# KNOWML: Improving Generalization of ML-NIDS with Attack Knowledge Graphs

Xin Fan Guo[†,‡§], Albert Meroño Penuela[†], Sergio Maffeis[‡], Fabio Pierazzi[§]

[†]*King's College London,* [‡]*Imperial College London,* [§]*University College London*

*Abstract*—Despite extensive research on Machine Learning–based Network Intrusion Detection Systems (ML-NIDS), their capability to detect diverse attack variants remains uncertain. Prior studies have largely relied on homogeneous datasets, which artificially inflate performance score and offer a false sense of security. Designing systems that can effectively detect a wide range of attack variants remains a significant challenge. The progress of ML-NIDS continues to depend heavily on human expertise, which can embed subjective judgments of system designers into the model, potentially hindering its ability to generalize across diverse attack types.

To address this gap, we propose KNOWML, a framework for knowledge-guided machine learning that integrates attack knowledge into ML-NIDS. KNOWML systematically explores the threat landscape by leveraging Large Language Models (LLMs) to perform automated analysis of attack implementations. It constructs a unified Knowledge Graph (KG) of attack strategies, on which it applies symbolic reasoning to generate KG-Augmented Input, embedding domain knowledge directly into the design process of ML-NIDS.

We evaluate KNOWML on 28 realistic attack variants, of which 10 are newly-collected for this study. Our findings reveal that baseline ML-NIDS models fail to detect several variants entirely, achieving F1-Scores as low as 0%. In contrast, our knowledge-guided approach achieves up to 99% F1-Score while maintaining a False Positive Rate below 0.1%.

## 1. Introduction

Machine Learning–based Network Intrusion Detection Systems (ML-NIDS) have been extensively studied for their ability to detect both evolving and previously unseen attacks. Recent work on unsupervised ML-NIDS has demonstrated promise in detecting wide range of attack variants [1], [2]. Despite this claimed potential, the true generalization capability of contemporary ML-NIDS to detect diverse attack variants remains uncertain [3], [4]. As noted by Sommer et al. [5]: "*strength of machine-learning tools is finding activity that is similar to something previously seen*". This observation raises concerns about the generalization ability of current systems. Most recent studies rely on widely-used benchmark datasets to demonstrate effectiveness, yet these datasets often lack diversity and fail to represent realistic attack scenarios [3], [6]. For example, Flood et al. [3] show that datasets such as ISCXIDS2012 [7], CIC-IDS2017 [8],

and CSE-CIC-IDS2018 [8] are generated using automated tools with fixed configurations. As a result, they exhibit low variability in attack structure, creating inflated performance scores and offering a false sense of the model's generalization capabilities. This raises a critical question: *Can current ML-NIDS systems genuinely generalize to detect a broad spectrum of attack variants in real-world environments?*

We posit that ML-NIDS inherently rely on assumptions introduced by system designers. These assumptions are embedded either in the manually-engineered features or in the model architecture itself. As a result, detection becomes constrained to patterns that reflect a limited understanding of attack behaviors. For instance, Kitsune [1] detects only anomalies based on specific traffic features such as jitter and packet size, and will fail on threats that do not rely on those particular characteristics. A straightforward response might be to monitor a wide range of traffic statistics to cover diverse attack types, but this may lead to learning spurious correlations.

To address this gap, we propose KNOWML, a novel framework that systematically integrates structured attack knowledge into the ML-NIDS pipeline. KNOWML leverages knowledge graphs (KGs) and symbolic reasoning to improve the model's understanding of diverse attack strategies. We use Large Language models (LLMs) to automatically extract the implementation logic of public attack implementations. These are unified in a KG where each node represents a specific attack strategy or technique, offering a structured and interpretable view of the threat landscape. Building on this representation, KNOWML applies symbolic reasoning over the graph to generate KG-Augmented Input, a knowledge-enriched input representation that embeds attack semantics into both the training and evaluation phases of the ML-NIDS.

Unlike prior work that treats KGs as standalone tools for explanation or post-hoc analysis [9], [10], KNOWML embeds KG-derived knowledge directly into the core learning and evaluation process via KG-Augmented Input. Our experiments show that ML-NIDS built on "limited" or narrowly defined assumptions fail to generalize, even to attack variants that nominally lie within their monitored feature space. Moreover, while contemporary research largely focuses on model architectures or dataset constraints [3], [11], our work adopts a knowledge-centric perspective. We examine how the absence of rich, structured attack knowledge fundamentally limits the generalization capabilities of current ML-

NIDS systems.

In summary, we make the following contributions:

- We propose KNOWML, a novel framework that improves the generalization of ML-NIDS by integrating structured knowledge of realistic attack strategies (§3). This framework is grounded in practical considerations and recent trends in attack developments (§2).
- We identify generalization failures in current ML-NIDS approaches and systematise them into three categories in §4.2. Our findings reveal that even well-studied attack families (TCP DoS, HTTP DoS, and SSH Brute Force) can evade detection when models are built on narrow or incomplete representations of the threat landscape. To investigate these limitations, we conduct both conceptual and empirical analyses across 28 attack variants, including 10 newly collected for this study.
- We show that KNOWML significantly improves model generalization across diverse attack variants. We evaluate its performance using three ML models (GMMs [12], Kitsune [1], and DA [13]), combined with three types of knowledge-based inputs, across four datasets. Our results demonstrate consistent improvements in detection performance, emphasizing the critical role of structured attack knowledge in building more and effective intrusion detection systems §5.2.

## 2. Our Motivations

This section outlines the motivation behind analyzing attack implementations to improve the generalization capability of ML-NIDS systems. Our approach stems from a key observation:

> **Motivation**
>
> Attack implementations reveal the tactical decisions made by adversaries and how specific configurations define the execution of an attack. These configurations form a tamper space—a set of possible manipulations—where each variant reflects a distinct attack strategy. By analysing a wide spectrum of real-world implementations, we can systematically enumerate these tamper spaces and approximate the range of tactics that attackers are likely to employ.

This idea of building a system on a spectrum of known attack variants to generalise on attack mutation comes from practical considerations on threat models and attack trends.

**Existing Threat Models are Narrow.** Current ML-NIDS research often assumes limited attacker behavior, leading to narrow threat models [14] and poor generalization. Real-world attackers adapt their strategies dynamically, drawing inspiration from existing attack methods. By systematically analyzing and enumerating existing attack strategies and

their corresponding tamper spaces, we can gain insights into how attacks might be adopted or re-used.

**Attack Trends Suggest Attack Re-Use.** An analysis of the NIST vulnerability database [15] highlights two trends. First, vulnerabilities exhibit *cyclical patterns*, past exploits resurface in new systems, particularly those serving similar functions. For example, *Nkiller2* (explained in §4.2.1) exemplifies this pattern: it first appeared in 2008 [16], recurred in 2009 [17], and re-emerged again in 2024 [18]. System updates, patches, and modifications may reintroduce older vulnerabilities. Second, attackers combine established techniques to discover new vulnerabilities. Recent HTTP DoS attacks demonstrate this trend by exploiting malformed header values [19] a technique documented as recently as May 2025 [20], leveraging combined strategies previously documented in earlier attacks, such as HTTP request smuggling [21].

These observations motivate our approach. By studying existing implementations, we equip detection systems with the context needed to recognise evolved or recombined attacks that share foundational traits with known exploits. In §5.2, we demonstrate how this knowledge integration enables detection of various attack categories, including more sophisticated non-throughput-based attacks.

## 3. KNOWML: Knowledge-Guided ML

In this section we present key components of KNOWML. Including, 1) KG Construction, which builds a unified graph of attack strategies, 2) Symbolic Reasoning, which extracts structural patterns from the graph, and 3) KG-Augmented Input, which incorporates extracted knowledge into the ML-NIDS detection process. We focus on unsupervised ML-NIDS, because of their potential for generalization beyond limited attacks in the training set [3], [22].

**Preliminaries.** A Knowledge Graph (KG) is a directed graph composed of subject-predicate-object $(s, p, o)$ triples, where each triple defines a relationship between entities.

Subject nodes describe entities or concepts, object nodes represent linked values, and predicates define the relationship between subjects and objects [23].

Formally, we define our generated KG as $G = (V, E)$. Nodes $V$ encompass four entity types: strategy nodes $S = \{s_1, s_2, ..., s_n\}$, cluster nodes $C = \{c_1, c_2, ..., c_p\}$, family nodes $F = \{f_1, f_2, ..., f_k\}$ and repository nodes $R = \{r_1, r_2, ..., r_m\}$. The nodes are connected through edges, represented as $E_{v_1 v_2} \subseteq v_1 \times v_2$ for bidirectional connections between node types, or as $v_1 E v_2$ for unidirectional connections between nodes of any of the four entity types.

### 3.1. Overview of KNOWML

Figure 1 presents an overview of the complete KNOWML pipeline, which comprises three main components:

- **KG Construction ❶-❸:** KNOWML constructs a unified KG by collecting open-source attack imple-
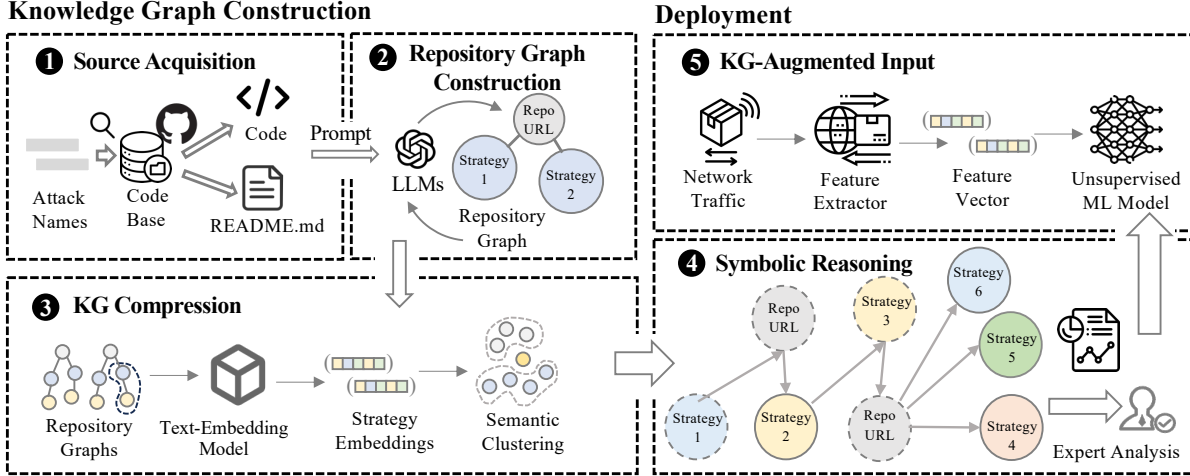
Figure 1: **Overview of the KNOWML pipeline.** The process begins with an input Attack Name and retrieves its corresponding implementation (Step 1). Next, a Repository Graph is constructed for the specific implementation (Step 2). These graphs are then compressed into a unified KG representing unique attack strategies (Step 3). This unified KG forms the basis for Symbolic Reasoning (Step 4). The output of Step 4 is used to generate KG-Augmented Input, which is then fed into the ML-NIDS for training and evaluation (Step 5).

mentations and extracting the corresponding strategies. This KG provides a structured representation of attack behaviors, enabling a comprehensive understanding of patterns across diverse implementations.

