

Private Yet Accurate: A Decentralized Approach to System Intrusion Detection

Jinghan Zhang*, Mati Ur Rehman*, Sharon Biju, Saleha Muzammil, Wajih Ul Hassan
University of Virginia
{qtk5be, wkw9be, fuu2ka, evz4sc, hassan}@virginia.edu

Abstract—We present MIRAGE, the first privacy-preserving Provenance-based IDS (PIDS) that integrates Federated Learning (FL) with graph representation learning to match centralized detection accuracy while preserving privacy and improving scalability. Building MIRAGE is non-trivial due to challenges in federating graph-based models across clients with heterogeneous logs, inconsistent semantic encodings, and temporally misaligned data. To address these challenges, MIRAGE introduces a novel process entity categorization-based ensemble, where specialized submodels learn distinct system behaviors and avoid aggregation errors. To enable privacy-preserving semantic alignment, MIRAGE designs a dual-server harmonization framework: one server issues encryption keys, and the other aggregates encrypted embeddings without accessing sensitive tokens. To remain robust to temporal misalignment across clients, MIRAGE employs inductive GNNs that eliminate the need for synchronized timestamps. Evaluations on DARPA datasets show that MIRAGE matches the detection accuracy of state-of-the-art PIDS and reduces network communication costs by 170×, processes datasets in minutes rather than hours.

I. INTRODUCTION

Intrusion Detection Systems (IDS) are crucial for countering Advanced Persistent Threats (APTs) in enterprises, as highlighted by major attacks [8, 9]. To strengthen defenses, many organizations turn to Managed Security Service Providers (MSSPs). A survey [7] of over 5,000 IT professionals found that 75% of companies use MSSPs. These providers integrate with client systems and typically configure them to transmit system logs to the cloud for centralized analysis. Figure 1 illustrates this MSSP architecture.

Recently, Provenance-based IDS (PIDS) [43, 48, 59, 62, 69, 78, 79, 83, 87] have proven highly effective by leveraging the rich contextual information in system logs. These systems transform system logs into provenance graphs and apply machine learning, particularly Graph Neural Networks (GNN), to learn benign behavior patterns. By monitoring these patterns, PIDS can detect deviations that may indicate potential security threats. Upon detecting such anomalies, the PIDS generates alerts for further investigation.

*Equal contribution.

Critical Limitations of Existing PIDS. Despite PIDS potential, the current operational mode of MSSPs and the state-of-the-art PIDS face significant challenges in enterprise settings, which are described below.

1. PRIVACY RISKS & CENTRALIZATION. Current PIDS depend on centralized infrastructure, requiring clients to transmit logs to a central server for aggregating large datasets and enabling deep learning models to capture benign behavior. This design introduces serious privacy risks [10, 35, 75], as logs often contain sensitive information, such as URLs, IP addresses, and application usage. These concerns are supported by a recent Datadog report [4]. Moreover, training on single-machine data is insufficient. Our experiments with FLASH [69] on the DARPA OpTC dataset [3] show a 40% F-score drop when using single-host data compared to multi-host data.

2. NETWORK OVERHEADS. Many PIDS face critical challenges with network overhead and bandwidth constraints. Transmitting large volumes of logs for intrusion detection imposes high costs on users and organizations. Modern systems can produce gigabytes of logs daily [44, 45]. Our analysis of FLASH [69] and KAIROS [23] using the OpTC dataset (Section V-C) highlights these issues. Organizations similar in scale to those in the OpTC dataset generate up to 1000 GB of logs daily, leading to significant network expenses. Users with limited bandwidth face difficulties uploading such large amounts efficiently.

3. SCALABILITY ISSUES. As the number of hosts increases, centralized PIDS suffer from log congestion, resulting in detection delays and high storage demands [28]. FLASH and KAIROS took 27.7 and 56.6 hours, respectively, to process a single day’s OpTC logs. These extended processing times compromise effective threat detection and response capabilities, particularly in large-scale environments.

Combine FL with PIDS: Opportunity & Challenges. A seemingly promising solution to mitigate privacy and scalability limitations in centralized PIDS is to integrate Federated Learning (FL) into their design. In such a federated PIDS framework, each client constructs its local provenance graph, encodes semantic attributes using local feature sets \mathcal{F}_i , and trains a local Graph Neural Network model GNN_i . These clients then transmit only their model updates, rather than raw logs, to a central server, where the updates are aggregated into a global model GNN_{global} . This approach preserves data locality and reduce bandwidth while enabling collaborative

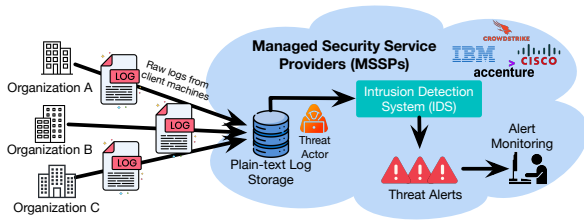


Fig. 1: The MSSP architecture for intrusion detection collects plain-text system logs in a centralized storage. These logs are then analyzed for intrusions. This setting risks privacy leaks if a threat actor or a curious analyst within the MSSP accesses the logs.

learning across multiple hosts. However, simply combining FL and PIDS introduces several new challenges that hinder effectiveness:

- C1 Feature Space Heterogeneity.** Inconsistent encoding of identical features across clients leads to difficulties in model convergence and reduces the effectiveness of federated averaging. PIDS, such as FLASH [69], utilize a Word2vec model to encode semantic attributes within a provenance graph. Due to the inherent randomness of the Word2vec algorithm, identical tokens t are encoded into different vectors $v_i(t)$ by each client i , using their local feature sets \mathcal{F}_i . This variability disrupts the convergence and efficacy of the aggregated global $\text{GNN}_{\text{global}}$ model from local client models leading to suboptimal anomaly detection performance, as detailed in [92].
- C2 Data Imbalance & Heterogeneity Among Clients.** Variations in data distribution and volume across clients, often non-IID [89], pose challenges in training a global model $\text{GNN}_{\text{global}}$ that performs uniformly well. Heterogeneous data resulting from different client applications can lead to suboptimal performance [68], and data imbalances can cause models to be biased towards clients with more extensive datasets, potentially overlooking unique patterns in less-represented clients [31]. This issue is critical for PIDS, as a biased $\text{GNN}_{\text{global}}$ may result in high false alarms, undermining detection accuracy.
- C3 Temporal Misalignment.** This issue arises in PIDS which uses Temporal Graph Networks (TGN) for threat detection. Aggregating temporal graph models from clients with temporally misaligned data challenges the creation of a cohesive global model. Systems like KAIROS [23] and ORTHRUS [49] employ TGN to trace the evolution of system provenance graphs. However, federated learning’s application across fragmented and misaligned client data impedes effective federated averaging, leading to improper alignment when aggregating models. This temporal constraint in TGN acts as a source of additional heterogeneity and hinders the formation of a cohesive global model $\text{GNN}_{\text{global}}$.

We present MIRAGE, a privacy-preserving PIDS that addresses the core challenges of applying FL to PIDS. Prior work [79] shows that graph-based models outperform traditional ML methods [25, 37] applied directly to system

logs [29, 81]. Building on this insight, MIRAGE introduces **Federated Provenance Graph Learning**, combining FL’s decentralized architecture with powerful graph representation learning to capture the intricate relationships among system entities. All major computations, including training, take place locally on client machines, leveraging their resources for real-time threat detection. This decentralized approach allows MIRAGE to scale effectively as more hosts are added, with each utilizing its own computing power and storage. *To the best of our knowledge, MIRAGE is the first system to achieve accurate, scalable, and privacy-preserving federated PIDS.*

1. PRIVACY-PRESERVING MODEL HARMONIZATION. To tackle **C1**, we implement a Word2vec harmonization scheme utilizing a dual-server architecture. In this setup, a central server issues encryption keys to clients, allowing them to securely encode Word2vec tokens. Subsequently, a utility server processes these encrypted tokens and corresponding vectors and to achieve a unified, privacy-preserving vector representation. This design ensures consistent semantic alignment across clients without exposing raw tokens, maintaining the privacy of sensitive process names, file paths, and network addresses.

2. ENTITY CATEGORIZATION ACROSS CLIENTS. To address **C2**, MIRAGE employs a novel categorization scheme for process entities, which enables consistent binning of similar activities across all clients in a privacy-preserving manner. The utility server collects encrypted process names and maps them into K_{cat} global categories using a mapping function. This standardized assignment ensures all clients use a shared category structure for training.

3. ENSEMBLE GNN LEARNING. Built on top of the categorization, we develop a GNN ensemble learning framework, where each submodel specializes in learning patterns from a specific process category. Clients align local process nodes to global bins, extract corresponding provenance subgraphs, and train category-specific models. These are aggregated across clients into K_{cat} global submodels, forming the ensemble $\text{GNN}_{\text{global}}$. This setup merges models with similar data distributions and preserves distinct local patterns.

4. DECENTRALIZED ANOMALY DETECTION. To avoid **C3**, MIRAGE employs an inductive GNN model [40], which has been shown by prior work [69, 79, 87] to offer good performance and is not affected by temporal dependencies. Each client performs decentralized inference using the global aggregated submodels, and a node is flagged as anomalous if it is consistently misclassified across all K_{cat} submodels. This detection leverages structural patterns without requiring temporal synchronization across clients, making it inherently robust to misaligned log timestamps. It also eliminates the need for global graph construction, allowing detection to proceed efficiently and independently on each client.

5. FORMAL PRIVACY ANALYSIS. MIRAGE provides formal privacy guarantees to support its dual-server architecture. We prove that (i) model updates sent to the central server provide indistinguishability under dataset perturbations (Central Server Privacy), (ii) the utility server cannot recover original tokens

from semantic vectors (Utility Server Privacy), and (iii) malicious clients cannot infer the presence of application-specific tokens from the global model (Client-Level Privacy).

