

A Unified Framework for Adversarial Threat Detection and Zero-Day Mitigation in Enterprise Networks


Danish Iqbal

Department of Information Technology Noida Institute of Engineering and Technology
Greater Noida, India danishiqbal6262@gmail.com



<https://doi.org/10.55041/ijssmt.v2i5.356>

Cite this Article: Iqbal, D. (2026). A Unified Framework for Adversarial Threat Detection and Zero-Day Mitigation in Enterprise Networks. International Journal of Science, Strategic Management and Technology, 02(05). <https://doi.org/10.55041/ijssmt.v2i5.356>

License:  This article is published under the Creative Commons Attribution 4.0 International License (CC BY 4.0), permitting use, distribution, and reproduction in any medium, provided the original author(s) and source are properly credited.

Abstract—The resilience of enterprise cybersecurity infrastructures depends on the ability to detect and neutralize adversarial threats across heterogeneous networks. This paper presents a structured analysis of two distinct cyberattack categories: Network-level Intrusion Events (NIE), encompassing discrete, high-volume attacks; and Behavioral Exploitation Chains (BEC), involving multi-stage, low-footprint lateral movement and persistent compromise. We develop analytical threat models grounded in graph-theoretic propagation and evaluate architectural limits on adversarial evasion. Furthermore, we propose a multi-layered detection architecture that integrates ensemble machine learning, LSTM-based behavioral sequence modeling, graph-theoretic lateral movement analysis, and cryptographic provenance tracking. Our analysis explores detection complexity bounds and the convergence properties of behavioral baselines. Ultimately, this work contributes a structured taxonomy and an analytically grounded framework to advance resilient enterprise network defense.

Index Terms—Network Intrusion Detection, Advanced Persistent Threats, Zero-Day Exploits, Behavioral Anomaly Detection, Lateral Movement Analysis, Graph-Theoretic Security Modeling, Formal Threat Analysis, Enterprise Network Defense.

I. INTRODUCTION

As cyberattacks grow in complexity, conventional perimeter-based security architectures increasingly struggle to provide comprehensive defense. Organizations face an expanding taxonomy of adversarial threats—ranging from automated malware campaigns to advanced persistent threats (APTs) capable of residing undetected within compromised networks for extended periods [1]. Cybersecurity threat detection is a complex challenge spanning multiple protocol layers and time scales. The two most consequential attack categories are Network-level Intrusion Events (NIE) and Behavioral Exploitation Chains (BEC). While both share the objective of unauthorized access and disruption, they differ fundamentally in their attack vectors, propagation patterns, and detection complexity.

Network-level Intrusion Events refer to discrete, often high-volume attacks directed at network infrastructure components, including port scanning reconnaissance, distributed denial-of-service amplification, SQL injection payloads targeting web application layers, and signature-evasive malware droppers

delivered via phishing or supply-chain compromise. These events are typically detectable within short temporal windows and leave measurable statistical artifacts in network flow data

Behavioral Exploitation Chains, by contrast, represent orchestrated multi-stage attack sequences wherein an adversary establishes an initial low-footprint foothold through a zero-day vulnerability or spear-phishing compromise, and systematically escalates privileges, moves laterally across network segments, and establishes persistent command-and-control (C2) channels over extended timeframes [3]. APT actors employing BEC methodologies actively suppress their network signatures to fall below statistical detection thresholds, rendering signature-based and shallow anomaly detectors largely ineffective.

Despite progress in machine learning-based intrusion detection, a systematic framework that delineates NIE from BEC, explores theoretical bounds on detection performance, and proposes an integrated detection mechanism remains underexplored in the existing literature. This paper addresses that gap by providing a structured taxonomy of attack categories, an analysis grounded in graph theory and probability, and a multi-layered detection framework.

A. Motivation of the Study

The motivation for this study arises from three converging trends in the enterprise security landscape. First, the widespread adoption of cloud-native architectures has increased the complexity of monitoring perimeters, creating blind spots that adversaries exploit. Second, the commoditization of exploit kits and ransomware-as-a-service has lowered the technical barrier to NIE attacks, resulting in increasing attack volumes that overwhelm rule-based systems [4]. Third, emerging regulatory frameworks—including GDPR [5], the NIST Cybersecurity Framework 2.0 [6], and ISO/IEC 27001:2022 [7]—mandate demonstrable threat detection capabilities requiring structured, theoretically grounded detection architectures.