- **Symbolic Reasoning ❹:** We introduce three inference that are applied to produce KG-Augmented Input: R1) Frequency Rule: Identifies commonly occurring attack strategies across multiple implementations, R2) Transitive Rule: Analyzes the relationships among strategies to infer composite attack paths through transitive connections, and R3) Cross-Family Rule: Discovers invariant strategies that persist across different attack families.

- **Knowledge-Augmented Input ❺**: This module transforms the symbolic insights from R1–R3 into structured inputs suitable for a ML-NIDS. A custom feature extractor processes network traffic streams and provides the KG-derived input into the ML.

## 3.2. KG Construction

Attackers often rely on automated tools or scripts to launch cyberattacks. To analyze their behavior and capture variations in attack methods, we construct a KG based on the implementation details of these scripts. Specifically, we hypothesize that the parameters defined in the implementation documentation (typically found in README.md files) represent distinct attack strategies. Each parameter configures or triggers a specific attack mechanism and corresponds to a unique attack vector. These vectors collectively contribute to observable differences in attack behavior. For example, consider a script invoked with the following parameters:

```
$python script.py --keep-alive --syn_flag
```

Each flag in this command specifies a particular attack characteristic, and together, they define a composite strategy represented within the KG. Note that we intentionally avoid examining the specific values assigned to these parameters (such as --keep-alive=0 or --keep-alive=1), since this approach could lead to overfitting by learning overly-specific patterns. Instead, our goal is to learn the existence of these configurable parameters as abstract attack strategies, to then suggest which features the ML-NIDS should monitor. To effectively capture diverse attack implementations, we developed a pipeline that systematically search open source repositories based on specified attack names (see ❶ in Figure 1). This module identifies relevant repositories by matching the attack name within repository descriptions or titles. It then retrieves both the implementation code and the associated README.md files. Subsequently, we generate individual Repository Graphs for each retrieved implementation and compress these graphs into a unified KG, as described in detail below.

**3.2.1. Repository Graph Construction ❷.** In this step, we build Repository Graphs that store the attack strategies implemented for investigated repositories. We define an ontology (KG schema) based on best practices in ontology engineering [23]. Specifically, our ontology includes the following classes and properties: i) *Attack*: the attack to which the strategy belong e.g., TCP DoS, ii) *Strategy*: A specific technique employed by attackers (parameter described above), and (ii) *Description*: A textual explanation describing how the strategy configures or modifies the attack; and iv) *Repository Identifier*: a Unique Resource Identifier (URI) that links each strategy to its source repositor. The is defined as a functional property, meaning each extracted strategy is associated with exactly one repository. We utilize LLMs

to automate the extraction of attack strategies from repositories. This extraction task is formally known as Named Entity Recognition (NER). The use of LLMs for this task is practical and justified by recent research demonstrating their effectiveness in various NER tasks, particularly in security contexts [24], [25], [26], [27]. Moreover, LLMs enable large-scale automated code analysis, previously infeasible due to extensive manual effort. In this paper, we introduce a framework designed to automate this NER process (financial cost considerations are discussed in §7). Specifically, we implement a GraphRAG-based approach [28] using the GPT-4o-mini model (for evaluation details and rationale behind this choice, see Section §4.1.1). The subsequent section addresses scalability considerations, recall degradation associated with increased content window sizes, and hallucination issues that KNOWML specifically aims to mitigate.

**Ensuring Logical Consistency and Scalability.** To ensure scalability and enable tractable symbolic reasoning of the defined KG, we designed our ontology according to OWL Lite standards [29], which guarantees polynomial-time reasoning and prevents computational intractability when processed by Description Logic (DL) reasoners. Additionally, our ontology was validated using HermiT [30] to ensure logical consistency through DL reasoning, verifying: (i) no logical contradictions exist; (ii) all concepts are satisfiable; and (iii) all axioms are mutually compatible.

**Recall Degradation with Increasing Context Window.** A significant challenge in extracting attack strategies arises from the unpredictable length of `README.md` files and source code, which can vary considerably. Research shows that language models experience reduced Recall performance when processing extended contexts [31], [32]. To mitigate this issue, we implemented an iterative "gleaning" method designed to offset potential recall degradation as the context window size increases. Specifically, we enhanced recall in entity detection by restricting the LLM to binary (`YES/NO`) responses about potentially missed entities, employing logit bias techniques [28].

**Managing Hallucination and Noise.** LLMs are known to occasionally hallucinate, generating inaccurate or fabricated information [33]. To address this challenge, we adopted a few-shot learning approach by providing explicit examples representing common extraction scenarios. These include: (i) *Structured inputs*: where extraction is straightforward, such as:

```
$ python script.py [-p1] [-p2]
options:
   -p1: <p1 description>
   -p2: <p2 description>
```

(ii) *Unstructured inputs*: where attack strategies are embedded within natural text; (iii) *Empty inputs*: Where code or `README.md` files do not explicitly mention parameters (e.g., simply `$ python script.py`); and (iv) *Irrelevant inputs*: where search results include unrelated repositories. For instance, a query for HTTP DoS may return unrelated projects, such as an HTTP-based Redis pooler [34]. In addition, we enforced structured output formats, requiring the LLM to follow predefined templates (in this case the ontology) with clearly specified fields [35]. This combined approach effectively minimises hallucination occurrences and improves extraction accuracy(see evaluation in §4.1.1).

**3.2.2. KG Compression ❸.** To eliminate redundant strategies, we compress the extracted Repository Graphs by clustering similar strategies, e.g., "Randomize Ports" and "Random Port Targeting". We perform semantic clustering using embeddings derived from each strategy's Name and Description, generated via OpenAI's text-embedding model [36]. These embeddings serve as the basis for grouping similar strategy nodes. Among available clustering algorithms, we select Hierarchical Agglomerative Clustering (*HAC*) [37] due to its ability to maximize inter-cluster distance. This property ensures that semantically distinct strategies remain in separate clusters. HAC produces a higher number of unique clusters than other algorithms [38]. Specifically, we employ HAC with complete linkage to maximize inter-cluster distance. After clustering, we select a single representative strategy per cluster. We define the representative as the strategy whose embedding has the minimum total distance to all other strategies in the same cluster. This selection ensures that the chosen representative best reflects the core semantics of the group. The selection is formally defined as:

$$ S_j = \arg\min \left\{ \sum_{i:s_i \in C_j} \|e_i - e_\ell\| \;\middle|\; s_\ell \in C_j \right\} \quad (1) $$

where $e_i$ is the embedding of the strategy $s_i \in C_j$, and $S_j$ denotes the representative strategy in cluster $C_j$ that has the minimum cumulative distance to all other strategies within the same cluster.

### 3.3. Symbolic Reasoning ❹

In this section, we introduce three symbolic inference rules and their underlying intuition. These rules form the basis of the KG-Augmented Input, illustrated in Figure 3 and further detailed in §3.4.

**3.3.1. Frequency Rule (R1).** To identify common attack strategies, we quantify their prevalence based on implementation frequency i.e., how often each strategy occurs relative to others. For a given strategy, $s'$, we define the frequency rule as:

$$ \text{Freq}(s') = \frac{|\{r \in R \mid (s', r) \in E_{SR}\}|}{|R|} \quad (2) $$

where $R$ is the set of repository nodes, and $s'$ denotes a particular strategy being evaluated.

This formula measures the relative occurrence of strategy $s'$ based on the number of repositories where it has been implemented. Specifically, the numerator counts the number

of unique repositories associated with strategy $s'$, while the denominator represents the total number of repositories examined.

**3.3.2. Transitive Rule (R2).** Transitive rules infer potential attack paths that a malicious actor might exploit. They address the question: "Given strategy A, what additional strategies could an attacker simultaneously employ?" This rule leverages the transitive property of strategies to identify combinations of attack methods, which implies configuring and activating multiple parameters concurrently to execute an attack. Formally, the transitivity property $E$ is defined as follows:

$E$ is *transitive* if, for all $v_1, v_2, v_3$, $(v_1, v_2) \in E$ and $(v_2, v_3) \in E$ implies $(v_1, v_3) \in E$.

Hence for given strategy $s_x$, the transitive rule $TC(s_x)$ is defined as:

$$TC(s_x) = \{ s_y \in S \mid (s_x, s_y) \in E_{\text{trans}} \}. \quad (3)$$

where $E_{\text{trans}}$ represents the transitive relation between strategies. A strategy $s_y$ belongs to $TC(s_x)$ iff there exists a path from $s_x$ to $s_y$. This can occur in two ways:

- *Single-hop path (via shared repository)* When two strategies $s_1$ and $s_2$ both belong to the same repository $r$, they form a transitive relation. For single hop this is usually the thought the $E_{RS}$ and $E_{SR}$ relation (illustrated by $S_1$ and $S_2$ in Figure 3). Formally:

$$\exists r \in R : (s_1, r) \in E_{SR} \wedge (r, s_2) \in E_{RS}. \quad (4)$$

- *Multi-hop path*: When strategies are connected through a chain of intermediate strategies, they form a multi-hop transitive path. For instance, if strategy $s_i$ is semantically similar to $s_j$ and both were assigned to the same cluster during the KG compression process, we then examine the neighbors of $s_j$ that share the same parent repository node. This process extends the logic of the single-hop path described above (illustrated by $S_6$ and $S_7$ in Figure 3). Formally, it is defined as:

$$\exists s_i, s_j, \dots, s_k \in S : \\ (s_x, s_i) \in E_{\text{trans}} \wedge (s_i, s_j) \in E_{\text{trans}} \wedge \quad (5) \\ \dots \wedge (s_k, s_y) \in E_{\text{trans}}.$$

**3.3.3. Cross-Family Rule (R3).** To determine if specific attack strategies consistently appear across various attack families (family invariant), we propose the cross-family rule. A strategy is considered invariant if it appears consistently across all $N$ attack families (i.e., it occurs exactly $N$ times). Identifying such invariant strategies clarifies the scope and criticality of attack surfaces requiring investigation. Note that cross-family strategies form a subclass of common strategies shared across all attack families, as illustrated in Figure 2. The cross family rule is formally defined as:

$$\mathcal{S}_{\text{inv}} = \{ s \in S \mid \quad |\{ f \in F \mid (s, f) \in E_{SF} \}| > 1 \} \quad (6)$$



For a single attack $R_3 \subset R_1$    For all attacks $R_3 = \{A_1 \cap A_2 \cap A_3\}$
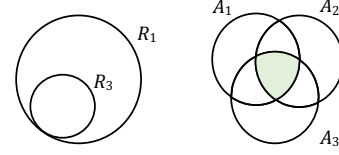
Figure 2: Illustration of Cross-Family Rule (R3), and its relationship to other rules.
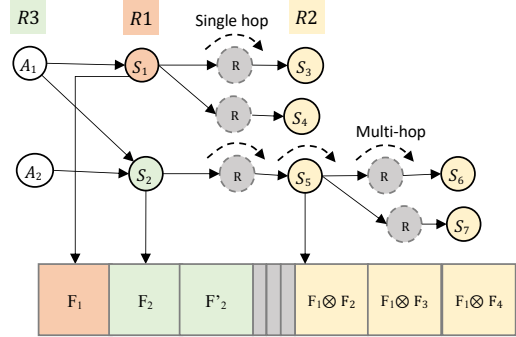


Figure 3: Application of Symbolic Reasoning Rules (§3.3) for KG-Augmented Input (§3.4).

where $S$ is the set of Strategy nodes, $F$ is the set of Family nodes, and $E_{SF}$ is the set of edges between elements of $S$ and $F$.