We evaluated MIRAGE using open-source datasets from DARPA, including E3 [5], E5 [1], and OpTC [15]. The evaluation covers four dimensions: (1) Accuracy, (2) Efficiency and scalability, (3) Robustness, and (4) Privacy. For accuracy, we compare MIRAGE against a vanilla privacy-preserving PIDS baseline that combines FL with a state-of-the-art (SOTA) PIDS. We find that MIRAGE significantly outperforms this baseline. Since existing SOTA PIDS already achieve near-perfect detection accuracy, *the primary goal of MIRAGE is to demonstrate that a privacy-preserving PIDS can be built while maintaining similar performance*. Our results show that MIRAGE achieves an average precision of 96% and recall of 97%, matching the accuracy of centralized SOTA PIDS. To evaluate efficiency and scalability, we show that MIRAGE achieves a 170-fold reduction in network communication costs compared to centralized systems. Its decentralized design enables much faster inference, bounded only by the slowest client. MIRAGE completes the full OpTC dataset in a few minutes, while existing PIDS require several hours.

II. RELATED WORK

ML-based IDS. ML techniques are widely used in threat detection. ProvDetector [78] applies Doc2Vec [54] to provenance graphs; Attack2Vec [73] uses temporally aware embeddings. DeepAid [41] classifies anomalous traffic, ProGrapher [83] combines Graph2Vec [66] and TextRCNN [53], and StreamSpot [62] clusters graph features. Other systems [12, 13, 39, 63] explore varied embeddings; some focus on malware [22, 47, 93]. DeepLog [29] uses RNNs on logs, Euler [51] combines GNNs and RNNs, and MAGIC [48] uses masked graph learning. MIRAGE is the first to combine federated learning with provenance-based IDS, addressing privacy, scalability, and heterogeneity. DistDet [28] detects APTs using Hierarchical System Event Trees (HSTs) to summarize local system activity and reduce network overhead. However, this summarization comes at the cost of expressiveness and privacy. HSTs must be transmitted in plaintext to a central server, exposing sensitive execution traces without any formal privacy guarantees. More critically, HSTs encode flat, linear event sequences and lack the structural richness of provenance graphs, limiting their ability to model complex causal and multi-entity interactions. Anomalies are detected by identifying sequences absent from a reference benign HST, a simplistic approach that fails to generalize in dynamic environments and results in frequent false positives.

Federating Learning in Threat Detection. Few IDS employ federated learning, with most research centered on Network Intrusion Detection. Examples include [60], proposing FL for IoT threat detection, and [33], which introduces a differentially private system for industrial IoT. [56] presents an efficient network intrusion detection framework, while [38] addresses non-IID data issues in FL for intrusion detection. MIRAGE advances this area by being the first system to integrate FL

with graph-based learning techniques for host-based threat detection using a categorization based GNN ensemble framework and secure Word2vec harmonization. XFedGraph-Hunter [76] employs FL and GNN for network intrusion detection, while FedHE-Graph [61] introduces an FL-based intrusion-detection mechanism in a single-server setting. Entente [82] applies FL to graph-based network intrusion detection but does not address semantic feature alignment across clients. Unlike MIRAGE, however, these approaches lack semantic harmonization which is essential for leveraging semantic feature vectors and achieving detection performance on par with centralized host-based intrusion detection systems.

Cryptographic Techniques. Cryptographic techniques, such as Multi-Party Computation (MPC) [26] and Fully Homomorphic Encryption (FHE) [18], offer strong privacy guarantees. MPC allows multiple parties to compute a function over their inputs while keeping those inputs private, and FHE enables computations on encrypted data. However, these methods face significant challenges, including increased system complexity and scalability issues [20, 30, 34], particularly with large datasets [65]. For instance, host intrusion detection systems often process terabytes of logs, making the computational overhead of MPC or FHE impractical for real-time threat detection [55]. In contrast, MIRAGE combines a scalable encryption technique with federated learning to ensure both scalability and privacy preservation.

Privacy-preserving FL. PrivateFL [84] tackles the heterogeneity caused by differential privacy (DP) in FL systems through personalized data transformations to protect model updates from inference attacks. Similarly, [16] applies FL and DP to detect browser fingerprinting, and [27] introduces secure federated averaging using homomorphic encryption. Poseidon [72] utilizes multiparty cryptography for privacy-preserving neural network training, while PpeFL [77] adopts local DP for FL, addressing privacy and model performance issues. Unlike these methods, MIRAGE introduces a privacy-preserving multi-model server architecture that avoids DP-induced noise, ensuring robust performance while resisting inference attacks.

III. THREAT MODEL & ASSUMPTIONS

Our threat model assumes that the central server operates with integrity, conducting the federated averaging process without malicious objectives. However, we recognize the risk that the central server could compromise the privacy of client logs if raw system logs are transmitted to it [56, 60]. We also consider the possibility of a curious central server attempting membership inference attacks by utilizing the model weights.

Realism of Dual-Server (Non-Colluding) Architecture. Following precedent in cryptographic and federated learning systems [71, 80, 90], we assume the presence of a non-colluding, trusted utility server. This assumption is not only standard but also practically feasible. The utility server processes only encrypted tokens and performs no learning or aggregation itself, making its trust requirement minimal. It can be securely

deployed within a Trusted Execution Environment (TEE) [64], monitored by a trusted third-party mediator [14], or hosted as a secure cloud compute instance [2] directly managed by the organization. In this architecture, the central server is responsible for aggregating and coordinating tasks, while the utility server remains under the organization’s control. This separation allows the utility server to be operated by trusted internal IT teams, specialized third-party security firms, or a dedicated cloud provider, all of whom prioritize data privacy and encryption. Non-collusion could also be ensured through strict legal agreements [71].

Client-Side Threat Considerations. For individual clients, we expect that attackers could disguise their harmful activities within benign data, making it difficult to distinguish between legitimate and harmful actions. Our model also considers the threat posed by zero-day vulnerabilities. Despite these challenges, we assume that the activities of attackers will be detectable in the system’s records (system logs).

Log Integrity. In line with prior studies on data provenance [45, 48, 58, 69, 74, 78, 79, 83, 87], our approach relies on the provenance collection system’s ability to accurately record all system activities and changes. Additionally, we ensure the integrity of audit logs is maintained through the use of established tamper-resistant storage solutions, such as [11, 67, 88]. Similar to other PIDS works, we assume the absence of any attack activities during the training phase.

IV. DESIGN

A. Problem Statement

MIRAGE aims to detect anomalous nodes in client-specific provenance graphs $\{PGClient_1, PGClient_2, \dots, PGClient_N\}$, which are constructed from system logs on a set of decentralized clients \mathcal{C} . These anomalies represent deviations from benign behavior and are indicative of potential security threats.

To achieve this, MIRAGE employs a novel adaptation of Federated Learning (FL) to collaboratively train a set of global GNN models $GNN_{\text{global}} = \{GNN_1, GNN_2, \dots, GNN_{K_{cat}}\}$ without requiring raw log data to leave client machines. Each client locally optimizes its model by minimizing a loss function $\mathcal{L}_i(w)$ and returns its updated model weights w_i to the central server. The global model is updated iteratively using federated averaging: $w \leftarrow \frac{1}{N} \sum_{i=1}^N w_i$, where N represents the total number of clients. This decentralized approach preserves privacy while enabling the detection of threats across diverse and distributed system logs. The key notations for our system are summarized in Table I.

B. Provenance Graph Constructor

MIRAGE builds a system provenance graph on each client using local system logs. It leverages built-in logging tools, such as Linux Audit [6], which record detailed system activities. These include process executions, file accesses, and network events. MIRAGE constructs a graph where nodes represent entities like processes, files, and sockets. Edges correspond to syscalls capturing interactions between them. Each

TABLE I: Key notations used in MIRAGE’s design.

Notation	Definition
N	Total number of clients in federated learning.
K_{cat}	Number of categories for process entities.
$\mathcal{C} = \{C_1, C_2, \dots, C_N\}$	Set of all client machines.
$PGClient_i = (\mathcal{V}_i, \mathcal{E}_i)$	Provenance graph for client i , with nodes \mathcal{V}_i and edges \mathcal{E}_i .
$GNN_{\text{global}} = \{GNN_1, \dots, GNN_{K_{cat}}\}$	Set of global GNN models, one per category.
$w_j^{(r)}$	Weights of global GNN for category j after round r .
$\mathcal{P}_{\text{global}}$	Global set of unique process entities.
$\psi(p)$	Categorization map assigning process p to category \mathcal{C}_j .
y_v	True label of node v .
\hat{y}_v^j	Predicted label of node v by submodel j .
v_k'	Semantic embedding for attribute k , produced by a client’s local Word2vec model.
$M_{w2v\text{-harm}}$	Harmonized Word2vec model that converts contextual attributes into vector space.
$\mathcal{L}^{(r)}$	Loss function value after round r .

node is enriched with attributes such as process names, command lines, file names, and IP addresses. Prior work [23, 69] shows these attributes help the model learn normal system behavior.

C. Word2vec Model Generation

After each client machine $C_i \in \mathcal{C}$ generates its local provenance graph $PGClient_i = (\mathcal{V}_i, \mathcal{E}_i)$, the node attributes in \mathcal{V}_i are converted into feature vectors for the subsequent graph learning phase. Existing systems, such as FLASH [69], have demonstrated the effectiveness of leveraging semantic attributes of nodes to enhance detection performance. System logs contain rich attributes associated with various entities, including process names, file paths, and network IP addresses. These attributes must be transformed into vector space to serve as node features for the global GNN models.

For this encoding, we employ the Word2vec model which processes different semantic attributes based on node types: **1) Process:** Process name and command-line arguments. **2) File:** File path. **3) Socket:** IP addresses and port. For each node $v \in \mathcal{V}_i$, we extract its 1-hop subgraph by collecting all neighbors of v . We transform this node subgraph into a “document” by combining:

- 1) The semantic attributes of v .
- 2) The types of causal events captured along the edges.
- 3) The relevant attributes of its neighbor nodes.