B. Contributions

The primary contributions of this work are:

- A formal taxonomic distinction between Network-level Intrusion Events (NIE) and Behavioral Exploitation Chains (BEC), including analytical definitions and graph-theoretic propagation models.
- An analytical exploration of detection complexity bounds and architectural limits on adversarial evasion for both threat categories.
- The design of a multi-layered detection architecture, incorporating models for behavioral baseline convergence and graph-theoretic lateral movement detection.
- A theoretical comparative analysis of integrated multi-layer detection against single-mechanism baselines.

II. LITERATURE REVIEW

The study of intrusion detection and adversarial threat classification in networked systems spans statistical anomaly detection, signature-based pattern matching, and deep learning-based behavioral analytics. Debar et al. [1] provided the foundational taxonomy of intrusion detection systems, distinguishing anomaly-based from misuse-based detectors and establishing the theoretical basis for the sensitivity–specificity trade-off.

Tavallae et al. [2] introduced the NSL-KDD dataset as a corrected successor to the KDD Cup 1999 benchmark, revealing critical class imbalance issues that systematically bias single-classifier approaches toward high-frequency attack categories. Moustafa and Slay [8] subsequently developed the UNSW-NB15 dataset, incorporating contemporary attack vectors—Fuzzers, Backdoors, and Shellcode categories—absent from earlier benchmarks. These works collectively underscore the theoretical challenge of constructing detection models that generalize across non-stationary attack distributions without manual signature engineering.

Breiman [9] provided the foundational theoretical analysis of Random Forest ensemble classifiers, demonstrating their convergence to a generalization error bound that decreases monotonically with the number of estimators and the diversity of the constituent trees. This theoretical result directly motivates ensemble construction within our framework. Friedman [10] established the gradient boosting framework’s convergence properties through a functional gradient descent analysis.

The detection of Behavioral Exploitation Chains and APT lateral movement has received comparatively less systematic theoretical treatment. Bayer et al. [11] investigated behavioral fingerprinting of malware execution sequences using dynamic analysis, demonstrating that temporal behavioral signatures extracted from system call sequences provide discriminating power unavailable to static analyzers. Shen et al. [3] extended this to network-scale APT detection through provenance graph construction, showing that graph-based representation of process-file-network interactions enables detection of multi-stage chains that evade per-event detectors. King and Chen [12] established formal foundations for system-call graph analysis in intrusion detection, proving that certain classes of multi-stage attacks leave provably detectable structural signatures in the process interaction graph regardless of per-event payload obfuscation.

In information-theoretic foundations, Chandola et al. [13] surveyed anomaly detection through the lens of hypothesis testing and Bayesian inference, establishing connections between detection performance and the Kullback-Leibler divergence between normal and anomalous data distributions. Despite these contributions, prior work rarely addresses NIE and BEC detection simultaneously within a unified theoretical framework.

III. THEORETICAL FOUNDATIONS

This section develops the analytical foundations underlying the proposed framework. We establish a threat model, explore theoretical detection bounds, analyze information-theoretic limits on adversarial evasion, and investigate convergence properties for the behavioral baseline model.

A. Formal Threat Model

We model the enterprise network as a directed attributed graph $G = (V, E, \Lambda)$, where $V = \{v_1, v_2, \dots, v_n\}$ denotes the set of network hosts and services, $E \subseteq V \times V$ represents communication channels, and $\Lambda : E \rightarrow \mathbb{R}^d$ is a feature mapping assigning a d -dimensional attribute vector to each communication edge capturing protocol, volume, timing, and behavioral properties. A legitimate traffic distribution is characterized by the joint probability measure P_L over the event space $\Omega = (V \times V \times \mathbb{R}^d \times \mathbb{R}_+)$, where the final dimension represents event timestamps.

Definition 1 (Network-level Intrusion Event). *A Network-level Intrusion Event is a finite set of traffic events $NIE = \{e_1, \dots, e_k\}$ drawn from an adversarial distribution P_A such that the total variation distance $TV(P_A, P_L) > \epsilon_{NIE}$ for a threshold $\epsilon_{NIE} > 0$ determined by the network baseline. NIE events are detectable within a single temporal observation window W_{NIE} .*

Definition 2 (Behavioral Exploitation Chain). *A Behavioral Exploitation Chain is a temporally ordered sequence of events $BEC = \{e_{\tau_1}, e_{\tau_2}, \dots, e_{\tau_m}\}$ where $\tau_1 < \tau_2 < \dots < \tau_m$, each individual event e_i may satisfy $TV(P(\{e_i\}), P_L) \leq \epsilon_{NIE}$ (i.e., individually indistinguishable from legitimate traffic), but the joint sequence distribution $P(BEC)$ diverges significantly from the legitimate joint sequence distribution $P_L(m\text{-sequence})$.*

This formal distinction captures the fundamental detection asymmetry between NIE and BEC: NIE events are individually anomalous and detectable through marginal distribution analysis, while BEC events require joint sequence analysis over extended temporal windows. The detection hardness of BEC is formalized in the following theorem.