## 3.4. KG-Augmented Input ❺

We construct the KG-Augmented Input by processing raw network traffic captures (pcaps) into structured feature vectors. These vectors serve as inputs for training and evaluating ML-NIDS. This approach enables the model to identify regions in the input space that are susceptible to adversarial manipulation, improving its ability to generalize across attack variants that employ different strategies or their combinations. The KG-Augmented Input is derived directly from the symbolic reasoning framework introduced in §3.3. The transformation of symbolic rules into features follows three rule-based mappings: Direct Feature Mapping (R1), Composite Feature Generation (R2), and Cross-Endpoint Feature Aggregation (R3). Each rule maps symbolic knowledge to concrete feature representations (see Figure 3), as detailed below.

**3.4.1. Direct Feature Mapping (R1).** This rule enables a one-to-one mapping between symbolic strategies and features. Let $\mathcal{S} = s_1, s_2, \dots, s_n$ denote the set of strategies identified by R1. Each strategy $s_i$ is directly translated into a corresponding feature $f_i$, for example, a strategy involving setting `--packet-size` yields a features such as `packet_size`. Formally, the mapping is defined as:

$$\phi_1 : \mathcal{S} \to \mathcal{F}_{\phi_1}, \quad \phi_1(s_i) = f_i \quad (7)$$

**3.4.2. Composite Feature Generation (R2).** R2 identifies relationships between strategies, allowing the construction of composite features. For each pair $(s_i, s_j)$ in the transitive closure $TC(s) \subseteq \mathcal{S} \times \mathcal{S}$, we generate a feature that combines $f_i$ and $f_j$. This captures interactions between strategies that jointly characterize behavior. The mapping is defined as:

$$\phi_2 : TC(s) \rightarrow \mathcal{F}_{\phi_2}, \quad \phi_2((s_i, s_j)) = f_i \otimes f_j \qquad (8)$$

Here, $\otimes$ denotes a composite operation such as multiplication or concatenation.

**3.4.3. Cross-Endpoint Feature Aggregation (R3).** R3 addresses invariant strategies that should be monitored consistently across endpoints. For each invariant strategy $s_i \in \mathcal{S}_{\text{inv}}$, we aggregate its corresponding feature $f_i$ across the set of endpoints $\mathcal{E}$. This aggregation captures persistent behavioral patterns. The mapping is defined as:

$$\phi_3 : \mathcal{F}_{\phi_1} \times \mathcal{E} \rightarrow \mathcal{F}_{\phi_3}, \quad f_i' = \text{Agg}_{e \in \mathcal{E}}(f_{i,e}) \qquad (9)$$

Here, $f_{i,e}$ is the value of feature $f_i$ at endpoint $e$, and Agg denotes an aggregation function that groups features based on the endpoint. This cross-endpoint aggregation enables the model to recognize globally consistent attack behavior, particularly for cross-family strategies, which tend to recur across different attack contexts. Aggregating these features based on destination allows the detection system to remain sensitive to subtle, distributed signs of compromise.

Given the above transformations, the final input vector $\boldsymbol{v}$ for ML is then:

$$\boldsymbol{v} = \mathcal{F}_{\phi_1} \cup \mathcal{F}_{\phi_2} \cup \mathcal{F}_{\phi_3} \qquad (10)$$

Note that the decision of composite function is guided by human expertise. Defining how strategies should be combined into meaningful features often involves subjective decisions, and no universal rule exists that applies to all cases. In this work, we include only semantically-meaningful composites that align with the domain context. For example, it was observed from the constructed KG that TCP DoS implementations allowed users to enable the `--frag` option alongside `--tpc_dos`, allowing an attack to be executed with fragmented packets and used composite feature `tcp_fragmented_count`. We further discuss the implications of function mapping in §7.

**3.4.4. Feature Extraction.** We adopt Welford's online algorithm [39], as well as the approach introduced in [1], to compute and maintain descriptive statistics for each feature $f_i \in \mathcal{F}_{\phi_1} \cup \mathcal{F}_{\phi_2}$. For each such feature, we compute the following five statistics: weight ($W$), cumulative sum ($CS$), mean ($\mu$), sum of squared residuals ($SSR$), and standard deviation ($\sigma$). These statistics are incrementally updated using Algorithm 1 (see Appendix). The efficiency of this method has been established in [1], where it was shown to require only a small memory footprint and offer constant-time updates, even when processing high-volume network streams. For features $f_i \in \mathcal{F}_{\phi_3}$, we perform aggregation

across destinations. The full update process proceeds as follows: At time $t$, when a packet arrives from host $H_1$ to host $H_2$, we first update the per-conversation statistics $f(H_1, H_2)$ for all features in $\mathcal{F}_{\phi_1}$ and $\mathcal{F}_{\phi_2}$ using Algorithm 1. Next, we update the aggregated features $f_i' \in \mathcal{F}_{\phi_3}$ based on destination-specific updates. The final input vector $\boldsymbol{v}$ at time $t$ is given by the concatenation $\boldsymbol{v} = F \parallel F'$ where $F$ contains the raw and composite features, and $F'$ contains the aggregated features (the entire update process is detailed in Algorithm 2 of the Appendix).

## 4. Evaluation: KG Construction and Analysis

This section presents the empirical evaluation KG construction described in §3.2. Followed by key findings from the constructed KG §4.2.

**Prototype.** We developed a prototype implementation of KNOWML comprising 5,578 lines of code, covering KG Construction, Symbolic Reasoning, and the extraction of KG-Augmented Inputs from network packet streams. For packet capture and preprocessing, we used `tshark` (version 3.2.3). The KG-augmented input module operates as a plug-and-play component. It supports both offline analysis (processing pcap files) and online analysis (training and inference on live traffic), ensuring compatibility with a variety of ML-NIDS. We saved the resulting knowledge graph in the `.graphml` format to support symbolic inference. The complete code and dataset will be publicly released upon publication of the paper. In the meantime, they are available upon request by contacting the first author.[1]

### 4.1. Knowledge Graph Construction

We constructed KG from three families: TCP DoS, HTTP DoS, and SSH Brute Force (refer to §7 for reason and discussion). In total we crawled 7,853 GitHub repositories using keyword search based on match in either name of the repository or its description using GitHub Search API [40]. To ensure broad coverage, we used both base terms and common variations of attack names. For example, "TCP Flood" and "TCP DoS", to account for differences in naming conventions (see Appendix Table 8 for the full keyword list).

The resulting KG includes 5,410 nodes and 27,876 edges. From this graph, we extracted 214 features to generate KG-Augmented Inputs for downstream ML tasks (see Table 7 in the Appendix). We discuss scalability, transformation cost, and design choices in §7.

**4.1.1. Evaluating LLM for KG construction.** To evaluate the effectiveness of LLMs in KG construction and to select a model that minimizes hallucination (i.e., maximizes Recall), we assessed their performance on strategy Named Entity Recognition (NER) tasks, as described in §3.2. We compared model outputs to human analysis by first reviewing all 7,853 implementations and selecting 150 representative

---

1. Xin Fan Guo, xinfan.guo@kcl.ac.uk

TABLE 1: Performance Results for Strategy NER task

| Model | Precision | Recall | F1-Score |
|---|---|---|---|
| GPT-4o mini | 0.4622 | **0.9058** | 0.6121 |
| GPT-3.5 Turbo | 0.4583 | 0.6413 | 0.5346 |
| Gollie 13B | 0.4231 | 0.0493 | 0.0884 |

repositories across four scenarios likely to occur during extraction: (1) structured inputs, (2) unstructured inputs, (3) empty inputs—used to test for hallucinations, a known LLM issue [33], and (4) irrelevant content. Our evaluation included both general-purpose and task-specific LLMs: GPT-3.5 and GPT-4o, which have demonstrated strong performance in security-related NER tasks [41], and Gollie [24], a local, domain-specific LLM optimized for NER. Table 1 reports Precision, Recall, and F1-Score for each model on the selected repositories. These metrics are standard in evaluating NER systems [24], [42], [43]. In our evaluation, Recall measures the proportion of LLM-extracted strategies that match human annotations, while Precision reflects how many additional strategies were identified. GPT-4o-mini achieved the highest Recall at 90.58%, indicating strong alignment with human-labeled entities. Although Precision was comparatively lower, this is not a concern. Most discrepancies result from variations in strategy phrasing (e.g., "Randomize Ports" vs. "Random Port Targeting"), where the underlying semantics remain intact. These can be addressed through the semantic clustering method outlined in §3.2.2. Based on the performance in Table 1, we selected GPT-4o-mini for KG construction. While LLMs are inherently non-deterministic [44], our results should stay reproducible when using the same GPT-4o-mini version (2024-07-18) with controlled temperature and seed settings (in our evaluation we always set it to 0) to eliminate output randomness.

## 4.2. Reflections on Limitations of ML-NIDS

Here, we present a discussion informed by insights derived from the constructed KG. We highlight key limitations in contemporary approaches, focusing on how restricted or incomplete knowledge can lead system designers to make flawed assumptions. §5 complements this discussion with empirical evidence that supports the concepts introduced here.

**4.2.1. Out-of-Dimension Attacks.** Our analysis of the constructed KG reveals attack strategies that are frequently overlooked by contemporary ML-NIDS. Existing approaches often depend on *limited a priori* assumptions, grounded in the designer's domain expertise and heuristics about which attacks are likely to occur and which features are most indicative of malicious behavior. These assumptions narrow the feature space and implicitly assume that future attacks will conform to predefined patterns.