These node-centric documents are then encoded into fixed-length vectors via Word2vec, which is trained on benign system logs. This approach effectively captures the semantic relationships between terms, producing dense embeddings that serve as node features for graph representation learning.

D. Privacy-preserving Model Harmonization

Addressing Token Variability Across Clients. Each client independently trains a Word2vec model using their local logs

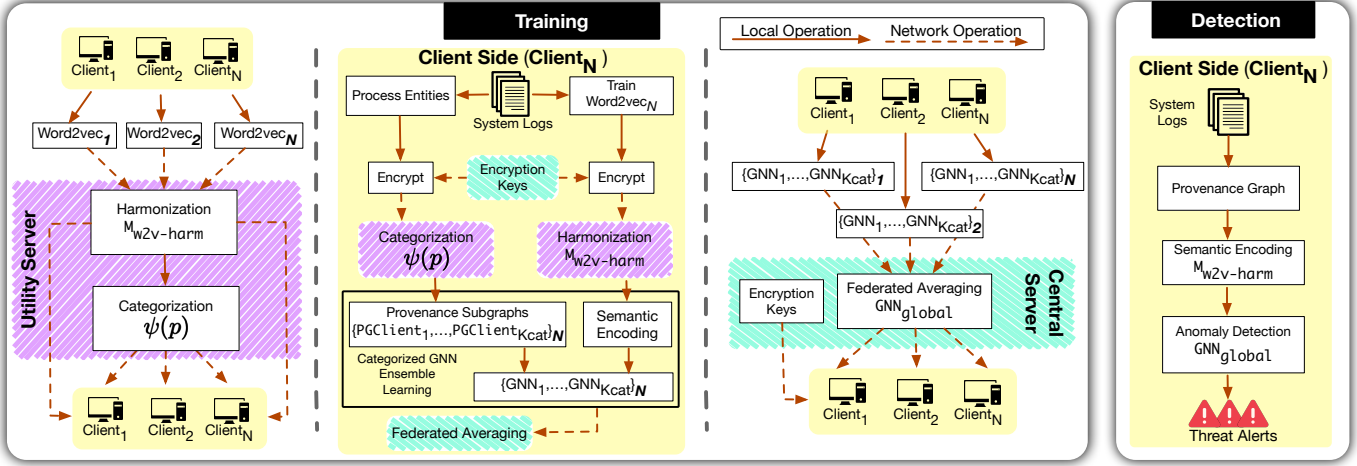


Fig. 2: High-level architecture of MIRAGE. In the training phase, our system builds local provenance graphs for each client and trains an ensemble of GNN models. Prior to this, we harmonize Word2vec models in a privacy preserving manner using encryption and dual-server architecture for semantic encoding. The local GNN models participate in federated learning to develop a global GNN_{global} model, which is then utilized for anomaly detection. Notations are defined in Table I.

for feature encoding. However, before these models can be utilized to encode text attributes effectively, we merge individual client Word2vec models into a unified model $M_{w2v\text{-harm}}$ for use across all clients. Without it, each client would produce a different encoding, v_a^i , for identical inputs, where i denotes the client. The variability in feature vectors, $\{v_a^1, v_a^2, \dots, v_a^N\}$, for the same attribute a across N clients, would compromise the consistency of client-based GNN models. To ensure uniformity, the feature vectors for overlapping attributes must be averaged across clients.

Our system harmonizes only the tokens that overlap in the Word2vec models across clients. Tokens that do not overlap are not averaged and remain unchanged. While clients share common activities due to standard system-level routines, some variations and patterns are unique to each client. Non-overlapping tokens preserve these unique patterns; however, they contribute to additional heterogeneity, which cannot be eliminated through harmonization. If these unique patterns are not accurately learned, they can lead to high false alarm rates. To capture these distinct local variations between clients, we have developed a novel ensemble GNN learning framework, detailed in Section IV-F, where each submodel specializes in distinct system patterns across clients, enhancing model precision as shown in Section V. *To better understand the degree of overlap in real-world settings, we now turn to an examination of client token overlap statistics.*

Overlap Statistics Across Clients. Our experiments with the OpTC dataset revealed that, on average, different hosts can have up to 75% overlap in process names, 41% in files, and 60% in network flows in the system logs. This finding aligns with related work [69], which shows that many system-level processes and files are common across different systems. In a centralized PIDS, these attributes are learned by a single semantic model, mapping them to the same embedding space. In contrast, a federated setting, where each client trains its own

model, can introduce model bias into the token embeddings, complicating the convergence for downstream GNN models using these vectors. To address this, it is essential to harmonize overlapping tokens to provide the model with a unified understanding of the semantic information present in system logs. *However, to achieve this harmonization securely and without compromising user privacy, we employ an encryption-based strategy, as described next.*

Harmonization with Encryption. Transferring tokens, containing sensitive data like process names, file names, and IP addresses, to a central server could breach user privacy. To mitigate this, we employ a trusted utility server. Initially, the central server uses Fernet symmetric key encryption [21, 46] to generate an encryption key, which is distributed to clients. Clients then encrypt their Word2vec token identifiers (vocabulary strings) using the shared key, $\tilde{t} = \text{Enc}(t, \text{key})$, and send pairs $(\tilde{t}, \mathbf{v}_t)$ (where \mathbf{v}_t is the plaintext embedding vector for token t) to the utility server. The utility server merges by matching \tilde{t} , averages the corresponding plaintext vectors, and dispatches the harmonized mapping back to clients (who may decrypt identifiers locally when needed), ensuring neither server learns the raw tokens assuming no collusion (Algorithm 2 in Appendix). Note that only token identifiers are encrypted; embedding vectors are transmitted in plaintext, which suffices because the servers never observe the underlying token strings and cannot link vectors to human-interpretable attributes without those identifiers.

E. Process Entity Categorization

During our initial experiments (Table V), we found that a single GNN_{global} cannot achieve good detection performance in an FL setting due to the diverse and heterogeneous distributions of clients data. The disparity in data distributions among N clients poses a significant challenge in effectively training a unified model that can generalize well across all clients, since

the global model struggles to capture the unique characteristics and patterns inherent in each client’s data. To address this limitation, we developed a comprehensive framework that organizes process entities across different clients into distinct groups. Each category is modeled by a dedicated submodel, enabling it to focus on the unique distribution and intricate patterns of its assigned category and thereby enhance overall detection performance.

The framework leverages a central utility server to orchestrate the categorization of process entities. Specifically, it proceeds in three main steps:

- 1) **Collection of Encrypted Process Names.** Each client $C_i \in \mathcal{C}$ transmits an encrypted list of its process names to the utility server. The server merges these into a global list of unique process entities, denoted by $\mathcal{P}_{\text{global}}$.
- 2) **Random Categorization.** The utility server applies the categorization map $\psi(p)$ to randomly partition all process names in $\mathcal{P}_{\text{global}}$ into K_{cat} categories. Although the partitioning is random, it is performed only once at the global level, ensuring that any process name p is consistently mapped to the same category j across all clients, i.e., $\psi(p) = j$.
- 3) **Distribution of Categories.** Finally, the resulting category assignments are disseminated back to each client, guaranteeing a unified and consistent mapping of processes into categorized bins across the federated network.

F. FL with Categorized GNN Ensemble Learning

Once the categorized bins of processes have been generated, we train specialized GNN models for each bin. Each GNN_{global} submodel focuses on a standardized subset of processes across clients, identified through entity categorization. This approach effectively captures distinct patterns and structures while reducing data heterogeneity. Our method proceeds in three main steps, as detailed in Algorithm 1.

Step 1: Initialization. As in Step 1 of Algorithm 1, the central server initializes the global GNN models, GNN_{global} , with random weights $w_j^{(0)}$. The server then transmits these initial weights $w_j^{(0)}$ to all clients in \mathcal{C} .

Step 2: Local Model Training. Each client C_i organizes its local processes into bins based on the assigned categories, $\psi(p)$ which it gets using the methodology detailed in section IV-E. From each category bin, the client constructs a *provenance subgraph* capturing the local relationships of processes within that category. The neighborhood hop length of these subgraphs matches the number of graph convolution layers in GNN_{global} , ensuring the model has sufficient neighborhood information for each node. Each client then leverages these categorized subgraphs, along with semantic feature vectors, to train the GNN models in an unsupervised manner. Similar to the approach in FLASH, the objective of each GNN submodel is to classify every node v into its corresponding type, yielding predictions \hat{y}_v^j . This localized, category-focused training ensures that each submodel effectively captures the unique structures and distributions within each category for the client’s data.

Algorithm 1: Training and federated aggregation of category-specific submodels.