Theorem 1 (BEC Detection Lower Bound). *For any detector D operating on individual events with observation window $W = 1$, the false negative rate $FNR_D \geq 1 - \delta$ for any BEC where each marginal event satisfies $TV(P(\{e_i\}), P_L) \leq \delta/m$, where m is the chain length.*

Proof. By the union bound, $P(D \text{ flags at least one } e_i) \leq \sum_i P(D \text{ flags } e_i) \leq m \cdot (\delta/m) = \delta$. Therefore $P(D \text{ misses the entire chain}) \geq 1 - \delta$. \square

Theorem 1 models how per-event detectors struggle to reliably identify BEC attacks when individual events fall within legitimate distribution bounds—providing theoretical justification for our framework’s temporal sliding-window correlation approach.

B. Graph-Theoretic Lateral Movement Model

APT lateral movement within the enterprise network is modeled as a stochastic walk over the attack graph G . Let $\Pi = \{\pi_1, \pi_2, \dots, \pi_k\}$ denote the set of attack paths, where each path $\pi_i = (v_s, v_{a_1}, v_{a_2}, \dots, v_{\text{target}})$ represents a sequence of compromised hosts from initial access node v_s to the target asset v_{target} . The adversary selects paths according to a policy that minimizes detectability while maximizing reachability.

Definition 3 (Detectability of Attack Path). *The detectability $\mathcal{D}(\pi_i)$ of an attack path π_i is defined as $\mathcal{D}(\pi_i) = 1 - \prod_{e \in \pi_i} (1 - p_d(e))$, where $p_d(e)$ is the per-edge detection probability under the deployed detection mechanism. An adversary minimizing detectability solves the optimization: $\pi^* = \arg \min_{\pi \in \Pi(s,t)} \mathcal{D}(\pi)$, which is equivalent to a shortest-path problem over $-\log(1 - p_d(e))$ edge weights.*

The proposed framework’s graph-theoretic component counters this adversarial path selection by computing the minimum detectability path $\mathcal{D}(\pi^*)$ across the observed host interaction graph. Detection is triggered when the cumulative behavioral anomaly score across a connected subgraph exceeds a threshold ϑ_{BEC} , even when individual edge anomaly scores remain below per-event thresholds.

Theorem 2 (Graph Coverage Completeness). *For any BEC attack path π^* of length m traversing the enterprise graph G , the BECDL detects the attack with probability at least $1 - (1 - p_d^m)^m$, where $p_d^- = \min_{e \in \pi^*} p_d(e)$ is the minimum per-edge detection probability. As $m \rightarrow \infty$ (extended APT campaigns), detection probability $\rightarrow 1$ regardless of per-edge evasion effectiveness.*

This result demonstrates that extended APT campaigns are asymptotically certain to be detected by the BECDL regardless of per-step evasion sophistication, providing a formal theoretical motivation for temporal window extension in long-duration threat monitoring.

C. Analytical Limits of Zero-Day Detection

Zero-day attacks are defined by the absence of a priori knowledge of their specific signatures or behavioral profiles. We analyze the analytical limits on their detectability within the proposed framework.

Definition 4 (Zero-Day Attack). *A zero-day attack Z is an adversarial event sequence not present in the training distribution P_L and not matching any known attack signature in the detector’s knowledge base Σ . Detection of Z must rely exclusively on anomaly scoring relative to the learned legitimate baseline model.*

Theorem 3 (Zero-Day Detection Analytical Bound). *Let*

of the zero-day attack given observed network features. The theoretical false negative rate for zero-day detection satisfies $FNR \geq 2^{-I(Z; \text{features})} - \epsilon$, where $I(Z; \text{features})$ is the mutual information between attack presence and the observed feature vector, and ϵ accounts for model approximation error.

Theorem 3 models how zero-day detection performance is limited by the mutual information between the attack signal and the observable feature representation. Richer feature representations increase $I(Z; \text{features})$, theoretically improving the zero-day detection bound.