This leads to a vulnerability we term Attack Strategy Bias (demonstrated in §5.2), where the ML-NIDS overfits to known strategies and fails to generalize to unconventional
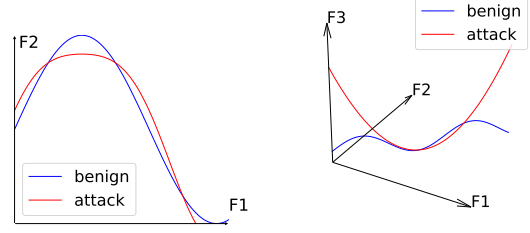


Figure 4: Illustration of Out-of-Dimension Attacks: (Left) Under constrained feature assumptions, benign and attack traffic overlap, leading to poor separability. (Right) Expanding the feature space with attack-aware knowledge enables better separation and detection of previously overlooked strategies.

ones. This limitation becomes evident when attackers exploit attack vectors not considered during model design. For example, Nkiller2 abuses TCP's flow-control mechanism by setting the window size to zero, halting data transmission. Since most systems do not monitor window size statistics, this attack proceeds undetected. Similarly, HTTP Header Overflow manipulates HTTP headers using excessive or malformed values. By operating in unmonitored portions of the input space, this technique bypasses detection models that focus on traffic volume or sized-based patterns. As shown in Figure 4, such strategies fall outside the expected feature distribution, underscoring the risk of static, narrow design assumptions in modern ML-NIDS.

**4.2.2. M-N End-Point Attacks.** ML-based NIDS often operate under *restrictive* assumptions about how attacks are launched, particularly regarding the configuration of source endpoint. These assumptions limit detection by focusing on isolated attacker-victim interactions and overlooking broader communication patterns. Our analysis of the KG reveals a critical gap: attackers can exploit endpoint diversity, such as IP spoofing or multi-source coordination to launch attacks that bypass models trained only on localized behavior.

For example, DoS attacks typically originate from a single source (1-to-1), whereas Distributed DoS (DDoS) attacks emerge from multiple sources targeting the same victim (M-to-1). Although current security literature distinguishes between these forms, most ML-NIDS architectures implicitly assume attacks are single-source and analyze each communication channel independently. This per-connection focus captures only local statistics and misses aggregate patterns that emerge at the destination level.

This design bias creates a critical vulnerability. When attackers distribute activity across many sources, each stream may appear benign in isolation. However, the cumulative effect is visible only through destination-based aggregation that reveals the full scope of the attack. Figure 5 illustrates how distributed threats can evade detection when models lack global traffic context. As with Out-of-Dimension Attacks, this reflects a deeper problem that is independent of ML. A reliance on static operational assumptions that

do not capture the diversity and adaptability of real-world adversarial behavior.
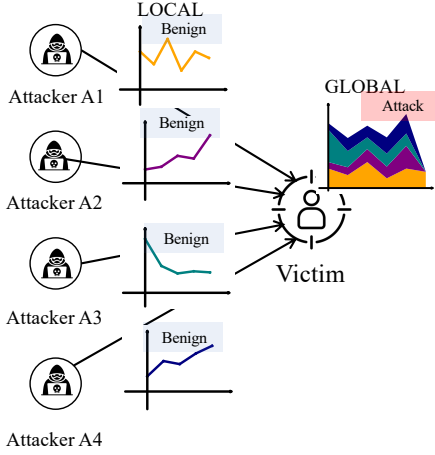


Figure 5: Distributed DoS Attack Detection: Single Conversation View (Local) Fails to Reveal Aggregate Malicious Behavior Visible in Destination-Based View (Global)

**4.2.3. Non-Throughput-Based Attacks.** An analysis of the KG revealed that attackers employ strategies that do not manifest in network throughput. Many existing methods for detecting DoS or Brute-force attacks employ *merely* rate or throughput based statistics [1], [6], [45], reflecting a limited understanding of the broader threat landscape. For example, attackers may use throttling mechanisms in TCP DoS attacks that synchronize with the Retransmission Timeout (RTO) clock, generating brief traffic bursts that even out over time (Low-rate TCP [46]). Similarly, attackers may use scripts that allow to set randoming payload padding to match the profile of benign traffic [47] (illustrated in Figure 6). In SSH Brute Force attacks, attackers may establish persistent SSH connections and attempt multiple connections within a single channel [48], bypassing detection systems that focus on counting distinct connections. These attack strategies are imperceptible to ML-NIDS models built on the assumption that DoS and Brute-force attacks manifest primarily in throughput statistics.

The following section provides empirical evidence showing that ML-NIDS, which encode different assumptions about attack types, perform significantly differently based on the knowledge they are built with.

## 5. Evaluation: KNOWML Detection Result

This section provides an empirical evaluation of the concepts introduced earlier and demonstrates the effectiveness of the KG-Augmented Input in detecting various categories of attacks. We begin by outlining the experimental setup, followed by a presentation of results organized by attack category.



Figure 6: Detection Bypass via Payload Padding: Attack traffic match Benign

### 5.1. Experimental Setup

**5.1.1. Datasets.** We evaluated KNOWML using four distinct datasets, each selected to address specific evaluation objectives.

**Network Environment Evaluation.** To assess detection performance across different network settings, we include one dataset representing enterprise networks and another focused on IoT environments. **CICIoT2023 (IoT-23)**: A dataset containing network traffic from 105 IoT devices, representative of IoT environments [49]. **CSE-CIC2017 (IDS-17)**: A dataset simulating an enterprise network environment, including diverse protocol coverage, 3 servers, and 7 hosts [8].

**Simulated Attack Variants (SIM).** This dataset simulates attack scenarios that a ML-NIDS might encounter in a given network environment (either IoT-23 or IDS-17). The selected attack strategies come from the KG, corresponding to well-documented attacks by NIST [15]. The simulation was conducted using the scapy library, a widely-used tool in network security research [50], [51], [52]. Implementation details are available in Appendix A.

**Captured Realistic Attack Variants (CAP).** This dataset was created by capturing real attack traffic to evaluate KNOWML and establish baseline performance against a variety of realistic attack variants. This cross-evaluation approach—where models are trained on benign data from one network environment (either IoT-23 or IDS-17) and evaluated on another—follows established methodologies from prior research [53], [54], [55]. Further details on the captured dataset are provided in Appendix B.

**5.1.2. Baselines.** We compared against two types of baselines: 1) Attack knowledge (we will refer to as Knowledge-Model, representing different prior knowledge bases) and 2) ML model (representing the model architecture). We consider each combination of KG-Model as a distinct baseline. This approach allows us to evaluate two critical aspects: first, the importance of Knowledge-Augmented Input for detection performance; second, whether complex ML models that extract meaningful representations from data can compensate for limited knowledge-based input.

**Knowledge-Models.** We evaluated knowledge models representing three distinct paradigms: 1) *Limited Knowledge*: **Kitsune (KIT)** [1] - KIT utilizes 100 features based on

TABLE 2: **Attack Variants used for evaluation.**

| Category | Variant Name | Dataset |
|---|---|---|
| Non-Throughput | Slowloris DoS* | IDS-17 |
| | Slowloris DDoS* | IoT-23 |
| | Slowhttptest DoS* | IDS-17 |
| | Fiberfox AVB | CAP |
| | Fiberfox BYPASS | CAP |
| | TCP Low-rate | SIM |
| | SSH Brute Force (P=0) | CAP |
| | SSH Brute Force (P=1) | CAP |
| M-N End-Point | HTTP DoS (1-1) | IoT-23 |
| | HTTP DDoS* (M-1) | IoT-23 |
| | Golden Eye DoS* | IDS-17 |
| | Hulk DoS* | IDS-17 |
| | Fiberfox GET (DoS) | CAP |
| | Fiberfox SLOW (DoS) | CAP |
| | Fiberfox STRESS (DoS) | CAP |
| | SYN DoS | IoT-23 |
| | SYN DDoS* | IoT-23 |
| | TCP DoS* | IoT-23 |
| | TCP DDoS | IoT-23 |
| | PSHACK DDoS* | IoT-23 |
| | RSTFIN DDoS* | IoT-23 |
| | SynonymousIP DDoS* | IoT-23 |
| | ACK Fragmentation DDoS* | IoT-23 |
| Out-of-Dimension | HTTP Overflow | SIM |
| | HTTP Mal Header | SIM |
| | Nkiller2 | IoT-23 |
| | SSH Brute Force (SSH BF)* | IoT-23 |
| | SSH Brute Force* | IDS-17 |

domain expertise to detect both known and unknown attacks. For example, the authors hypothesize that increased jitter in IP surveillance camera streams could indicate man-in-the-middle attacks. We selected KIT to examine how practitioners design features based on partial attack knowledge, lacking a comprehensive understanding of all potential attack variants. The corresponding ML-model, KIT-ML, serves as a baseline for evaluating model architectures. 2) *Standardized Approach*: **Standardized Feature Set (SFS)** [56] - SFS proposes a feature set derived from the CISCO NetFlow standards [57], aiming to facilitate cross-model comparisons. We selected SFS to investigate whether features developed for general network monitoring, rather than explicit attack detection, can still be effective for identifying attacks. 3)*Random Selection*: **CICFlowMeter (CIC)** [8] - Despite over 2,000 citations in the ML-NIDS community, the 87 features proposed by CIC lack a documented rationale for their design. These features, ranging from basic network statistics (e.g., flow byte rate) to protocol-specific counters (e.g., `ece_cnt` for ECE flag packets), are not explained in terms of their relevance to attack detection. We selected CIC to establish a performance baseline when features are chosen without a clear understanding of attacks or detection objectives.

**ML models.** We selected three ML model architectures that represent distinct approaches in ML-NIDS research, including: i) Gaussian Mixture Models (**GMM**) [12], [58], [59] ii) Denoising Autoencoder (**DA**) [13], and iii) Kitsune (**KIT-ML**) [1] extends the autoencoder approach by grouping correlated features into subgroups before training. Each subgroup uses a separate autoencoder, increasing model complexity. We selected KIT-ML (an ensemble of DA) to test whether sophisticated architectures that mine meaningful representations from feature correlations can substitute for KNOWML.

Note that we also attempted to include Whisper [6] (CCS'21), which claims effectiveness against evasion attacks. However, the implementation achieved nearly 0% Recall even at 10% FPR. Despite contacting the authors, we could not reproduce the published results. We therefore excluded Whisper from our evaluation while noting this does not reflect on the original work's merit.

**5.1.3. Metrics.** We used $F_1$-Score, Precision and Recall as our primary evaluation metrics to its widespread adaptation in literature [6], [45], [60] and unbiased nature [22].

**5.1.4. Hyperparameter Selection.** We conduct an extensive grid-search to identify optimal parameters for each model to enable fair comparison and to utilize full capabilities of baseline approaches [11], [22]. To avoid data snooping [11], we use a time-aware dataset split of 80% for training (only on benign data), 10% for validation, and 10% for testing. Note that for the grid search deliberately excluded some attacks to find an optimal model. This approach simulates realistic deployment: models optimize performance on known attacks (marked with * in Table 2) but must detect unseen variants from the same attack families during evaluation. Attack used for training are used for training are *distinctly* different from attacks use for evaluation (e.g., low rate attack excluded). This methodology evaluates the model's ability to generalize beyond its training distribution—an essential requirement for deploying ML-NIDS in real-world environments. Our grid search evaluated 94 parameter combinations for each of the 4 KG-Model baselines (including KNOWML). We tested three FPR thresholds (0%, 0.01%, and 0.1%) to identify optimal anomaly detection boundaries. This yielded 3,3384 total experiments (3 ML models × 94 parameters × 3 thresholds x 4 Knowledge-Models). Table 6 in the Appendix details all parameter ranges. We trained baseline methods exclusively on benign traffic and selected models that achieved the lowest False Positive Rate (FPR) while maintaining the highest True Positive Rate (TPR).