Inputs : Categorization map ψ ; number of categories K_{cat} ; client set \mathcal{C} ; number of FL rounds R

Output: Trained global GNN models GNN_{global}

```

1 foreach category  $j \in \{1, \dots, K_{\text{cat}}\}$  do
2   |  $w_j^{(0)} \leftarrow \text{InitializeRandomWeights}()$ 
3 end
4 Broadcast  $\{w_1^{(0)}, \dots, w_{K_{\text{cat}}}^{(0)}\}$  to all clients  $C_i \in \mathcal{C}$ 
5 for  $r \leftarrow 1$  to  $R$  do
6   | foreach client  $C_i \in \mathcal{C}$  do
7     | foreach category  $j \in \{1, \dots, K_{\text{cat}}\}$  do
8       |  $\mathcal{V}_j^{(i)} \leftarrow \{p \in C_i \mid \psi(p) = j\}$ 
9       |  $\text{PG}_j^{(i)} \leftarrow \text{ExtractSubgraph}(\text{PG}_{\text{Client } i}, \mathcal{V}_j^{(i)})$ 
10      |  $w_j^{(i,r)} \leftarrow \text{LocalTrain}(w_j^{(r-1)}, \text{PG}_j^{(i)})$ 
11      end
12      Send  $\{w_1^{(i,r)}, \dots, w_{K_{\text{cat}}}^{(i,r)}\}$  to the Central Server
13    end
14    foreach category  $j \in \{1, \dots, K_{\text{cat}}\}$  do
15      |  $w_j^{(r)} \leftarrow \text{FederatedAveraging}(\{w_j^{(i,r)} \mid C_i \in \mathcal{C}\})$ 
16    end
17    Broadcast  $\{w_1^{(r)}, \dots, w_{K_{\text{cat}}}^{(r)}\}$  to all clients  $C_i \in \mathcal{C}$ 
18  end
19 return  $GNN_{\text{global}} \leftarrow \{w_1^{(R)}, w_2^{(R)}, \dots, w_{K_{\text{cat}}}^{(R)}\}$ 

```

Step 3: Federated Averaging and Category-Based Aggregation. As in Step 2 of Algorithm 1, the server collects the updated parameters from all N clients for each category-specific submodel. It then applies the federated averaging algorithm to merge these parameters, computing $\bar{w} = \frac{1}{N} \sum_{i=1}^N w_i$ for each category’s global submodel. Because every submodel is trained on processes belonging to the same category across all clients, the data distributions within each submodel are more homogeneous, effectively addressing data heterogeneity. This category-based aggregation integrates knowledge gained from multiple clients, leading to robust, specialized global models. It is important to note that the central server receives only the model weights from clients; no raw system log data is transmitted to it. Steps 1–3 of Algorithm 1 are repeated for R rounds, concluding once there is no further reduction in the training loss $\mathcal{L}^{(r)}$. The resulting ensemble of global submodels, each attuned to a distinct process category, provides a comprehensive solution that balances the diverse data distributions across clients while preserving privacy.

G. GNN-based Anomaly Detection

MIRAGE employs a node level detection methodology focusing on identifying irregular nodes through the comparison of their expected and observed types. This approach is grounded in a detailed analysis of both the surrounding structures and inherent properties of the nodes, to define normal pattern baselines for various node types. Typically, entities with malicious intentions display neighborhood structures and characteristics deviating from these established norms as shown in previous work [23, 69, 83]. In operational phases, the detection of anomalies that diverge from the pre-

established node distribution patterns often results in their misclassification. The emergence of nodes misclassified in the system’s output is indicative of potential security issues.

MIRAGE performs threat detection in a decentralized manner on clients’ provenance graphs ($\{PGClient_1, PGClient_2, \dots, PGClient_N\}$) in an organization. For a given provenance graph $PGClient_i$, MIRAGE uses the $GNN_{\text{global}} \{GNN_1, GNN_2, \dots, GNN_{K_{cat}}\}$ trained using federated learning. Each submodel performs inference on the client’s full provenance graph, utilizing the nodes’ features X_v and the graph’s adjacency matrix A to predict each node v ’s label \hat{y}_v^j . Let us define a misclassification indicator $\mathcal{M}(v, j)$ as

$$\mathcal{M}(v, j) = \begin{cases} 1, & \text{if } \hat{y}_v^j \neq y_v, \\ 0, & \text{otherwise.} \end{cases}$$

Accordingly, a node v is identified as an anomaly if it is misclassified by *all* submodels, i.e.,

$$A(v) = \begin{cases} 1, & \text{if } \sum_{j=1}^{K_{cat}} \mathcal{M}(v, j) = K_{cat}, \\ 0, & \text{otherwise.} \end{cases}$$

This indicates that none of the submodels recognize the neighborhood structure or features displayed by v , prompting its classification as an anomaly. To regulate the frequency of alerts, we define a threshold T similar to FLASH. This parameter sets a threshold on the likelihood of a classification being considered valid, with a higher value of T implying stronger confidence and increasing the probability of identifying anomalies.

V. EVALUATION

Our evaluation experiments are conducted on a machine running Ubuntu 18.04.6 LTS, equipped with a 10-core Intel CPU, NVIDIA RTX 2080 GPU, and 120 GB of memory. In our experiments, we set the federated learning rounds and the number of categorized GNN to 10 per host. Each model is trained for 20 epochs per round. We use regularization and dropout layers in our models to avoid overfitting. To evaluate MIRAGE, we address the following research questions:

- **RQ1.** How does MIRAGE compare to vanilla privacy-preserving PIDS in terms of detection performance? (Section V-A)
- **RQ2.** How does MIRAGE compare to existing PIDS in terms of detection performance? (Section V-B)
- **RQ3.** How scalable is MIRAGE in an enterprise-level setting with multiple host machines? (Section V-C)
- **RQ4.** What is the resource consumption of various components of MIRAGE and its end-to-end processing time on a client machine? (Section V-D)
- **RQ5.** What does the ablation study reveal about MIRAGE’s effectiveness across key factors? (Appendix B)

Implementation. MIRAGE is developed in Python with around 6000 lines of code. It leverages the PyTorch and Torch Geometric libraries to implement the federated provenance graph learning framework. The graph learning framework

uses the GraphSAGE [40] family of GNN. Our architecture consists of two graph convolution layers with a Tanh activation function in between. The last layer uses a softmax function to output class probabilities for the nodes. For implementing the Word2vec model, we employ the Gensim library. Secure communication between clients and the utility server is ensured through Python’s Cryptography module. The federated averaging, semantic vector harmonization, and entity categorization modules are implemented as individual Python functions on the central and utility servers.

Datasets. We have utilized the DARPA E3 [5], E5 [1], and OpTC [3] datasets for our evaluation. These datasets contain real-world APT attacks observed in enterprise networks. The attack traces they include are stealthy, span several days, and mirror the characteristics of actual APTs. Consequently, achieving strong detection accuracy on these datasets indicates that our system can deliver comparable performance in real-world deployments. Furthermore, these datasets incorporate logs of various sizes from Linux, FreeBSD, Android, and Windows operating systems. Our system’s robust detection accuracy across all datasets demonstrates its effective generalization to heterogeneous platforms with differing log sizes, reaching performance levels on par with state-of-the-art centralized PIDS. Notably, the datasets capture APT attacks of varying stealthiness, with the proportion of malicious nodes ranging from 0.05% in OpTC to 6% in E5, and E3. The low infiltration rate in OpTC underscores the system’s capacity for detecting highly stealthy adversaries, whereas the more widespread attacks in E3 and E5 highlight the resilience of our approach when confronted with a higher density of malicious nodes. Collectively, these datasets serve as a strong benchmark to evaluate the scalability and adaptability of our system. Each DARPA dataset is accompanied by ground truth documents that aid in distinguishing benign events from malicious ones. For this evaluation, we employ attack labels from existing systems such as ThreaTrace, KAIROS, and FLASH. ATLASv2 [70] is another recent dataset containing APT attack traces, but we did not evaluate it because it only includes data from two hosts, making it unrepresentative of a typical federated learning scenario.

Detectors for comparison. To benchmark our system, we compare against state-of-the-art PIDS. MAGIC [48] applies masked graph representation learning to identify threats. FLASH [69], another node-level system, leverages semantic feature vectors and an embedding recycling database for enhanced detection and efficiency. As shown in Table III, FLASH surpasses and thus serves as our primary baseline. We also include ORTHRUS [49] and KAIROS [23], which use temporal graph networks to capture system behavior over time. We exclude Streamspot [62] and ThreaTrace [79] as they are outperformed by FLASH and KAIROS. While DISDET [28], Prographer [83], and Shadewatcher [87] are notable, we exclude them due to closed-source code or proprietary components limiting reproducibility. Moreover, FLASH and ORTHRUS have already demonstrated superior performance

TABLE II: Comparison of MIRAGE against vanilla privacy-preserving PIDS as baseline.

Dataset	Baseline							MIRAGE						
	Precision	Recall	F1-score	TP	FP	FN	TN	Precision	Recall	F1-score	TP	FP	FN	TN
E3-CADETS	0.64	0.99	0.78	12846	7243	4	699725	0.97	0.99	0.98	12846	389	6	706,577
E3-TRACE	0.85	0.99	0.92	67357	11390	20	2404623	0.95	0.99	0.97	67357	3196	26	2,412,811
E3-THEIA	0.69	0.99	0.81	25314	11337	38	3493999	0.95	0.99	0.97	25313	1211	49	3,504,115
OpTC	0.36	0.99	0.53	649	1142	1	1286213	0.90	0.92	0.91	596	65	54	1,287,290
E5-CADETS	0.76	0.99	0.86	40098	12356	10	1258638	0.96	0.99	0.97	40096	1844	12	1,269,150
E5-THEIA	0.73	0.99	0.84	54810	20767	270	1250910	0.99	0.99	0.99	54803	277	197	1,271,480
E5-ClearScope	0.79	0.97	0.88	13670	3491	397	79857	0.99	0.97	0.99	13679	3	388	83345

over Prographer [83] and Shadewatcher [87]. We provide more details on why DISDET is unsuitable for comparison in Section II. It is important to note that, similar to existing works (e.g., KAIROS, Shadewatcher, and Prographer), MIRAGE considers only three node types in provenance graphs: *processes*, *files*, and *sockets*. However, in the E3 dataset, FLASH has also been evaluated using additional node types. Therefore, we executed FLASH using these three node types to report the results in Table III.

A. RQ1: Detection Performance Against a Vanilla Privacy-Preserving PIDS

We conducted experiments to analyze the detection performance of MIRAGE compared to a vanilla privacy-preserving PIDS, constructed by naively applying FL to the FLASH system. To simulate this setup, we operated FLASH in a decentralized manner: each client locally trained Word2vec to encode semantic features and then trained their own GNN models. These models were subsequently aggregated into a global model using the standard federated averaging algorithm.