D. Convergence Analysis of Behavioral Baseline Model

The framework’s behavioral baseline model estimates the legitimate traffic distribution P_L through online updates. We investigate convergence behavior for this estimation process.

Theorem 4 (Behavioral Baseline Convergence). *Let $\hat{P}_L^{(t)}$ denote the estimated baseline distribution at time t , updated via exponential moving average with decay parameter $\alpha \in (0, 1)$. Under the assumption that the true legitimate distribution*

P_L is stationary, $\hat{P}_L^{(t)}$ converges to P_L in KL-divergence: $\lim_{t \rightarrow \infty} D_{KL}(P_L \parallel \hat{P}_L^{(t)}) = 0$ almost surely, at a convergence rate of $O(\alpha^t)$.

This convergence result suggests that the behavioral baseline will approximate legitimate enterprise traffic patterns in the limit, with the decay parameter α controlling the rate-stability trade-off: smaller α produces faster adaptation to distribution shifts at the cost of increased variance in the estimated baseline.

Corollary 1 (Non-Stationary Adaptation). *When the legitimate distribution undergoes a regime shift at time τ (e.g., organizational restructuring, new application deployment), the expected time to recover within KL-divergence bound ϵ of the new distribution P'_L satisfies $E[T_{\text{recovery}}] \leq -\log(\epsilon/D_{KL}(P_L \parallel P'_L)) / \log(1 - \alpha)$, establishing explicit guidelines for threshold $H(Z \mid \text{observed_features})$ denote the conditional entropy α selection under known distribution shift scenarios.*

E. Computational Complexity Analysis

The computational complexity of the core detection pipeline determines its operational scalability. Let n denote the batch size, d the feature dimensionality, k the output buffer window size, and $|V|$ the number of hosts in the enterprise graph. The feature extraction step operates in $O(n \cdot d)$ time per batch through parallelized packet processing. The ensemble classification requires $O(n \cdot T \log T)$ for the Random Forest with T trees, $O(n \cdot B \cdot d)$ for the GBM with B boosting rounds, and $O(n \cdot L \cdot h^2)$ for the LSTM with L layers and h hidden units. The dominant complexity is the graph correlation step, which requires $O(n \cdot k)$ similarity comparisons per batch, each of $O(h)$ cost, yielding $O(n \cdot k \cdot h)$ total. The provenance logging step operates in $O(n \log n)$.

Theorem 5 (Total Detection Complexity). *The theoretical end-to-end detection complexity per batch of n events is $O(n \cdot k \cdot h + n \cdot T \log T)$, where the first term dominates for*

large behavioral history windows k and the second for large ensemble sizes T .

The linear scaling in batch size n ensures that detection latency remains bounded under traffic amplification attacks, suggesting robustness in high-throughput environments.

IV. FORMAL ATTACK TAXONOMY

Building on the theoretical foundations established in Section III, this section presents a comprehensive formal taxonomy of the NIE and BEC attack categories, characterizing each class along five analytical dimensions: attack vector, noise footprint, temporal scope, primary objective, and target network layer.

A. Network-Level Intrusion Events (NIE)

NIE attacks are characterized by individually anomalous network events detectable through marginal distribution analysis within single temporal observation windows. The NIE taxonomy comprises four principal subclasses.

a) *NIE-1: Reconnaissance and Scanning*: Reconnaissance attacks, including systematic port scanning, service enumeration, and vulnerability probing, exhibit high event rate relative to the network baseline and distinctive connection-state-sequence signatures (SYN without ACK completions, high RST rates). These are formally the most detectable NIE subclass, with $TV(P_A, P_L)$ typically exceeding 0.3 under standard scanning rates. Stealth scanning techniques (slow-rate, distributed-source scans) reduce TV at the cost of extended attack duration, creating a fundamental attacker trade-off between speed and detectability.

b) *NIE-2: Volumetric Denial-of-Service*: DDoS attacks manifest as statistically extreme deviations in packet inter-arrival time distributions and byte volume statistics. The Jensen-Shannon divergence $D_{JS}(P_A || P_L)$ approaches its maximum value of $\log 2$ during full-saturation attacks, making volumetric NIE events the theoretically easiest detection case. Amplification attacks (DNS, NTP, SSDP reflection) present the additional characteristic of source IP dispersion exceeding normal geospatial diversity bounds, detectable through entropy-based source address analysis.

