificial intelligence risk management framework (ai rmf 1.0),” National Institute of Standards and Technology, Tech. Rep., 2023, nIST AI 100-1. [Online]. Available: <https://nvlpubs.nist.gov/nistpubs/ai/nist.ai.100-1.pdf>
- [215] European Parliament and Council, “Regulation (eu) 2024/1689 of the european parliament and of the council of 13 june 2024 laying down harmonised rules on artificial intelligence (artificial intelligence act),” Official Journal of the European Union, July 12 2024, oJ L, 2024/1689. [Online]. Available: <http://data.europa.eu/eli/reg/2024/1689/oj>
- [216] International Organization for Standardization and International Electrotechnical Commission, “Information technology — artificial intelligence — overview of trustworthiness in artificial intelligence,” ISO/IEC JTC 1/SC 42, Tech. Rep. ISO/IEC TR 24028:2020, May 2020, edition 1.
- [217] R. Raman, R. Kowalski, K. Achuthan *et al.*, “Navigating artificial general intelligence development: Societal, technological, ethical, and brain-inspired pathways,” *Scientific Reports*, vol. 15, p. 8443, 2025.
- [218] S. B. Hakim, M. Adil, A. Velasquez, and H. H. Song, “Ansr-dt: An adaptive neuro-symbolic learning and reasoning framework for digital twins,” *arXiv preprint arXiv:2501.08561*, 2025.