We evaluated MIRAGE on the DARPA E3, E5, and OpTC datasets. E3 includes scenarios like Cadets, Theia, and Trace, each representing logs from a single host. We treated each scenario as a separate client, trained local GNN models, and performed 10 rounds of federated averaging. The global model was then evaluated on the corresponding attack logs. E5 is similarly structured but each scenario includes logs from three hosts. For the OpTC dataset, we used 10 hosts selected randomly for inclusion in our federated provenance graph learning experiment. Additionally, we conducted an enterprise-level analysis using all 1000 hosts from the OpTC dataset, the results of which are detailed in Section V-C. In all datasets, we trained on benign logs and evaluated on attack logs that appear chronologically after. For example, OpTC contains six days of benign and three days of attack logs, and we evaluated MIRAGE across all attack days. Detection metrics matched those used by prior node-level detectors, including ThreaTrace, MAGIC, and FLASH.

The results are shown in Table II. MIRAGE consistently outperforms the vanilla FL FLASH across all detection metrics. This performance gain is due to MIRAGE’s use of Word2vec harmonization and categorization-based ensemble learning, which effectively handle data heterogeneity. In contrast, the

vanilla FL FLASH lacks mechanisms to address such heterogeneity, resulting in degraded detection performance. We observed similar degradation with other centralized PIDS as well, which can be attributed to the absence of specialized mechanisms for handling heterogeneity in decentralized model training settings.

B. RQ2: Detection Performance Against PIDS

We conducted experiments to assess how MIRAGE compares with other systems in terms of detection performance. We used the same experimental setup and datasets described in RQ1. Table III shows that MIRAGE’s performance is comparable to MAGIC, FLASH, KAIROS, and ORTHRUS despite data heterogeneity, diverse log patterns, and data imbalance. KAIROS’s evaluation uses coarser time-window granularity compared to the node-level granularity of FLASH and MIRAGE, yet our results remain competitive. In some scenarios, MIRAGE exhibits slightly lower precision than centralized PIDS systems due to decentralized training, where each client learns from a limited and locally biased view of system behavior. Without full global context, certain rare benign patterns may appear anomalous to local models, leading to modest increases in false positives.¹ Beyond detection performance, MIRAGE offers privacy-preserving features and decentralized, scalable operation, emphasizing its practicality in real-world deployments.

C. RQ3: Scalability in Enterprise settings

We evaluated MIRAGE’s scalability using the OpTC dataset, which simulates a large enterprise environment with numerous hosts. We compared MIRAGE to centralized systems like FLASH and KAIROS, focusing on network and processing overhead.

Network Overhead: We estimate this cost for the FLASH and KAIROS system using the OpTC dataset. Each host within OpTC produces approximately 1GB of audit logs daily, equating to nearly one million audit events. For an organization with

¹We excluded KAIROS on E3 Trace and ORTHRUS from both E3 Trace and OpTC benchmarks due to missing results in their original publications. Further, unlike FLASH, both KAIROS and ORTHRUS are built on Temporal Graph Network (TGN) architectures, which are highly sensitive to time-window selection and require extensive hyper-parameter tuning to achieve stable performance. Despite multiple efforts, we were unable to reproduce satisfactory results on these datasets using their open-source implementations.

TABLE III: Comparison of MIRAGE with SOTA PIDS. Prec.: Precision; Rec.: Recall; F1.: F1-Score. While FLASH performs slightly better, MIRAGE offers strong privacy and scalability through decentralization. Refer to SOTA PIDS papers for their FP/FN details. The fraction in parentheses for Mirage indicates how many systems (out of the total compared) it outperforms or matches on that metric.

Dataset	MAGIC [48]			FLASH [69]			KAIROS [23]			ORTHRUS [49]			MIRAGE		
	Prec.	Rec.	F1.	Prec.	Rec.	F1.	Prec.	Rec.	F1.	Prec.	Rec.	F1.	Prec.	Rec.	F1.
E3-CADETS	0.94	0.99	0.97	0.99	0.99	0.99	0.80	1.00	0.89	0.52	0.37	0.43	0.97 (3/4)	0.99 (3/4)	0.98 (3/4)
E3-TRACE	0.99	0.99	0.99	0.99	0.99	0.99	-	-	-	-	-	-	0.95 (0/2)	0.99 (2/2)	0.97 (0/2)
E3-THEIA	0.98	0.99	0.99	0.99	0.99	0.99	0.91	1.00	0.95	0.81	0.40	0.54	0.95 (2/4)	0.99 (3/4)	0.97 (2/4)
OpTC	-	-	-	0.93	0.92	0.93	0.84	1.00	0.91	-	-	-	0.90 (1/2)	0.92 (1/2)	0.91 (1/2)
E5-CADETS	0.00	1.00	0.00	0.97	0.99	0.98	1.00	1.00	1.00	0.17	0.02	0.03	0.96 (2/4)	0.99 (3/4)	0.97 (2/4)
E5-THEIA	0.00	0.00	0.00	0.99	0.99	0.99	0.67	1.00	0.80	0.87	0.19	0.31	0.99 (4/4)	0.99 (3/4)	0.99 (4/4)
E5-ClearScope	0.00	1.00	0.00	0.99	0.96	0.98	0.67	1.00	0.80	0.33	0.08	0.13	0.99 (4/4)	0.97 (3/4)	0.99 (4/4)

1,000 hosts, the total daily log volume would be 1,000 GB. This data volume necessitates transmission over the network to a central server operating the PIDS. In contrast, MIRAGE achieves a significant reduction in these overheads. The only network expenses for MIRAGE arise from the transmission of the GNN and Word2vec models; the GNN model is roughly 13KB, while the average Word2vec model is 6 MB. Thus, the communication cost with the central server would be 12.70 MB, and for the utility server, it would be 5.86 GB. Hence, the network latency of model communications in MIRAGE is minimal compared to centralized systems where the raw system logs need to be sent over the network. Ultimately, MIRAGE results in a 170-fold decrease in communication costs compared to centralized PIDS.

Processing Overhead: Additionally, FLASH processes one million events in about 100 seconds, implying that processing events from 1,000 clients would necessitate approximately 27.7 hours. In contrast, KAIROS processes 57,000 events in 11.6 seconds, leading to a processing time of 204 seconds for one million events. Consequently, processing data from 1,000 clients with KAIROS would require around 56.6 hours. Compared to existing systems, MIRAGE processes client logs in a decentralized manner and the total processing time is bounded by the client with the most log data; it will only take approximately 3 minutes to run inference on the complete OpTC dataset. The central and utility servers conduct a simple mean operation on the models, taking only a few seconds.

D. RQ4: Resource consumption & Overheads

Resource Consumption. We conducted experiments to analyze the resource consumption of the central, utility server, and client-side modules of MIRAGE. We modeled the resource utilization on a client machine using different batches of audit events of varying sizes. For the central and utility servers, we studied resource consumption by varying the number of clients to understand the demands of federated averaging and semantic vector harmonization. The results, depicted in Figure 3, indicate that MIRAGE’s resource consumption is moderate. Specifically, MIRAGE can process up to 100,000 audit events simultaneously while consuming less than 900 MB of memory and utilizing less than 20% of CPU resources.

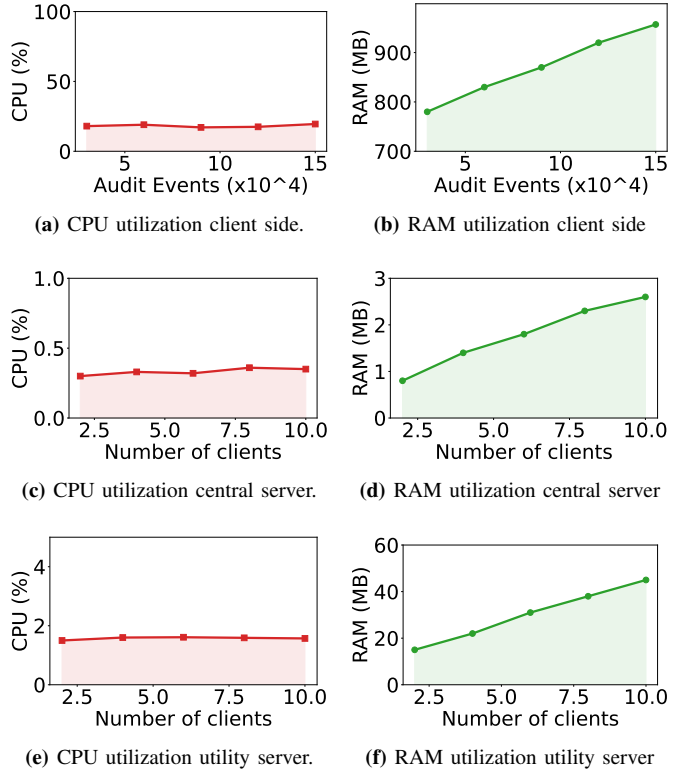


Fig. 3: Resource consumption of various components of MIRAGE.

This performance suggests that MIRAGE does not significantly burden the client machine, especially considering the typically low event throughput on such machines. Additionally, our analysis of the host data in the OpTC dataset shows that, on average, each client generates approximately 100,000 audit log events within a three-hour period. For the central and utility servers, the resource usage is minimal, demonstrating that our architecture is scalable and suitable for large organizations with many clients.

Processing Time Analysis. We conducted experiments to study the end-to-end processing time of our system for a client machine. For this, we used batches of audit events of various sizes, conducting end-to-end inference with MIRAGE

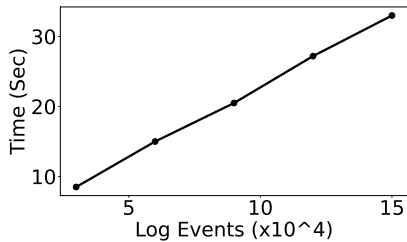


Fig. 4: Processing time for various audit event sizes evaluated using OpTC dataset.

to measure the time taken to process these events on a client machine. The results, illustrated in Figure 4, demonstrate that MIRAGE processes events with notable efficiency. For example, it requires approximately 23 seconds to process a batch of 100,000 events. Given our previous analysis of host logs in the OpTC dataset, which indicated that each host generates an average of 100,000 events in three hours, MIRAGE can process 24 hours worth of log data on a client in merely 3 minutes. This level of efficiency ensures that our system is highly effective, preventing any potential log congestion.