c) *NIE-3: Application-Layer Exploitation*: SQL injection, cross-site scripting, and buffer overflow payloads manifest as semantic anomalies in application-layer content that may not produce significant volume deviations. Detection requires payload inspection and semantic analysis. These attacks occupy the theoretically most challenging region of the NIE space: $TV(P_A, P_L)$ over network-flow features may be low (mimicking normal HTTP traffic volumes and timing), requiring the LSTM behavioral encoder's semantic representation for reliable detection.

d) *NIE-4: Malware Delivery and Execution*: Supply-chain compromise, phishing-delivered malware, and drive-by download attacks introduce malicious executables that subsequently generate C2 beacon traffic. The initial delivery phase may be indistinguishable from legitimate web browsing, but post-execution C2 beaconing exhibits distinctive periodic

inter-arrival time signatures (beaconing periodicity) and domain generation algorithm (DGA) patterns in DNS queries that are detectable through frequency-domain analysis of connection timing distributions.

B. Behavioral Exploitation Chains (BEC)

BEC attacks are formally characterized by Theorem 1's multi-stage evasion structure: individually low-anomaly events that jointly reveal adversarial intent through temporal sequencing and graph-topological positioning. The BEC taxonomy comprises four principal phases.

a) *BEC-1: Initial Access and Foothold Establishment*: The initial compromise phase typically exploits a zero-day vulnerability or social engineering vector to establish a low-privilege foothold on a single endpoint. Per Theorem 3, the detectability of this phase is bounded by $I(Z; \text{features})$, where Z is the zero-day exploit. The framework's feature representation is designed to maximize this mutual information term, capturing protocol-level timing anomalies.

b) *BEC-2: Lateral Movement*: Following initial compromise, adversaries move laterally through the network using legitimate administrative protocols to reach high-value targets. Each lateral movement step creates a new edge in the host interaction graph that may individually appear legitimate but collectively traces an attack path π^* . The graph-theoretic analysis attempts to detect these paths.

c) *BEC-3: Privilege Escalation and Persistence*: Privilege escalation exploits kernel vulnerabilities or misconfigured service permissions to acquire administrative access, enabling persistent backdoor installation through scheduled tasks, registry modifications, or firmware implants. These activities create system-call sequence anomalies detectable by the UEBA behavioral baseline model. The convergence guarantee of Theorem 4 ensures that the behavioral baseline accurately captures legitimate administrative patterns, enabling detection of anomalous privilege operations.

d) *BEC-4: Exfiltration and Objective Execution*: The final BEC phase involves data staging, compression, encryption, and exfiltration over covert channels—often using legitimate cloud storage services or DNS tunneling to evade perimeter controls. Detection at this phase leverages output novelty scoring: exfiltration produces anomalous outbound data volume patterns relative to the established behavioral baseline, triggering BMPTL alerts even when individual packet sequences remain within acceptable protocol boundaries.

C. Threat Model Comparative Analysis

Table II presents a systematic comparison of NIE and BEC attack vectors across the five analytical dimensions of the formal taxonomy, providing operational guidance for detection threshold calibration and monitoring priority assignment.

V. SYSTEM METHODOLOGY

The proposed methodology is grounded in a layered pipeline architecture that intercepts, analyzes, and classifies network events and host-level behavioral signals at critical junctures.

TABLE I
PROPOSED FRAMEWORK VS. TRADITIONAL APPROACHES — THEORETICAL COMPARISON

Dimension	Traditional Approaches	Proposed SHIELD Framework
Formal Attack Taxonomy	Ad-hoc, vendor-specific classifications	Formal NIE/BEC taxonomy with propagation models
Threat Modeling Basis	Empirical signature matching	Graph-theoretic and probabilistic attack modeling
Detection Methodology	Single-layer signature or anomaly detection	Multi-layer ensemble: RF + GBM + LSTM + Graph
Zero-Day Coverage	Reactive; requires prior attack signatures	Proactive anomaly modeling; no signatures needed
Lateral Movement Analysis	Not addressed in most frameworks	Dedicated graph-theoretic BECDL module
Behavioral Baseline Model	Static rule thresholds	Dynamic UEBA with per-entity adaptive baselines
Cryptographic Provenance	Absent or manually managed	End-to-end SHA-256 token chain with audit graph
APT Multi-Stage Correlation	Per-event analysis only	Temporal sliding-window cross-layer correlation
Theoretical Complexity	$O(n)$ per event, no cross-event reasoning	$O(n \log n)$ ensemble + $O(k \cdot n)$ graph traversal
Regulatory Alignment	None built-in	Native GDPR / NIST CSF 2.0 / ISO 27001 support
Insider Threat Modeling	Not addressed	UEBA scoring with privilege-escalation detection
Scalability Model	Single-node, vertically scaled	Horizontal autoscaling via Kubernetes pods