## 5.2. Detection Results

This section presents results that assesses utility of KG-Augmented Input for detection compare to baseline approaches. The results are presented and discussed by category.
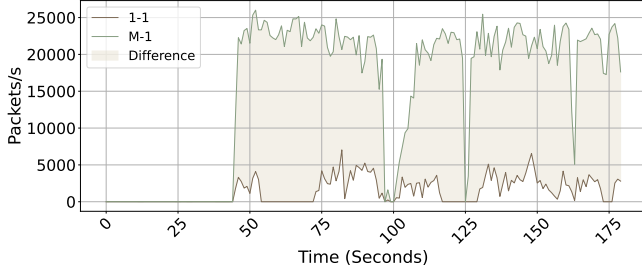
Figure 7: Wireshark I/O Graph showing packets/s with RST-FLAG set on IoT-23 dataset. The green line represents the distribution from single attacker to a victim host (1-1), while the green line indicates distribution by monitoring all attacker to victim host (M-1).

**5.2.1. M-N End-Point Attacks.** Table 3 presents performance against M-N Endpoint Attacks. Different KG-Model baselines demonstrate markedly different performance results. The evaluated attack variants possess significant increases in network throughput. The results in Table 3 break the assumption that such attacks should be easy to detect given their homogeneous nature and significant rise in network throughput. Several Knowledge-Model baselines completely fail to detect certain variants, showing 0% $F1$-score in some cases. For example, SFS demonstrates 0% $F1$-score on DDoS-RSTFIN attack across all ML-model types. This occurs because Knowledge-Model baselines are built on different *prior* assumptions about how attacks might be launched (as discussed in §4.2), restricting ML model's learning context. Investigating CIC and SFS reveals that they provide single conversation statistics, whereas KIT monitors aggregation of source-based statistics. All assume that attacks are launched from a *single* host. This limitation becomes especially evident when observing CIC and SFS results on distributed attacks. CIC demonstrates 75.09% $F1$-score on TCP DoS (GMM) but 0% on TCP DDoS across all algorithms. This dramatic performance degradation exemplifies **Attack Strategy Bias**, when Knowledge-Models embed specific attack strategy assumptions, they overfit to those patterns and fail against alternative strategies. In addition this also demonstrates that ML-model type and complexity play no major role. The prior assumptions built into Knowledge-Models fundamentally limit the perspective and context available to ML-NIDS, thereby restricting their detection capabilities against diverse M-N-Endpoint attacks. Figure 7 illustrates this limitation: the data distribution clearly shows how ML-NIDS misses obvious anomalies (green line) when the underlying model assumes single-host attacks (brown line). KNOWML's KG-Augmented Input accounts for diverse attack launching strategies and attack variants, outperforming baseline approaches in most scenarios with detection performance ≥94%. These results demonstrate the critical importance of building models on comprehensive attack knowledge rather than restrictive assumptions.

**Takeaway.** Current ML-NIDS models often rely on *restrictive* assumptions about how attacks are launched (e.g., single-channel attacks).

As a result, models tend to overfit to specific aggregation or attack strategies, and fail to detect variants that fall outside these patterns—*Attack Strategy Bias*. In contrast, KNOWML incorporates a diverse range of attack strategies into ML-NIDS, enabling broader generalization and, in most cases, outperforming baselines built on narrow assumptions.

**5.2.2. Out-of-Dimension Attacks.** Table 4 presents the detection performance against Out-of-Dimension attacks across different environments. It can be observed that existing Knowledge-Models completely fail to detect most Out-of-Dimension attacks, achieving 0% $F_1$-Score regardless of model type. These results expose a fundamental limitation of ML-NIDS systems built on limited knowledge. The KIT model, for example, depends solely on manually crafted heuristics and expert intuition, while the SFS model constructs its feature set based on general network monitoring metrics rather than attack-relevant characteristics. Both approaches prove inadequate against Out-of-Dimension threats. Only CIC (KG-model based on *Random Selection*) showed limited detection capability on Nkiller2 attacks, exactly demonstrating its random selection nature. It showed arbitrary coverage on Out-of-Dimension attacks, with performance on Nkiller2 attacks achieving at best 86.96% $F_1$-score in IoT networks and 68.97% $F_1$-score in enterprise-like networks. It can be also observed that model type and complexity play only secondary roles in detection performance. ML models unarguably contributed to detection, for example, switching from DA to KIT-ML improves CIC's Nkiller detection on IoT-23 by approximately 6%. However, it is the underlying KG-model that determines detection success. KNOWML's KG-Augmented Input produces significant performance improvements, validating our hypothesis that semantically relevant input enhances detection capabilities. ML models built on limited attack knowledge suffer severe generalization failures when confronting diverse attack variants in production environments, particularly Out-of-Dimension attacks. KNOWML achieves consistent improvements across all network environments and model types, with F1-score increases up to 99.9%. Figure 8 provides a detailed performance comparison between KNOWML and the best-performing baseline (CIC), demonstrating KNOWML's superiority over contemporary approaches.

**Takeaway.** Existing ML-NIDS are based on personal knwoledge of the system designers, and fail to generalize to variants of well-known DoS and Brute-Force attacks.

As a result, such models struggle to generalize to Out-of-Dimension attacks. In contrast, knowledge-based systems like KNOWML, which incorporate information about the broader attack space, demonstrating superior average performance on Out-of-Dimension attacks.

**5.2.3. Non-Throughput Based Attacks.** Finally, Table 5 presents performance against Non-Throughput Based

TABLE 3: Detection Accuracy of KNOWML and Baselines Under **M-N End-Point Attacks** in $F_1$-Score at FPR$\leq 0.1\%$. **Bold** values indicate the best-performance **for each ML model**. Improvements achieved by KNOWML are marked with ↑.

| Attack | GMM | | | | DA | | | | KIT-ML | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | SFS | CIC | KIT | KNOWML | SFS | CIC | KIT | KNOWML | SFS | CIC | KIT | KNOWML |
| DDoS-HTTP (IoT-23) | 0.8427 | 0.0022 | 0.0003 | ↑ **0.9475** | 0.8502 | 0.0101 | 0.2909 | ↑ **0.9476** | 0.8418 | 0.0033 | 0.8506 | ↑ **0.9476** |
| DDoS-SYN (IoT-23) | 0.0000 | 0.9099 | 0.0000 | ↑ **0.9979** | 0.0000 | 0.0004 | 0.9620 | **0.9979** | 0.0000 | 0.0000 | 0.9873 | ↑ **0.9979** |
| DDoS-TCP (IoT-23) | 0.0000 | 0.0000 | 0.4294 | ↑ **0.9997** | 0.0000 | 0.0000 | 0.9992 | ↑ **0.9997** | 0.0000 | 0.0000 | 0.9986 | ↑ **0.9999** |
| DDoS-PSHACK (IoT-23) | 0.0000 | 0.2405 | 0.0000 | ↑ **0.9250** | 0.0000 | 0.0008 | 0.9656 | ↑ **0.9753** | 0.0000 | 0.0078 | 0.9831 | ↑ **0.9772** |
| DDoS-RSTFIN (IoT-23) | 0.0000 | 0.7057 | 0.4383 | ↑ **1.0000** | 0.0000 | 0.0000 | 0.9983 | ↑ **1.0000** | 0.0000 | 0.0000 | 0.9950 | ↑ **1.0000** |
| DDoS-SynonymousIP (IoT-23) | 0.0000 | 0.0000 | 0.9771 | ↑ **1.0000** | 0.0000 | 0.8000 | 0.9999 | ↑ **1.0000** | 0.0000 | 0.8001 | 0.9999 | ↑ **1.0000** |
| DDoS-ACK Fragmentation (IoT-23) | **0.9801** | 0.8534 | 0.0007 | 0.9438 | **0.9803** | 0.8442 | 0.5375 | 0.9438 | **0.9802** | 0.8237 | 0.9476 | 0.9438 |
| DoS-HTTP (IoT-23) | 0.6676 | 0.1453 | 0.9173 | ↑ **0.9834** | 0.7277 | 0.0742 | 0.9991 | ↑ **0.9835** | 0.6905 | 0.0256 | 0.9748 | ↑ **0.9001** |
| DoS-SYN (IoT-23) | 0.0439 | 0.0146 | 0.0010 | ↑ **0.9003** | 0.0897 | 0.0001 | **0.9156** | 0.9000 | 0.0930 | 0.0000 | 0.9990 | ↑ **0.9835** |
| DoS-TCP (IoT-23) | 0.1499 | 0.7509 | 0.3583 | ↑ **0.9999** | 0.3273 | 0.3329 | 0.9786 | ↑ **0.9999** | 0.3357 | 0.0000 | 0.9988 | ↑ **0.9999** |
| DoS Hulk (IDS-17) | 0.0228 | 0.1843 | 0.0000 | ↑ **0.9178** | 0.0406 | 0.2616 | 0.9032 | ↑ **0.9677** | 0.0270 | 0.5604 | 0.8641 | ↑ **0.9679** |
| DoS GoldenEye (IDS-17) | 0.3609 | 0.2990 | 0.0000 | ↑ **0.9934** | 0.3122 | 0.4376 | 0.0005 | ↑ **0.9981** | 0.3052 | 0.7792 | 0.0000 | ↑ **0.9953** |
| Fiberfox GET (CAP) | 0.2203 | 0.1623 | 0.6740 | ↑ **0.9766** | 0.0100 | 0.3205 | **0.9957** | 0.9742 | 0.0320 | 0.1619 | **0.9997** | 0.9742 |
| Fiberfox SLOW (CAP) | 0.9999 | 0.1828 | 0.8242 | ↑ **0.9825** | 0.0000 | 0.2514 | 0.9959 | ↑ **0.9796** | 0.0000 | 0.1610 | 0.9997 | ↑ **0.9796** |
| Fiberfox STRESS (CAP) | **0.9988** | 0.2492 | 0.0000 | 0.9605 | 0.6269 | 0.9914 | 0.9869 | ↑ **0.9291** | 0.8814 | 0.9930 | **0.9992** | 0.9295 |

TABLE 4: Detection Accuracy of KNOWML and Baselines Under **Out-of-Dimension Attacks** in $F_1$-Score at FPR$\leq 0.1\%$. **Bold** values indicate the best-performing baseline **for each attack**. ↑ indicates the improvement of KNOWML on the same model with the best-performing baseline.