VI. DISCUSSION & FUTURE WORK

Unseen Tokens in Semantic Encoder. MIRAGE uses Word2vec models to encode semantic attributes, but these models only generate embeddings for previously seen tokens, degrading feature quality for new nodes with unseen tokens. Node representations in MIRAGE combine semantic attributes with neighboring system calls, but for unseen tokens, system calls dominate, as in ThreaTrace [79]. More frequent Word2vec retraining or adopting subword-level models like fastText [50], which build embeddings from subword units, can improve coverage. Tokenization methods like Byte-Pair Encoding [17] further aid by breaking unknown words into smaller components.

Scalability Bottlenecks. As MIRAGE scales to large enterprise networks, two key bottlenecks emerge: (i) the utility server may face throughput limits due to processing encrypted tokens from all clients, and (ii) aggregating frequent model updates from many clients risks state explosion, increasing bandwidth and computational costs. MIRAGE mitigates these via lightweight models, ensemble learning over categorized system activities, and submodel aggregation. To further scale, we plan to explore hierarchical federation [91], model compression [24], selective update aggregation [85], and robust communication protocols [52].

GraphSage Over Others. We adopt GraphSAGE as our base GNN due to its inductive capability and efficient minibatch training, both well-suited for federated settings. Compared to GCN, which requires full-graph access, and TGN, which imposes a temporal constraint that complicates federated learning, GraphSAGE offers a better balance of scalability, flexibility, and performance.

Alert Investigation. Validating alerts in systems like MIRAGE is critical for reliability and avoiding alert fatigue [42]. Tra-

ditionally, security analysts manually review alerts based on activities within local provenance graphs, but this process is time-consuming, error-prone, and lacks scalability and privacy. Privacy-preserving techniques such as Secure Multi-party Computation (SMC) [36], Homomorphic Encryption (HE) [86], and Zero-Knowledge Proofs (ZKP) [32] can address these challenges. SMC enables collaborative alert validation while preserving input privacy, HE ensures confidentiality during encrypted data analysis, and ZKP allows one party to prove an alert’s validity without revealing additional information. We identify privacy-preserving alert verification as a promising research direction. We leave it to future work to develop methods for privately sharing alert data with a central server, enabling analysts to perform attack analysis.

Comparison with Existing FL Solutions. We also evaluated our approach against existing federated learning solutions designed to address data heterogeneity, specifically FedProx [57] and FedOpt [19]. FedProx adds a proximal term to penalize deviations from the global model, while FedOpt employs a server-side optimizer for aggregating client updates. In our experiments across the OpTC, E3, and E5 datasets, both methods showed improvements over standard FedAvg but did not match the performance of MIRAGE. FedProx’s emphasis on global consistency limited its ability to capture client-specific patterns, resulting in higher false alarm rates. FedOpt, while effective at aggregating updates, struggled with the highly disparate updates common in heterogeneous system logs. Our ensemble-based approach, which combines process entity categorization with semantic vector harmonization, proved more effective at preserving client-specific patterns while still enabling meaningful aggregation across diverse environments.

VII. CONCLUSION

We introduced MIRAGE, a privacy-preserving intrusion detector that combines federated provenance graph learning with secure Word2vec harmonization and process entity categorization-based GNN ensemble training. MIRAGE addresses data heterogeneity, preserves client privacy, and reduces network overhead. Evaluations on large datasets show that MIRAGE matches centralized PIDS in accuracy and scales better.

VIII. ACKNOWLEDGMENT

We thank the reviewers for their insightful feedback. This material is based upon work supported by the National Science Foundation (NSF) under Award Nos. 2339483 and 2530655, and by the Commonwealth Cyber Initiative (CCI).

REFERENCES

- [1] DARPA Engagement 5. <https://github.com/darpa-i2o/Transparent-Computing/tree/master>.
- [2] What is an Instance in Cloud Computing? <https://scaleyourapp.com/what-is-an-instance-in-cloud-computing-a-thorough-guide/>.
- [3] DARPA OPTC. <https://github.com/FiveDirections/OpTC-data>.
- [4] Audit Logging Overview. <https://www.datadoghq.com/knowledge-center/audit-logging/>.
- [5] DARPA Engagement 3. <https://github.com/darpa-i2o/Transparent-Computing/blob/master/README-E3.md>.

- [6] The Linux audit daemon. <https://linux.die.net/man/8/auditd>.
- [7] IT security economics. <https://www.kaspersky.com/blog/it-security-economics-2020-part-4/>.
- [8] NotPetya Attack. https://en.wikipedia.org/wiki/Petya_and_NotPetya.
- [9] The SolarWinds Cyber-Attack: What You Need to Know. <https://www.cisecurity.org/solarwinds/>.
- [10] R. Aghili, H. Li, and F. Khomh. An empirical study of sensitive information in logs, 2024. URL <https://arxiv.org/abs/2409.11313>.
- [11] A. Ahmad, S. Lee, and M. Peinado. Hardlog: Practical tamper-proof system auditing using a novel audit device. In *2022 IEEE Symposium on Security and Privacy (SP)*. IEEE Computer Society, 2022.
- [12] S. Aljawarneh, M. Aldwairi, and M. B. Yassein. Anomaly-based intrusion detection system through feature selection analysis and building hybrid efficient model. *Journal of Computational Science*, 25, 2018.
- [13] A. Alsaheel, Y. Nan, S. Ma, L. Yu, G. Walkup, Z. B. Celik, X. Zhang, and D. Xu. Atlas: A sequence-based learning approach for attack investigation. In *USENIX Security Symposium*, 2021.
- [14] J. Alwen, J. Katz, Y. Lindell, G. Persiano, A. Shelat, and I. Visconti. Collusion-free multiparty computation in the mediated model. In *Advances in Cryptology-CRYPTO 2009: 29th Annual International Cryptology Conference, Santa Barbara, CA, USA, August 16-20, 2009. Proceedings*. Springer, 2009.
- [15] M. M. Anjum, S. Iqbal, and B. Hamelin. Analyzing the usefulness of the darpa optic dataset in cyber threat detection research. In *Proceedings of the 26th ACM Symposium on Access Control Models and Technologies*, 2021.
- [16] M. S. M. S. Annamalai, I. Bilogrevic, and E. De Cristofaro. Fp-fed: Privacy-preserving federated detection of browser fingerprinting. *arXiv preprint arXiv:2311.16940*, 2023.
- [17] A. Araabi, C. Monz, and V. Niculae. How effective is byte pair encoding for out-of-vocabulary words in neural machine translation? *arXiv preprint arXiv:2208.05225*, 2022.
- [18] F. Armknecht, C. Boyd, C. Carr, K. Gjøsteen, A. Jäschke, C. A. Reuter, and M. Strand. A guide to fully homomorphic encryption. *Cryptology ePrint Archive*, 2015.
- [19] M. Asad, A. Moustafa, and T. Ito. Fedopt: Towards communication efficiency and privacy preservation in federated learning. *Applied Sciences*, 10(8), 2020.
- [20] G. Asharov, Y. Lindell, T. Schneider, and M. Zohner. More efficient oblivious transfer and extensions for faster secure computation. In *Proceedings of the 2013 ACM SIGSAC conference on Computer & communications security*, 2013.
- [21] M. U. Bokhari and Q. M. Shallal. A review on symmetric key encryption techniques in cryptography. *International journal of computer applications*, 147(10), 2016.
- [22] S. S. Chakkaravarthy, D. Sangeetha, and V. Vaidehi. A survey on malware analysis and mitigation techniques. *Computer Science Review*, 32, 2019.
- [23] Z. Cheng, Q. Lv, J. Liang, Y. Wang, D. Sun, T. Pasquier, and X. Han. Kairos: Practical intrusion detection and investigation using whole-system provenance. In *IEEE Symposium on Security and Privacy (SP)*, 2024.
- [24] T. Choudhary, V. Mishra, A. Goswami, and J. Sarangapani. A comprehensive survey on model compression and acceleration. *Artificial Intelligence Review*, 53, 2020.
- [25] K. Chowdhary and K. Chowdhary. Natural language processing. *Fundamentals of artificial intelligence*, 2020.
- [26] R. Cramer, I. B. Damgård, et al. *Secure multiparty computation*. Cambridge University Press, 2015.
- [27] V. A. Dasu, S. Sarkar, and K. Mandal. Prov-fl: Privacy-preserving round optimal verifiable federated learning. In *Proceedings of the 15th ACM Workshop on Artificial Intelligence and Security*, 2022.
- [28] F. Dong, L. Wang, X. Nie, F. Shao, H. Wang, D. Li, X. Luo, and X. Xiao. *DISTDET: A Cost-Effective distributed cyber threat detection system*. In *32nd USENIX Security Symposium (USENIX Security 23)*, 2023.
- [29] M. Du, F. Li, G. Zheng, and V. Srikumar. DeepLog: Anomaly detection and diagnosis from system logs through deep learning. In *ACM Conference on Computer and Communications Security (CCS)*, 2017.
- [30] W. Du and M. J. Atallah. Secure multi-party computation problems and their applications: a review and open problems. In *Proceedings of the 2001 workshop on New security paradigms*, 2001.
- [31] M. Duan, D. Liu, X. Chen, R. Liu, Y. Tan, and L. Liang. Self-balancing federated learning with global imbalanced data in mobile systems. *IEEE Transactions on Parallel and Distributed Systems*, 32(1), 2020.
- [32] U. Fiege, A. Fiat, and A. Shamir. Zero knowledge proofs of identity. In *Proceedings of the nineteenth annual ACM symposium on Theory of computing*, 1987.
- [33] O. Friha, M. A. Ferrag, M. Benbouzid, T. Berghout, B. Kantarci, and K.-K. R. Choo. 2df-ids: Decentralized and differentially private federated learning-based intrusion detection system for industrial iot. *Computers & Security*, 127, 2023.
- [34] C. Gentry. Fully homomorphic encryption using ideal lattices. In *Proceedings of the forty-first annual ACM symposium on Theory of computing*, 2009.
- [35] S. Ghiasvand and F. M. Ciorba. Anonymization of system logs for preserving privacy and reducing storage. In *Advances in Information and Communication Networks: Proceedings of the 2018 Future of Information and Communication Conference (FICC), Vol. 2*. Springer, 2019.
- [36] O. Goldreich. Secure multi-party computation. *Manuscript. Preliminary version*, 78(110), 1998.
- [37] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11), 2020.
- [38] W. Guo, Z. Yao, Y. Liu, L. Zhang, L. Li, T. Li, and B. Wu. A new federated learning model for host intrusion detection system under non-iid data. In *2023 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*. IEEE, 2023.
- [39] M. Gyanchandani, J. Rana, and R. Yadav. Taxonomy of anomaly based intrusion detection system: a review. *International Journal of Scientific and Research Publications*, 2, 2012.
- [40] W. Hamilton, Z. Ying, and J. Leskovec. Inductive representation learning on large graphs. *Advances in neural information processing systems*, 30, 2017.
- [41] D. Han, Z. Wang, W. Chen, Y. Zhong, S. Wang, H. Zhang, J. Yang, X. Shi, and X. Yin. Deepaid: Interpreting and improving deep learning-based anomaly detection in security applications. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*, 2021.
- [42] W. U. Hassan, S. Guo, D. Li, Z. Chen, K. Jee, Z. Li, and A. Bates. NoDoze: Combatting threat alert fatigue with automated provenance triage. In *Network and Distributed System Security (NDSS)*, 2019.
- [43] W. U. Hassan, A. Bates, and D. Marino. Tactical provenance analysis for endpoint detection and response systems. In *IEEE Symposium on Security and Privacy (S&P)*. IEEE, 2020.
- [44] M. N. Hossain, J. Wang, R. Sekar, and S. D. Stoller. Dependence-preserving data compaction for scalable forensic analysis. In *USENIX Security Symposium*, 2018.
- [45] M. A. Inam, Y. Chen, A. Goyal, J. Liu, J. Mink, N. Michael, S. Gaur, A. Bates, and W. U. Hassan. Sok: History is a vast early warning system: Auditing the provenance of system intrusions. In *2023 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2023.
- [46] E. G. Ismail, A. Chahboun, and N. Raissouni. Fernet symmetric encryption method to gather mqtt e2e secure communications for iot devices. 2020.
- [47] T. Isohara, K. Takemori, and A. Kubota. Kernel-based behavior analysis for android malware detection. In *2011 seventh international conference on computational intelligence and security*. IEEE, 2011.
- [48] Z. Jia, Y. Xiong, Y. Nan, Y. Zhang, J. Zhao, and M. Wen. Magic: Detecting advanced persistent threats via masked graph representation learning. *arXiv preprint arXiv:2310.09831*, 2023.
- [49] B. Jiang, T. Bilot, N. El Madhoun, K. Al Agha, A. Zouaoui, S. Iqbal, X. Han, and T. Pasquier. Orthus: Achieving high quality of attribution in provenance-based intrusion detection systems. In *Security Symposium*