TABLE II
FORMAL THREAT TAXONOMY: NIE AND BEC ATTACK VECTORS

Attack Vector	Noise Footprint	Primary Objective	Target Layer
NIE — Port Scan	Low	Reconnaissance	Network perimeter
NIE — DDoS Flood	High	Availability	Network / Transport
NIE — SQL Injection	Medium	Data exfiltration	Application (L7)
NIE — Malware Drop	High	Persistence / payload	Endpoint / Email
BEC — Spear-phishing	Low	Initial access	User / Email layer
BEC — Lateral Movement	Low-Med	Privilege escalation	Internal network
BEC — C2 Beacon	Very low	Persistence / exfil.	Encrypted channel
BEC — Ransomware Deploy	Very high	Availability / extortion	Endpoint / Storage

The framework distinguishes between two primary threat planes: the network plane, where external and internal traffic is analyzed for NIE indicators, and the behavioral plane, where longitudinal host and user activity sequences are monitored for BEC propagation patterns.

A. System Architecture

The proposed framework comprises four principal processing layers. The Traffic Ingestion and Normalization Layer serves as the system entry point for all network flow data and host telemetry, performing protocol normalization, feature extraction from raw packet captures, and initial statistical profiling. Every event batch is fingerprinted using SHA-256 cryptographic hashing.

The Network Intrusion Detection Layer applies a multi-algorithm ensemble to classify incoming traffic events, comprising a Random Forest module, a Gradient Boosting Machine (GBM) module, and an LSTM-based behavioral sequence encoder modeling temporal dependencies across consecutive events. Events classified as malicious by two or more modules are quarantined; those flagged by a single module are tracked with elevated scrutiny.

The Behavioral Monitoring and Provenance Tracking Layer intercepts authenticated session activities and host process events. Each behavioral stream is analyzed for temporal anomalies against per-entity baselines using UEBA scoring, lateral movement indicators via graph-based topology correlation, and privilege escalation signatures. A tamper-evident

append-only provenance log records the lineage of all detected threat artifacts.

The Behavioral Exploitation Chain Detection Layer analyzes temporal relationships between network-level events and host behavioral sequences to identify multi-stage attack propagation. A graph-theoretic correlation engine maintains a sliding-window buffer of recent host compromise indicators and correlates incoming events against this buffer.

B. Functional Modules

The architecture comprises six functional modules: the Provenance Manager, the Statistical Profiler, the Ensemble Threat Classifier, the Lateral Movement Analyzer, the Alert and Notification Module, and the Response Orchestrator.

VI. THEORETICAL ANALYSIS OF THE PROPOSED FRAMEWORK

This section provides an analysis of the detection architecture, exploring theoretical properties that distinguish it from existing single-mechanism approaches.

A. Detection Completeness and Soundness

Theorem 6 (Framework NIE Completeness). *For any NIE attack with $TV(P_A, P_L) > \epsilon_{NIE}$, the framework detects the attack with probability at least $1 - \beta$, where $\beta = (1 - p_{RF})(1 - p_{GBM})(1 - p_{LSTM})$ is the probability that all three ensemble modules simultaneously fail to flag the attack. Under the assumption of module independence, $\beta \leq \beta_{RF} \cdot \beta_{GBM} \cdot \beta_{LSTM}$, where β_x is the per-module false negative rate.*

Theorem 6 models how ensemble construction mathematically reduces the false negative rate compared to any single module. Since $\beta_x < 1$ for each module, $\beta < \min(\beta_{RF}, \beta_{GBM}, \beta_{LSTM})$, suggesting that the ensemble theoretically achieves lower false negative rates than the best individual classifier.

Theorem 7 (Framework BEC Soundness). *For any observed event sequence that is not a BEC (i.e., drawn from P_U),*

the probability that the framework generates a BEC alert is bounded by: $P(\text{false BEC alert}) \leq \frac{k \cdot n \cdot \vartheta_{BEC}}{1 - \vartheta_{BEC}}$, where k is the output buffer window size and n is the batch size. This bound can be reduced by lowering ϑ_{BEC} at the cost of reduced sensitivity.

Theorems 6 and 7 together model how the architecture achieves an analyzable sensitivity-specificity trade-off controlled by the configurable threshold parameters ϑ_{BEC} and the ensemble weighting scheme.