| Knowledge-Model | ML model | IoT Network (IoT-23) | | | | Enterprise Network (IDS-17) | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | HTTP Mal | HTTP Overflow | Nkiller2 | SSH BF | HTTP Mal | HTTP Overflow | Nkiller2 | SSH BF |
| CIC | DA | 0.0000 | 0.0000 | 0.7792 | 0.0000 | 0.0082 | 0.0000 | 0.2182 | 0.0000 |
| | GMM | 0.0000 | 0.0000 | 0.8696 | 0.0000 | 0.2588 | 0.0000 | 0.5660 | 0.0000 |
| | KIT-ML | 0.0000 | 0.0000 | 0.8451 | 0.0000 | 0.3582 | 0.0000 | 0.6897 | 0.0000 |
| KIT | DA | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0610 | 0.0621 | 0.0980 | 0.0028 |
| | GMM | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| | KIT-ML | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| SFS | DA | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0135 | 0.0317 | 0.0000 | 0.0002 |
| | GMM | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| | KIT-ML | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 | 0.0000 |
| KNOWML | DA | ↑ **0.9985** | 0.9770 | 0.9190 | ↑ **0.9750** | 0.9630 | 0.9630 | 0.9319 | 0.4476 |
| | GMM | 0.9844 | ↑ **0.9817** | 0.8841 | 0.9727 | ↑ **0.9974** | ↑ **0.9974** | ↑ **0.9820** | ↑ **0.9607** |
| | KIT-ML | ↑ **0.9985** | 0.9800 | ↑ **0.9276** | ↑ **0.9750** | 0.9956 | 0.9956 | 0.9046 | 0.2243 |

**Mimicry Attacks.** When ML models are built on strong *a priori* assumptions about what kind of attacks ML-NIDS might encounter during production, they fail to generalize on Non-Throughput Based Attacks. More concretely, observe how performance for KIT decreases from 97.48% on HTTP DoS (Table 3) to 0-51.98% F1-score (Table 5) on variants of Non-Throughput Based HTTP DoS variants such as Fiberfox AVB, BYPASS, or SlowLoris. Note that this category is different from Out-of-Dimension attacks where attacker operate out of predefined spaces, here attacker still operate *within* the same space. Remember that KIT was created and evaluated primarily on throughput-based variants that rely on significant rise in rate and jitter statistics. Any attacks that resemble benign traffic across those dimensions will be missed by KIT-baseline regardless of ML-model type. Limited detection capability was also observed for SFS baseline methods, which employ standardized approaches

designed for network monitoring (CISCO Netflow), and CIC baseline methods that rely on random selection. SFS achieves on average 22.49% $F1$-score attack variants, and CIC 20.69% $F1$-Score on average, substantiating our claim that non-attack knowledge-based approaches fail to generalize on more complex attack variants that employ mimicry techniques. KNOWML, which accounts for diverse attack variants and strategies, outperforms the baseline approaches and improves performance regardless of model type, with average 91.70%-93.28% $F1$-score. Also note note that for this category of attacks, increasing model complexity proves unnecessary (from DA to KIT-ML (Ensemble of DA)), demonstrating the utility of KNOWML even on simple models.

**Takeaway.** ML-NIDS models often embed the designer's prior beliefs about possible attack variants. A limited understanding of threats, such as assuming that DoS and Brute-

TABLE 5: Detection Accuracy of KNOWML and Baselines Under **Non-Throughput Attacks** in $F_1$-Score at FPR$\leq 0.1\%$. **Bold** values highlight the best-performing baseline **for each ML model**. Improvements achieved by KNOWML are marked with $\uparrow$.

| Attack | GMM | | | | DA | | | | KIT-ML | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | SFS | CIC | KIT | KNOWML | SFS | CIC | KIT | KNOWML | SFS | CIC | KIT | KNOWML |
| Low rate TCP (SIM) | 0.0000 | **0.8696** | 0.0000 | 0.7699 | 0.0000 | 0.7792 | 0.0000 | $\uparrow$ **0.9389** | 0.0000 | 0.8451 | 0.0000 | $\uparrow$ **0.9277** |
| Fiberfox AVB (CAP) | **0.9622** | 0.0000 | 0.0000 | 0.8523 | 0.0042 | 0.0000 | 0.0000 | $\uparrow$ **0.8994** | 0.0042 | 0.0000 | 0.0000 | 0.8638 |
| Fiberfox BYPASS (CAP) | 0.3205 | 0.0000 | 0.2751 | $\uparrow$ **0.9185** | 0.0000 | 0.5509 | 0.0045 | $\uparrow$ **0.9647** | 0.0000 | 0.0000 | 0.0030 | $\uparrow$ **0.8431** |
| SSH Patator (P=0) (CAP) | 0.0000 | 0.0000 | 0.0002 | $\uparrow$ **0.9998** | 0.0000 | 0.0000 | 0.0000 | $\uparrow$ **0.9998** | 0.0000 | 0.0000 | 0.0000 | $\uparrow$ **0.9998** |
| SSH Patator (P=1) (CAP) | 0.4175 | 0.0000 | 0.0000 | $\uparrow$ **0.9994** | 0.0000 | 0.0000 | 0.0000 | $\uparrow$ **0.9994** | 0.0000 | 0.0000 | 0.0000 | $\uparrow$ **0.9994** |
| DoS-Slowloris (IoT-23) | 0.0568 | 0.0000 | 0.0000 | $\uparrow$ **0.9803** | 0.0000 | 0.0000 | 0.0028 | $\uparrow$ **0.9714** | 0.0000 | 0.0000 | 0.0000 | $\uparrow$ **0.9818** |
| DoS-Slowhttptest (IDS-17) | 0.0031 | 0.0000 | 0.0000 | $\uparrow$ **0.9214** | 0.0002 | 0.3117 | 0.0018 | $\uparrow$ **0.9652** | 0.0000 | 0.0281 | 0.0000 | $\uparrow$ 0.9647 |
| DDoS-SlowLoris (IDS-17) | 0.0392 | 0.0027 | 0.0031 | $\uparrow$ **0.9803** | 0.0303 | 0.0131 | 0.0015 | $\uparrow$ **0.7548** | 0.0302 | 0.0001 | 0.5198 | $\uparrow$ **0.7560** |
| Average | 0.2249 | 0.1090 | 0.0348 | **0.9212** | 0.0043 | 0.2069 | 0.0010 | **0.9383** | 0.0043 | 0.1092 | 0.0654 | **0.9170** |



(a) KNOWML on IoT-23

(b) CIC on IoT-23

(c) KNOWML on IDS-17
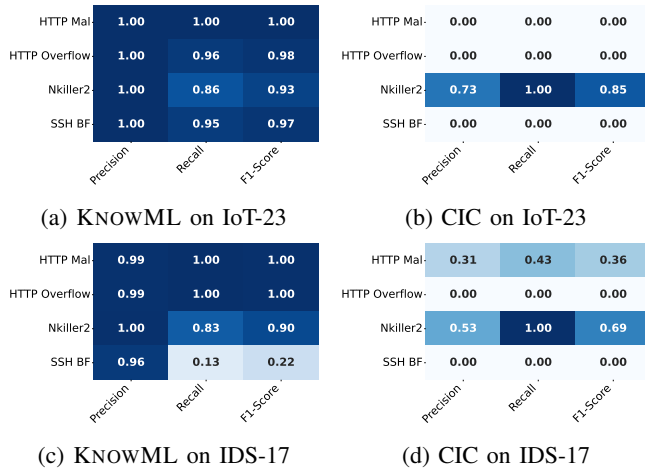
(d) CIC on IDS-17

Figure 8: Performance comparison between Best Performing Baseline and KNOWML with KIT-ML

Force attacks manifest through jitter statistics, gets encoded into the ML models. This restricts the model's ability to generalize effectively to non-throughput-based attacks. KNOWML, which accounts for diverse attack variants, is capable of generalizing to Non-Throughput variants.

# 6. Related Work

This section clarifies the novelty of KNOWML compared with three key areas of closely related work: Knowledge Graph Applications in Security, Automatic Feature Engineering, and Research on the Limitations of ML-NIDS.

**Knowledge Graphs (KG) Applications in Security.** KGs have seen widespread use in various security-related tasks [61]. These include: 1) knowledge management, where KGs enhance visual analytics tools to support expert decision-making [62]; 2) threat discovery, where KGs help identify potential zero-day attacks—often by modeling them as Attack Graphs to investigate Advanced Persistent Threats (APTs) [10]; and 3) attack investigation, where KGs support both attack attribution [63] and the exploration of possible attack paths [64]. In most of these efforts, KGs serve as auxiliary resources. For example, SEPSES uses a KG to explain Snort rules [9], but does not apply the KG directly within NIDS tasks.

In contrast, KNOWML integrates KGs directly into the training and evaluation stages of the ML-NIDS pipeline. This represents a fundamental shift from using KGs solely as external knowledge to embedding them within the core detection process.

A few studies have explored this direct integration. For instance, Yang et al. [65] construct a KG from the NSL-KDD dataset by modeling co-occurrence relationships between features. These relationships are weighted based on frequency and embedded into a CNN-BiLSTM architecture. However, KNOWML differs significantly. It does not build the KG from a specific dataset, thus avoiding the dataset dependency that often limits generalization due to the homogeneity of common benchmarks [3]. Moreover, KNOWML employs symbolic rules to extract semantical relevant information to enhance detection capabilities of ML-NIDS.

**Automated Feature Engineering.** Prior work has explored automatic feature extraction for machine learning applications. These efforts fall into two main categories: 1) *Representation learning*, which focuses on automatically generating feature representations—for example, searching for optimal data representation for classification using methods such as reinforcement learning (RL) [66]; 2) *Feature mining*, which extracts semantically meaningful information from external sources, as seen in systems like FeatureSmith [67] and ChainSmith [68]. These approaches primarily adopt data-driven strategies, which can lead to overfitting on the training dataset or limit applicability to specific domains. FeatureSmith and ChainSmith were created for malware detection with goal create a minimal feature sets sufficient to recognize known malware families on the Drebin [69] dataset. In contrast, KNOWML addresses a different problem domain. It automates feature generation by analyzing attacker behavior at the implementation level and produces KG-augmented Inputs designed to capture a wide spectrum of attack variants. Unlike prior approaches, KNOWML does

not aim to identify a minimal feature subset for a fixed set of threats in a provided dataset.

**Works exploring limitation of ML-NIDS.** Recent critiques of ML-NIDS research have focused on two main perspectives: i) dataset inadequacies [3], [22], [70] and ii) flawed ML design practices [11], [71]. Flood et al. [3] identified "bad smells"—questionable practices in benchmark NIDS datasets—highlighting issues such as *poor dataset diversity*, that creates highly-dependent features. These limitations prevent ML models from extracting meaningful information, which restricts their ability to generalize, and results in artificially-inflated performance scores.

Prior critiques focus primarily on models and datasets. In contrast, this work adopts a knowledge-centric perspective. We argue that limited understanding of attacker behavior leads to flawed design assumptions in ML-NIDS, weakening detection capabilities. KNOWML addresses this by constructing knowledge-based representations of attacks, enriching input features with semantic context and improving detection across diverse attack variants.