(USENIX Sec'25). USENIX, 2025.

- [50] A. Joulin, E. Grave, P. Bojanowski, and T. Mikolov. Bag of tricks for efficient text classification. *arXiv preprint arXiv:1607.01759*, 2016.
- [51] I. J. King and H. H. Huang. Euler: Detecting network lateral movement via scalable temporal link prediction. In *Network and Distributed System Security Symposium*, 2022.
- [52] M. Laclau, L. Renou, and X. Venel. Robust communication on networks. *arXiv preprint arXiv:2007.00457*, 2020.
- [53] S. Lai, L. Xu, K. Liu, and J. Zhao. Recurrent convolutional neural networks for text classification. In *Proceedings of the AAAI conference on artificial intelligence*, volume 29, 2015.
- [54] Q. Le and T. Mikolov. Distributed representations of sentences and documents. In *International conference on machine learning*. PMLR, 2014.
- [55] K. H. Lee, X. Zhang, and D. Xu. Loggc: garbage collecting audit log. In *CCS*, 2013.
- [56] J. Li, X. Tong, J. Liu, and L. Cheng. An efficient federated learning system for network intrusion detection. *IEEE Systems Journal*, 2023.
- [57] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith. Federated optimization in heterogeneous networks. *Proceedings of Machine learning and systems*, 2, 2020.
- [58] Q. Liu, M. Shoaib, M. U. Rehman, K. Bao, V. Hagenmeyer, and W. U. Hassan. Accurate and scalable detection and investigation of cyber persistence threats. *arXiv preprint arXiv:2407.18832*, 2024.
- [59] Q. Liu, K. Bao, W. U. Hassan, and V. Hagenmeyer. Hades: Detecting active directory attacks via whole network provenance analytics. *IEEE Transactions on Dependable and Secure Computing*, 2025.
- [60] D. Man, F. Zeng, W. Yang, M. Yu, J. Lv, and Y. Wang. Intelligent intrusion detection based on federated learning for edge-assisted internet of things. *Security and Communication Networks*, 2021, 2021.
- [61] A. A. Mansour Bahar, K. S. Ferrahi, M.-L. Messai, H. Seba, and K. Amrouche. Fedhe-graph: Federated learning with hybrid encryption on graph neural networks for advanced persistent threat detection. In *Proceedings of the 19th International Conference on Availability, Reliability and Security*, pages 1–10, 2024.
- [62] E. Manzoor, S. M. Milajerdi, and L. Akoglu. Fast memory-efficient anomaly detection in streaming heterogeneous graphs. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2016.
- [63] Z. K. Maseer, R. Yusof, N. Bahaman, S. A. Mostafa, and C. F. M. Foozy. Benchmarking of machine learning for anomaly based intrusion detection systems in the cicsids2017 dataset. *IEEE access*, 9, 2021.
- [64] F. McKeen, I. Alexandrovich, I. Anati, D. Caspi, S. Johnson, R. Leslie-Hurd, and C. Rozas. Intel® software guard extensions (intel® sgx) support for dynamic memory management inside an enclave. In *Proceedings of the Hardware and Architectural Support for Security and Privacy 2016*. 2016.
- [65] A. J. Menezes, P. C. Van Oorschot, and S. A. Vanstone. *Handbook of applied cryptography*. CRC press, 2018.
- [66] A. Narayanan, M. Chandramohan, R. Venkatesan, L. Chen, Y. Liu, and S. Jaiswal. graph2vec: Learning distributed representations of graphs. *arXiv preprint arXiv:1707.05005*, 2017.
- [67] R. Paccagnella, P. Datta, W. U. Hassan, A. Bates, C. Fletcher, A. Miller, and D. Tian. Custos: Practical tamper-evident auditing of operating systems using trusted execution. In *Network and distributed system security symposium*, 2020.
- [68] L. Qu, Y. Zhou, P. P. Liang, Y. Xia, F. Wang, E. Adeli, L. Fei-Fei, and D. Rubin. Rethinking architecture design for tackling data heterogeneity in federated learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022.
- [69] M. U. Rehman, H. Ahmadi, and W. U. Hassan. Flash: A comprehensive approach to intrusion detection via provenance graph representation learning. In *IEEE Symposium on Security and Privacy (S&P)*, 2024.
- [70] A. Riddle, K. Westfall, and A. Bates. Atlasv2: Atlas attack engagements, version 2, 2023.
- [71] A. Roy Chowdhury, C. Wang, X. He, A. Machanavajhala, and S. Jha. Crypte: Crypto-assisted differential privacy on untrusted servers. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*, 2020.
- [72] S. Sav, A. Pyrgelis, J. R. Troncoso-Pastoriza, D. Froelicher, J.-P. Bossuat, J. S. Sousa, and J.-P. Hubaux. Poseidon: Privacy-preserving federated neural network learning. *arXiv preprint arXiv:2009.00349*, 2020.
- [73] Y. Shen and G. Stringhini. Attack2vec: Leveraging temporal word embeddings to understand the evolution of cyberattacks. In *USENIX Security Symposium*, 2019.
- [74] M. Shoaib, A. Suh, and W. U. Hassan. Principled and automated approach for investigating ar/vr attacks. In *USENIX Security Symposium*, 2025.
- [75] A. J. Slagell, K. Lakkaraju, and K. Luo. Flaim: A multi-level anonymization framework for computer and network logs. In *LISA*, volume 6, 2006.
- [76] N. D. H. Son, H. T. Thi, P. T. Duy, and V.-H. Pham. Xfedgraph-hunter: An interpretable federated learning framework for hunting advanced persistent threat in provenance graph. In *International Conference on Information Security Practice and Experience*, pages 546–561. Springer, 2023.
- [77] B. Wang, Y. Chen, H. Jiang, and Z. Zhao. Ppefl: Privacy-preserving edge federated learning with local differential privacy. *IEEE Internet of Things Journal*, 2023.
- [78] Q. Wang, W. U. Hassan, D. Li, K. Jee, X. Yu, K. Zou, J. Rhee, Z. Chen, W. Cheng, C. A. Gunter, et al. You are what you do: Hunting stealthy malware via data provenance analysis. In *Network and Distributed System Security (NDSS)*, 2020.
- [79] S. Wang, Z. Wang, T. Zhou, H. Sun, X. Yin, D. Han, H. Zhang, X. Shi, and J. Yang. Threatrace: Detecting and tracing host-based threats in node level through provenance graph learning. *IEEE Transactions on Information Forensics and Security*, 17, 2022.
- [80] C. Wu, F. Wu, L. Lyu, T. Qi, Y. Huang, and X. Xie. A federated graph neural network framework for privacy-preserving personalization. *Nature Communications*, 13(1), 2022.
- [81] B. Xia, J. Yin, J. Xu, and Y. Li. Loggan: A sequence-based generative adversarial network for anomaly detection based on system logs. In *Science of Cyber Security: Second International Conference, SciSec 2019, Nanjing, China, August 9–11, 2019, Revised Selected Papers 2*. Springer, 2019.
- [82] J. Xu, C. Li, Y. Zheng, and Z. Li. Entente: Cross-silo intrusion detection on network log graphs with federated learning. *arXiv preprint arXiv:2503.14284*, 2025.
- [83] F. Yang, J. Xu, C. Xiong, Z. Li, and K. Zhang. Prographer: An anomaly detection system based on provenance graph embedding. 2023.
- [84] Y. Yang, B. Hui, H. Yuan, N. Gong, and Y. Cao. {PrivateFL}: Accurate, differentially private federated learning via personalized data transformation. In *32nd USENIX Security Symposium (USENIX Security 23)*, 2023.
- [85] D. Ye, R. Yu, M. Pan, and Z. Han. Federated learning in vehicular edge computing: A selective model aggregation approach. *IEEE Access*, 8, 2020.
- [86] X. Yi, R. Paulet, E. Bertino, X. Yi, R. Paulet, and E. Bertino. *Homomorphic encryption*. Springer, 2014.
- [87] J. Zeng, X. Wang, J. Liu, Y. Chen, Z. Liang, T.-S. Chua, and Z. L. Chua. Shadewatcher: Recommendation-guided cyber threat analysis using system audit records. In *IEEE Symposium on Security and Privacy (SP)*, 2022.
- [88] R. Zhao, M. Shoaib, V. T. Hoang, and W. U. Hassan. Rethinking tamper-evident logging: A high-performance, co-designed auditing system. In *ACM Conference on Computer and Communications Security (CCS)*, 2025.
- [89] Y. Zhao, M. Li, L. Lai, N. Suda, D. Civin, and V. Chandra. Federated learning with non-iid data. *arXiv preprint arXiv:1806.00582*, 2018.
- [90] Y. Zhao, H. Zhou, and Z. Wan. Superfl: Privacy-preserving federated learning with efficiency and robustness. *Cryptology ePrint Archive*, 2024.
- [91] Z. Zhong, W. Bao, J. Wang, X. Zhu, and X. Zhang. Flee: A hierarchical federated learning framework for distributed deep neural network over cloud, edge, and end device. *ACM Transactions on Intelligent Systems*