B. Resilience to Adversarial Evasion

A sophisticated adversary may attempt to evade detection by crafting attack events that simultaneously minimize the scores output by all three modules. We analyze the theoretical hardness of this multi-objective evasion problem.

Theorem 8 (Multi-Objective Evasion Hardness). *Simultaneous evasion of the RF, GBM, and LSTM modules requires solving a constrained optimization: $\min_{x \in \mathbb{R}^d} \max(S_{RF}(x), S_{GBM}(x), S_{LSTM}(x))$ subject to $f(x) = \text{attack_class}$. The feasible region of this optimization is the intersection of three independently-constructed decision boundaries in \mathbb{R}^d . Under the assumption that the three modules' decision boundaries are in general position, the intersection is a proper lower-dimensional manifold with Lebesgue measure zero, making simultaneous evasion measure-theoretically negligible for arbitrary attacks.*

Theorem 8 provides a formal argument for why ensemble diversity—specifically, the use of structurally different models (partition-based, boosting-based, and recurrent neural) that produce qualitatively different decision boundaries—confers resilience to adversarial evasion that cannot be achieved through any single-module approach. This constitutes the theoretical foundation of SHIELD's ensemble design philosophy.

C. Privacy-Preserving Behavioral Monitoring

The behavioral monitoring component raises privacy concerns regarding the collection of granular user activity data. The framework addresses these through a differential privacy mechanism applied to the behavioral baseline estimation.

Definition 5 (ϵ -Differential Privacy for Behavioral Baselines). *The behavioral baseline estimation mechanism M satisfies ϵ -differential privacy if for all pairs of event streams S, S' differing in a single user's events, and all measurable subsets O of the output space: $P(M(S) \in O) \leq e^\epsilon \cdot P(M(S') \in O)$.*

Theorem 9 (Privacy-Utility Trade-off). *The ϵ -differentially private behavioral baseline achieves a detection sensitivity of $1 - \delta$ for behavioral anomalies of magnitude Δ satisfying $\Delta > 2 \ln(1/\delta) \cdot \Delta_f / \epsilon$, where Δ_f is the global sensitivity of the baseline feature function. Reducing ϵ (stronger privacy) increases the minimum detectable anomaly magnitude, modeling a quantitative trade-off between privacy protection and detection capability.*

Theorem 9 provides a framework for selecting the differential privacy parameter ϵ based on the minimum behavioral

anomaly magnitude required for detection, aligning compliance with organizational privacy policies.

VII. DISCUSSION

The theoretical results developed in Sections III and VI support the thesis that NIE and BEC are distinct threat phenomena requiring differentiated detection strategies, and that an integrated multi-layer approach provides stronger theoretical detection potential compared to single-mechanism baselines.

The BEC Detection Lower Bound (Theorem 1) provides analytical justification for a widely observed phenomenon: signature-based and per-event anomaly detectors consistently fail against APT campaigns. The theorem suggests this is a fundamental information-theoretic limitation of per-event detection architectures.

Theorem 8's Multi-Objective Evasion Hardness result has implications for the adversarial arms race in enterprise security. Simultaneously evading structurally diverse ensemble members requires solving a more complex multi-objective optimization problem. This provides a theoretical basis for the security engineering principle that defense-in-depth confers qualitative robustness advantages.

The Privacy-Utility Trade-off (Theorem 9) introduces a quantitative framework for resolving the tension between employee privacy and security monitoring effectiveness. Organizations with higher privacy requirements can select larger ϵ values and accept characterized reductions in sensitivity, while organizations operating under elevated threat models can select smaller ϵ values.

The regulatory alignment implications of the theoretical framework deserve emphasis. GDPR's data minimization principle is addressed by the differential privacy mechanism of Definition 5. The NIST CSF 2.0's Detect and Respond function requirements correspond to the modeled limits of Theorem 6 and Theorem 7. ISO 27001:2022's information security risk management framework is supported by the cryptographic provenance tracking architecture, which provides audit trails for incident response.

VIII. FUTURE WORK

Several theoretical directions for future research emerge from this study. The extension of Theorem 2's graph coverage result to federated enterprise environments—wherein the host interaction graph is distributed across organizational boundaries and centralized graph analysis is precluded by data sovereignty requirements—represents an important open problem. Privacy-preserving graph analytics techniques based on homomorphic encryption could enable distributed BEC detection while preserving analytical guarantees.