## 7. Discussion

**Cost of Knowledge Graph Construction.** KNOWML provides a scalable and cost-efficient solution for extracting knowledge from large-scale codebases such as GitHub. For instance, extracting data (including text embeddings) from 2,333 repositories (TCP DoS) using GPT-4o-mini costs approximately $13. KNOWML also significantly reduces human labor. If we assume domain expert in network security would require about 10 minutes to analyze a single repository, to analyze all 7,853 repositories examined in this study would take approximately 157 workdays (assuming an 8-hour workday and no interruptions). In contrast, KNOWML completes the same task in just 26 hours and 7 minutes, demonstrating substantial savings in both manual effort and financial cost.

**Manual Validation for Feature Mapping.** The process of mapping strategies into features for detection, as discussed in §3.4, still requires expert input and manual validation. This step is susceptible to subjective interpretation. Experts may disagree on what constitutes an appropriate feature—for example, how to construct a composite feature representing strategies that simultaneously set the TCP ACK flag and use IP fragmentation (see Composite Feature Generation in §3.4.2). KNOWML introduces a general framework for structural strategy analysis and, within this paper, the authors propose a set of features that best capture the intended strategy and corresponding attack behaviors. We do not claim these feature mappings to be definitive or universally applicable. Rather, we encourage future research to explore alternative mappings within the proposed framework. As language models continue to improve, we plan to investigate the use of LLMs to automate this component in future work.

**Choice of Attack Families.** In this paper we examine the following families: TCP DoS, HTTP DoS, and SSH Brute Force. These were selected for two reasons. First, they represent some of the most frequently encountered attack types in real-world deployments [72]. Second, they are widely studied, allowing us to assess what additional knowledge KNOWML can uncover beyond existing literature. While it is easier to show gaps in ML-NIDS using obscure or uncommon attacks, doing so with well-known and heavily researched variants demonstrates the practical utility of our approach on more challenging problems.

**Considerations on Sophisticated Adversaries.** In this paper, we evaluate KNOWML specifically for detection tasks. Although we do not directly assess its resilience to adversarial machine learning [4], exploring this direction presents an important avenue for future work. Our intuition is that KNOWML raises the difficulty for adversaries by modeling the *problem space* [73] with the parameters that attackers manipulate to launch an attack. We plan to investigate how knowledge-based systems can help mitigate adversarial threats in future research.

**Other Code Repositories.** In this study, we use GitHub as the primary source for knowledge extraction. In future work, we aim to extend KNOWML to other sources such as Exploit-DB [74], which offers code for disclosed vulnerabilities. Adapting to such sources would require changes to the prompt engineering techniques used in knowledge graph construction, but the remaining components of the system should remain reusable.

## 8. Conclusion

In this work, we introduced KNOWML, a framework for systematically analyzing attack strategies based on their implementation logic. KNOWML offers actionable insights that enhance the performance of ML-based NIDS through symbolic reasoning and Knowledge-Augmented Input.

Our prototype demonstrates that KNOWML helps overcome common pitfalls in ML-NIDS design, particularly those caused by limited domain knowledge or assumptions made by system developers. The results show that KNOWML achieves competitive detection performance with a low false positive rate. Furthermore, it exposes limitations in existing ML-NIDS approaches by successfully identifying realistic attack variants, including those from families often considered easy to detect. These findings highlight both the utility of our approach and the need for more robust knowledge integration in future ML-NIDS designs.

## 9. Acknowledgement

## References

[1] Y. Mirsky, T. Doitshman, Y. Elovici, and A. Shabtai, "Kitsune: an ensemble of autoencoders for online network intrusion detection," *arXiv preprint arXiv:1802.09089*, 2018.

[2] M. Soltani, B. Ousat, M. J. Siavoshani, and A. H. Jahangir, "An adaptable deep learning-based intrusion detection system to zero-day attacks," *Journal of Information Security and Applications*, vol. 76, p. 103516, 2023.

[3] R. Flood, G. Engelen, D. Aspinall, and L. Desmet, "Bad design smells in benchmark nids datasets," in *2024 IEEE 9th European Symposium on Security and Privacy (EuroS&P)*. IEEE, 2024, pp. 658–675.

[4] G. Apruzzese, P. Laskov, E. Montes de Oca, W. Mallouli, L. Brdalo Rapa, A. V. Grammatopoulos, and F. Di Franco, "The role of machine learning in cybersecurity," *Digital Threats: Research and Practice*, vol. 4, no. 1, pp. 1–38, 2023.

[5] R. Sommer and V. Paxson, "Outside the closed world: On using machine learning for network intrusion detection," in *2010 IEEE symposium on security and privacy*. IEEE, 2010, pp. 305–316.

[6] C. Fu, Q. Li, M. Shen, and K. Xu, "Realtime robust malicious traffic detection via frequency domain analysis," in *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*, 2021, pp. 3431–3446.

[7] A. Shiravi, H. Shiravi, M. Tavallaee, and A. A. Ghorbani, "Toward developing a systematic approach to generate benchmark datasets for intrusion detection," *computers & security*, vol. 31, no. 3, pp. 357–374, 2012.

[8] I. Sharafaldin, A. H. Lashkari, A. A. Ghorbani *et al.*, "Toward generating a new intrusion detection dataset and intrusion traffic characterization." *ICISSp*, vol. 1, pp. 108–116, 2018.

[9] E. Kiesling, A. Ekelhart, K. Kurniawan, and F. Ekaputra, "The sepses knowledge graph: an integrated resource for cybersecurity," in *International Semantic Web Conference*. Springer, 2019, pp. 198–214.

[10] J. Zhang, J. Zheng, Z. Zhang, T. Chen, Y.-a. Tan, Q. Zhang, and Y. Li, "Att&ck-based advanced persistent threat attacks risk propagation assessment model for zero trust networks," *Computer Networks*, vol. 245, p. 110376, 2024.

[11] D. Arp, E. Quiring, F. Pendlebury, A. Warnecke, F. Pierazzi, C. Wressnegger, L. Cavallaro, and K. Rieck, "Dos and don'ts of machine learning in computer security," in *31st USENIX Security Symposium (USENIX Security 22)*, 2022, pp. 3971–3988.

[12] P. An, Z. Wang, and C. Zhang, "Ensemble unsupervised autoencoders and gaussian mixture model for cyberattack detection," *Information Processing & Management*, vol. 59, no. 2, p. 102844, 2022.

[13] A. Thakkar and R. Lohiya, "A review on machine learning and deep learning perspectives of ids for iot: recent updates, security issues, and challenges," *Archives of Computational Methods in Engineering*, vol. 28, no. 4, pp. 3211–3243, 2021.

[14] D. Wagner and P. Soto, "Mimicry attacks on host-based intrusion detection systems," in *Proceedings of the 9th ACM Conference on Computer and Communications Security*, 2002, pp. 255–264.

[15] National Institute of Standards and Technology, "Nist: National institute of standards and technology," U.S. Department of Commerce, 2023, accessed on September 23, 2024. [Online]. Available: https://www.nist.gov/

[16] MITRE Corporation, "CVE-2008-4609: Tcp implementation vulnerability," https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2008-4609.

[17] N. V. Database, "CVE-2009-1926: Tcp/ip null pointer dereference vulnerability," https://nvd.nist.gov/vuln/detail/CVE-2009-1926, 2009, accessed: 2024-10-20.

[18] National Institute of Standards and Technology, "CVE-2024-32984: Yamux stream multiplexer unbounded pending frames vulnerability," https://nvd.nist.gov/vuln/detail/CVE-2024-32984.

[19] "Cve-2024-35296 detail," https://nvd.nist.gov/vuln/detail/CVE-2024-35296, July 2024, accessed: 2025-02-26.

[20] Apache Software Foundation, "CVE-2025-31650: Improper Input Validation in Apache Tomcat," https://nvd.nist.gov/vuln/detail/CVE-2025-31650, 2025, published: 2025-04-28, Last Modified: 2025-05-06. [Online]. Available: https://nvd.nist.gov/vuln/detail/CVE-2025-31650

[21] "CVE-2020-35863 detail," https://nvd.nist.gov/vuln/detail/CVE-2020-35863, December 2020, accessed: 2025-02-26.

[22] G. Apruzzese, P. Laskov, and J. Schneider, "Sok: Pragmatic assessment of machine learning for network intrusion detection," 2023. [Online]. Available: https://arxiv.org/abs/2305.00550

[23] E. F. Kendall and D. L. McGuinness, *Ontology engineering*. Morgan & Claypool Publishers, 2019.

[24] O. Sainz, I. García-Ferrero, R. Agerri, O. L. de Lacalle, G. Rigau, and E. Agirre, "Gollie: Annotation guidelines improve zero-shot information-extraction," *arXiv preprint arXiv:2310.03668*, 2023.

[25] H. B. Giglou, J. D'Souza, F. Engel, and S. Auer, "Llms4om: Matching ontologies with large language models," *arXiv preprint arXiv:2404.10317*, 2024.

[26] B. Min, H. Ross, E. Sulem, A. P. B. Veyseh, T. H. Nguyen, O. Sainz, E. Agirre, I. Heintz, and D. Roth, "Recent advances in natural language processing via large pre-trained language models: A survey," *ACM Comput. Surv.*, vol. 56, no. 2, Sep. 2023. [Online]. Available: https://doi.org/10.1145/3605943

[27] Z. Li, J. Zeng, Y. Chen, and Z. Liang, "Attackg: Constructing technique knowledge graph from cyber threat intelligence reports," in *European Symposium on Research in Computer Security*. Springer, 2022, pp. 589–609.

[28] D. Edge, H. Trinh, N. Cheng, J. Bradley, A. Chao, A. Mody, S. Truitt, and J. Larson, "From local to global: A graph rag approach to query-focused summarization," 2024. [Online]. Available: https://arxiv.org/abs/2404.16130

[29] M. Dean and G. Schreiber, "Owl web ontology language guide: Feature synopsis," https://www.w3.org/TR/2004/REC-owl-features-20040210/#s4, 2004, w3C Recommendation, Accessed: 2025-06-01.

[30] K. Dentler, R. Cornet, A. Ten Teije, and N. De Keizer, "Comparison of reasoners for large ontologies in the owl 2 el profile," *Semantic Web*, vol. 2, no. 2, pp. 71–87, 2011.

[31] Y. Kuratov, A. Bulatov, P. Anokhin, D. Sorokin, A. Sorokin, and M. Burtsev, "In search of needles in a 10m haystack: Recurrent memory finds what llms miss," *arXiv preprint arXiv:2402.10790*, 2024.

[32] N. F. Liu, K. Lin, J. Hewitt, A. Paranjape, M. Bevilacqua, F. Petroni, and P. Liang, "Lost in the middle: How language models use long contexts," *Transactions of the Association for Computational Linguistics*, vol. 12, pp. 157–173, 2024.