and Technology (TIST), 13(5), 2022.

- [92] T. Zhou, J. Zhang, and D. H. Tsang. Fedfa: Federated learning with feature anchors to align features and classifiers for heterogeneous data. *IEEE Transactions on Mobile Computing*, 2023.
- [93] M. F. Zolkipli and A. Jantan. An approach for malware behavior identification and classification. In *2011 3rd international conference on computer research and development*, volume 1. IEEE, 2011.

APPENDIX A DETAILED PRIVACY-PRESERVING WORD2VEC HARMONIZATION ALGORITHM

Algorithm 2 presents our privacy-preserving Word2Vec harmonization pipeline across N clients. Each client encrypts token identifiers with a shared symmetric key and sends encrypted token–vector pairs to a utility server, which merges overlaps by averaging vectors under encrypted identifiers. The resulting harmonized model $M_{w2v\text{-harm}}$ is then distributed back to all clients.

APPENDIX B RQ5: ABLATION STUDY

In this ablation study, we analyze the impact of key parameters within MIRAGE.

Algorithm 2: Privacy-preserving model harmonization.

```

Inputs : Client Word2Vec models
           {Word2vec1, Word2vec2, ..., Word2vecN}.
Output: Harmonized Word2Vec model  $M_{w2v\text{-harm}}$  sent to clients.
/* Central server distributes symmetric encryption key
   to each client. */
1 foreach client  $C_i$  do
2   | Send key to  $C_i$ 
3 end
/* Clients encrypt token identifiers. */
4 foreach client model Word2vec $i$  do
5   | foreach token  $t$  in Word2vec $i$  do
6     |  $\tilde{t} \leftarrow \text{Enc}(t, \text{key})$  Send  $(\tilde{t}, \text{Word2vec}_i[t])$  to Utility Server
7   | end
8 end
/* Utility server merges models using encrypted
   identifiers. */
9  $\text{TokenVectors} \leftarrow \text{InitializeEmptyDictionary}()$ 
10  $\text{TokenCounts} \leftarrow \text{InitializeEmptyDictionary}()$ 
11 foreach received pair  $(\tilde{t}, \text{Vector})$  do
12   | if  $\text{TokenVectors.HasKey}(\tilde{t})$  then
13     |  $\text{TokenVectors}[\tilde{t}] \leftarrow \text{TokenVectors}[\tilde{t}] + \text{Vector}$ 
14     |  $\text{TokenCounts}[\tilde{t}] \leftarrow \text{TokenCounts}[\tilde{t}] + 1$ 
15   | else
16     |  $\text{TokenVectors}[\tilde{t}] \leftarrow \text{Vector}$ 
17     |  $\text{TokenCounts}[\tilde{t}] \leftarrow 1$ 
18   | end
19 end
/* Average the vectors for overlapping tokens. */
20 foreach token  $\tilde{t}$  in  $\text{TokenVectors.Keys}()$  do
21   |  $\text{TokenVectors}[\tilde{t}] \leftarrow \text{TokenVectors}[\tilde{t}] / \text{TokenCounts}[\tilde{t}]$ 
22 end
23  $M_{w2v\text{-harm}} \leftarrow \text{NewModel}(\text{TokenVectors}, \text{TokenVectors.Keys}())$ 
24 foreach client  $C_i$  do
25   | Send  $M_{w2v\text{-harm}}$  to  $C_i$ 
26 end
27 return Harmonized model  $M_{w2v\text{-harm}}$  has been dispatched to all clients.

```

A. Effect of federated averaging rounds

We employed the DARPA E3 dataset to examine the impact of federated averaging rounds on detection performance. Our

methodology involved training the model over a range of federated averaging rounds and subsequently evaluating the model’s detection capabilities. The outcomes are depicted in Figure 5, which shows that detection performance improves up to a certain number of rounds before declining due to overfitting. Notably, this inflection point is also characterized by a minimal decrease in training loss, suggesting that the model has reached its learning capacity. This observation proves to be a valuable metric for determining the optimal moment to stop training, thereby preventing overfitting and ensuring optimal model performance.

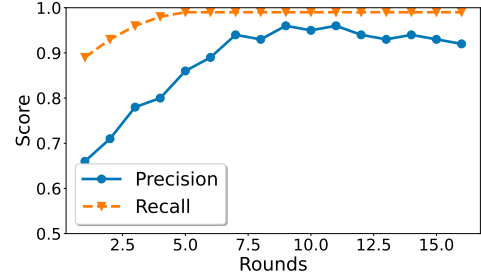


Fig. 5: Federated averaging rounds vs detection performance using E3 dataset.

B. Effectiveness of Word2vec harmonization

We evaluated the effectiveness of our Word2vec vector harmonization scheme through two experiments using the OpTC, E3 and E5 datasets. In the first experiment, each client utilized its locally trained Word2vec model to encode semantic features during the training process. In the second experiment, we harmonized the individually trained models into a central Word2vec model using the utility server architecture, as explained in Section IV. Then each client used this centralized model for generating semantic features. Table IV presents the average results for the DARPA E3, E5 and OpTC datasets. By employing the harmonized models, we achieved significantly better detection outcomes. This is because these local models encode different information for identical tokens across hosts. Such variability leads to heterogeneity in the feature space for the GNN model, impairing the model’s ability to generalize and converge effectively during the federated learning process, thereby yielding suboptimal results. However, through our novel, privacy-preserving aggregation of these semantic models, we have addressed this issue by merging the information in these models together to give the GNN more generalized input vectors.

C. Effectiveness of categorized graph learning

To examine the effectiveness of our process categorization based GNN ensemble technique, we conducted experiments comparing our ensemble technique against a single model approach in a federated setting. Specifically, for the ensemble method, we designated the number of categories to be 10. This approach standardizes all processes across various hosts into 10 distinct categories, ensuring that the GNN models learn

similar distributions regardless of the host. Such standardization is crucial for mitigating the impact of heterogeneous distributions and data imbalance during the federated averaging process. Table V reports the average results for the DARPA E3, E5 and OpTC datasets. The results indicate that utilizing categorized ensemble models yields superior performance. This improvement is attributed to each sub-model’s enhanced ability to concentrate on different patterns of system activity from different clients, thereby reducing the likelihood of these patterns becoming conflated during the federated averaging process.

TABLE V: Effectiveness of categorization-based ensemble learning.

Dataset	Type	Prec.	Rec.	F-Score
OpTC	Single	0.86	0.91	0.88
	Ensemble	0.90	0.92	0.91
E3	Single	0.89	0.99	0.94
	Ensemble	0.96	0.99	0.97
E5	Single	0.92	0.96	0.94
	Ensemble	0.98	0.98	0.98

TABLE IV: Effectiveness of Word2vec harmonization.

Dataset	Type	Prec.	Rec.	F-Score
OpTC	Local	0.58	0.97	0.73
	Harmonized	0.90	0.92	0.91
E3	Local	0.83	0.96	0.89
	Harmonized	0.96	0.99	0.97
E5	Local	0.81	0.94	0.87
	Harmonized	0.98	0.98	0.98