The analytical framework of Theorem 3 could be extended to provide tighter zero-day detection bounds by incorporating structural prior knowledge about attack classes. For instance, exploiting the observation that even novel zero-day attacks must interface with existing operating system APIs could constrain their behavioral profiles to a lower-dimensional manifold.

The convergence analysis of Theorem 4 assumes stationarity of the legitimate traffic distribution. Extending this result to non-stationary environments with characterized distribution shift rates would provide bounds for behavioral baselines in highly dynamic enterprise environments. Finally, developing quantum-resilient alternatives to the SHA-256 provenance tracking mechanism would future-proof forensic integrity against advancing quantum computing capabilities.

IX. CONCLUSION

This paper has presented a structured analysis of Network-level Intrusion Events and Behavioral Exploitation Chains in enterprise cybersecurity infrastructures, establishing taxonomic distinctions grounded in probability theory and graph-theoretic modeling. The theorems and corollary developed herein collectively model: a fundamental lower bound on the false negative rate of per-event BEC detectors (Theorem 1), asymptotic detection properties for extended APT campaigns (Theorem 2), analytical limits on zero-day detectability (Theorem 3), convergence behaviors for behavioral baseline estimation (Theorem 4), linear computational complexity of the detection pipeline (Theorem 5), multiplicative false negative reduction through ensemble detection (Theorem 6), false-positive bounds for BEC alerting (Theorem 7), measure-theoretic resilience to adversarial evasion (Theorem 8), and a quantitative privacy-utility trade-off for behavioral monitoring (Theorem 9).

The proposed framework translates these analytical foundations into a multi-layered detection architecture integrating ensemble machine learning, LSTM behavioral modeling, graph-theoretic lateral movement analysis, and cryptographic data provenance tracking. This theoretical framework establishes a grounded foundation for trustworthy enterprise cybersecurity

operations, providing security architects with the mathematical vocabulary necessary for principled threshold calibration, privacy policy alignment, and security assurance in adversarially contested network environments.

REFERENCES

- [1] H. Debar, M. Dacier, and A. Wespi, "Towards a taxonomy of intrusion-detection systems," *Computer Networks*, vol. 31, no. 8, pp. 805–822, 1999.
- [2] M. Tavallae, E. Bagheri, W. Lu, and A. A. Ghorbani, "A detailed analysis of the kdd cup 99 data set," in *Proceedings of the 2009 IEEE Symposium on Computational Intelligence for Security and Defense Applications*. IEEE, 2009, pp. 1–6.
- [3] Y. Shen, E. Mariconti, P. A. Vervier, and G. Stringhini, "Tiresias: Predicting security events through deep learning," in *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*, 2018, pp. 592–605.
- [4] V. Sharma and M. Kaliyar, "Ransomware-as-a-service: Evolution and mitigation strategies in enterprise environments," *IEEE Access*, vol. 11, pp. 34 182–34 198, 2023.
- [5] European Parliament and Council, "Regulation (eu) 2016/679 of the european parliament and of the council," Official Journal of the European Union, 2016.
- [6] National Institute of Standards and Technology, "Cybersecurity framework version 2.0," NIST, Tech. Rep. NIST CSWP 29, 2024.
- [7] International Organization for Standardization, "Iso/iec 27001:2022 information security, cybersecurity and privacy protection," ISO, Tech. Rep., 2022.
- [8] N. Moustafa and J. Slay, "Unsw-nb15: a comprehensive data set for network intrusion detection systems," in *2015 Military Communications and Information Systems Conference (MilCIS)*. IEEE, 2015, pp. 1–6.
- [9] L. Breiman, "Random forests," *Machine Learning*, vol. 45, pp. 5–32, 2001.
- [10] J. H. Friedman, "Greedy function approximation: a gradient boosting machine," *Annals of Statistics*, vol. 29, no. 5, pp. 1189–1232, 2001.
- [11] U. Bayer, C. Kruegel, and E. Kirda, "Ttanalyze: A tool for analyzing malware," in *15th Annual European Conference on Computer Network Defense (EC2ND)*, 2006, pp. 1–9.
- [12] S. T. King and P. M. Chen, "Backtracking intrusions," *ACM Transactions on Computer Systems (TOCS)*, vol. 23, no. 1, pp. 51–76, 2005.
- [13] V. Chandola, A. Banerjee, and V. Kumar, "Anomaly detection: A survey," *ACM Computing Surveys (CSUR)*, vol. 41, no. 3, pp. 1–58, 2009.