[33] L. Huang, W. Yu, W. Ma, W. Zhong, Z. Feng, H. Wang, Q. Chen, W. Peng, X. Feng, B. Qin *et al.*, "A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions," *ACM Transactions on Information Systems*, vol. 43, no. 2, pp. 1–55, 2025.

[34] Hiett, "serverless-redis-http," https://github.com/hiett/serverless-redis-http, 2020, accessed: 2025-06-02.

[35] OpenAI, "Structured outputs with language models," https://cookbook.openai.com/examples/structured_outputs_intro, 2023, accessed: 2025-06-02.

[36] ——, "Embedding models," 2024, accessed: 2025-06-06. [Online]. Available: https://platform.openai.com/docs/guides/embeddings/embedding-models

[37] D. Müllner, "Modern hierarchical, agglomerative clustering algorithms," *arXiv preprint arXiv:1109.2378*, 2011.

[38] J. Diaz-Rodriguez, "k-llmmeans: Summaries as centroids for interpretable and scalable llm-based text clustering," *arXiv preprint arXiv:2502.09667*, 2025.

[39] B. P. Welford, "Note on a method for calculating corrected sums of squares and products," *Technometrics*, vol. 4, no. 3, pp. 419–420, 1962.

[40] GitHub, "GitHub REST API - Search," 2022, accessed: 2025-06-02. [Online]. Available: https://docs.github.com/en/rest/search/search?apiVersion=2022-11-28

[41] Y. Hu, F. Zou, J. Han, X. Sun, and Y. Wang, "Llm-tikg: Threat intelligence knowledge graph construction utilizing large language model," *Computers & Security*, vol. 145, p. 103999, 2024.

[42] S. Bogdanov, A. Constantin, T. Bernard, B. Crabbé, and E. Bernard, "Nuner: entity recognition encoder pre-training via llm-annotated data," *arXiv preprint arXiv:2402.15343*, 2024.

[43] Á. García-Barragán, A. González Calatayud, O. Solarte-Pabón, M. Provencio, E. Menasalvas, and V. Robles, "Gpt for medical entity recognition in spanish," *Multimedia Tools and Applications*, pp. 1–20, 2024.

[44] S. Sawadogo, A. Sabane, R. Kafando, A. K. Kabore, and T. F. Bissyande, "Revisiting the non-determinism of code generation by the gpt-3.5 large language model," in *2025 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER)*. IEEE, 2025, pp. 36–44.

[45] R. Doriguzzi-Corin, S. Millar, S. Scott-Hayward, J. Martínez-del Rincón, and D. Siracusa, "Lucid: A practical, lightweight deep learning solution for ddos attack detection," *IEEE Transactions on Network and Service Management*, vol. 17, no. 2, pp. 876–889, 2020.

[46] A. Kuzmanovic and E. W. Knightly, "Low-rate tcp-targeted denial of service attacks: the shrew vs. the mice and elephants," in *Proceedings of the 2003 conference on Applications, technologies, architectures, and protocols for computer communications*, 2003, pp. 75–86.

[47] A. Kachayev, "Fiberfox," https://github.com/kachayev/fiberfox, accessed: 2025-04-18.

[48] Lanjelot, "Patator: A multi-purpose brute-forcer," https://github.com/lanjelot/patator/blob/master/src/patator/patator.py, 2025, accessed: 2025-04-20.

[49] E. C. P. Neto, S. Dadkhah, R. Ferreira, A. Zohourian, R. Lu, and A. A. Ghorbani, "Ciciot2023: A real-time dataset and benchmark for large-scale attacks in iot environment," *Sensors*, vol. 23, no. 13, p. 5941, 2023.

[50] G. Apruzzese, A. Fass, and F. Pierazzi, "When adversarial perturbations meet concept drift: an exploratory analysis on ml-nids," in *Proceedings of the 2024 Workshop on Artificial Intelligence and Security*, 2024, pp. 149–160.

[51] X. Jiang, S. Liu, A. Gember-Jacobson, A. N. Bhagoji, P. Schmitt, F. Bronzino, and N. Feamster, "Netdiffusion: Network data augmentation through protocol-constrained traffic generation," *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, vol. 8, no. 1, pp. 1–32, 2024.

[52] D. Han, Z. Wang, Y. Zhong, W. Chen, J. Yang, S. Lu, X. Shi, and X. Yin, "Evaluating and improving adversarial robustness of machine learning-based network intrusion detectors," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 8, pp. 2632–2647, 2021.

[53] G. Apruzzese, L. Pajola, and M. Conti, "The cross-evaluation of machine learning-based network intrusion detection systems," *IEEE Transactions on Network and Service Management*, vol. 19, no. 4, pp. 5152–5169, 2022.

[54] I. Arnaldo and K. Veeramachaneni, "The holy grail of' systems for machine learning" teaming humans and machine learning for detecting cyber threats," *ACM SIGKDD Explorations Newsletter*, vol. 21, no. 2, pp. 39–47, 2019.

[55] C. G. Cordero, E. Vasilomanolakis, A. Wainakh, M. Mühlhäuser, and S. Nadjm-Tehrani, "On generating network traffic datasets with synthetic attacks for intrusion detection," *ACM Transactions on Privacy and Security (TOPS)*, vol. 24, no. 2, pp. 1–39, 2021.

[56] M. Sarhan, S. Layeghy, and M. Portmann, "Towards a standard feature set for network intrusion detection system datasets," *Mobile networks and applications*, pp. 1–14, 2022.

[57] Cisco Systems, Inc., "Cisco ios netflow," 2025, accessed: 2025-05-25. [Online]. Available: https://www.cisco.com/c/en/us/products/ios-nx-os-software/ios-netflow/index.html

[58] H. Zhang, L. Huang, C. Q. Wu, and Z. Li, "An effective convolutional neural network based on smote and gaussian mixture model for intrusion detection in imbalanced dataset," *Computer Networks*, vol. 177, p. 107315, 2020.

[59] S. K. Amalapuram, A. Tadwai, R. Vinta, S. S. Channappayya, and B. R. Tamma, "Continual learning for anomaly based network intrusion detection," in *2022 14th International Conference on COMmunication Systems & NETworkS (COMSNETS)*. IEEE, 2022, pp. 497–505.

[60] Y. Mirsky, "Kitsune-py: A network intrusion detection system based on incremental statistics (afterimage) and an ensemble of autoencoders (kitnet)," https://github.com/ymirsky/Kitsune-py, 2018, accessed: 2025-04-15.

[61] X. Zou, "A survey on application of knowledge graph," in *Journal of Physics: Conference Series*, vol. 1487, no. 1. IOP Publishing, 2020, p. 012016.

[62] S. Noel, E. Harley, K. H. Tam, M. Limiero, and M. Share, "Cygraph: graph-based analytics and visualization for cybersecurity," in *Handbook of statistics*. Elsevier, 2016, vol. 35, pp. 117–167.

[63] Z. Zhu, R. Jiang, Y. Jia, J. Xu, and A. Li, "Cyber security knowledge graph based cyber attack attribution framework for space-ground integration information network," in *2018 IEEE 18th International Conference on Communication Technology (ICCT)*. IEEE, 2018, pp. 870–874.

[64] S. Noel, E. Harley, K. H. Tam, and G. Gyor, "Big-data architecture for cyber attack graphs," *MITRE case*, no. 14-3549, 2014.

[65] X. Yang, G. Peng, D. Zhang, and Y. Lv, "An enhanced intrusion detection system for iot networks based on deep learning and knowledge graph," *Security and Communication Networks*, vol. 2022, no. 1, p. 4748528, 2022.

[66] U. Khurana, H. Samulowitz, and D. Turaga, "Feature engineering for predictive modeling using reinforcement learning," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.

[67] Z. Zhu and T. Dumitraş, "Featuresmith: Automatically engineering features for malware detection by mining the security literature," in *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, 2016, pp. 767–778.

[68] Z. Zhu and T. Dumitras, "Chainsmith: Automatically learning the semantics of malicious campaigns by mining threat intelligence reports," in *2018 IEEE European symposium on security and privacy (EuroS&P)*. IEEE, 2018, pp. 458–472.

[69] D. Arp, M. Spreitzenbarth, M. Hubner, H. Gascon, K. Rieck, and C. Siemens, "Drebin: Effective and explainable detection of android malware in your pocket." in *Ndss*, vol. 14, no. 1, 2014, pp. 23–26.

[70] G. Engelen, V. Rimmer, and W. Joosen, "Troubleshooting an intrusion detection dataset: the cicids2017 case study," in *2021 IEEE Security and Privacy Workshops (SPW)*, 2021, pp. 7–12.

[71] L. D'hooge, M. Verkerken, B. Volckaert, T. Wauters, and F. De Turck, "Establishing the contaminating effect of metadata feature inclusion in machine-learned network intrusion detection models," in *International Conference on Detection of Intrusions and Malware, and Vulnerability Assessment*. Springer, 2022, pp. 23–41.

[72] CrowdStrike, Inc., "Common cyberattacks: Definitions, examples and prevention tips," https://www.crowdstrike.com/en-us/cybersecurity-101/cyberattacks/common-cyberattacks/, 2025, accessed: 2025-06-04.

[73] F. Pierazzi, F. Pendlebury, J. Cortellazzi, and L. Cavallaro, "Intriguing properties of adversarial ml attacks in the problem space," in *2020 IEEE symposium on security and privacy (SP)*. IEEE, 2020, pp. 1332–1349.

[74] "The exploit database," https://www.exploit-db.com/, 2025, accessed: 2025-06-04.

[75] (2009) CVE-2009-1926 detail. National Vulnerability Database (NVD), NIST. [Online]. Available: https://nvd.nist.gov/vuln/detail/CVE-2009-1926

[76] National Vulnerability Database, "CVE-2024-39316: Regular Expression Denial of Service (ReDoS) in Rack::Request::Helpers," https://nvd.nist.gov/vuln/detail/CVE-2024-39316, 2024, accessed: 2025-04-10.

[77] ——, "CVE-2021-31166: HTTP Protocol Stack Remote Code Execution Vulnerability," https://nvd.nist.gov/vuln/detail/CVE-2021-31166, 2021, accessed: 2025-04-18.

[78] Swisskyrepo, "PayloadsAllTheThings: Command Injection," https://github.com/swisskyrepo/PayloadsAllTheThings/blob/master/Command%20Injection/README.md, 2025, accessed: 2025-04-12.

# Appendix A.
# Simulated Attacks

We simulated 6 attack scenarios for both IoT-23 and IDS-17 datasets, including the aforementioned Nkiller2 attack. In total, we simulated 6 attacks. Each attack simulates a real-world *known* vulnerability. The simulations are conducted by taking an hour of benign traffic. For each simulation there was least and was send at at least 2 times inrease in rate compared to the original benign traffic (if not stated otherwise), the goals here is to shows that even with an increase in the average byte rate those attacks cannot be still captured dues to insufficient flow definitions. The configuration and set up details is given as follows:

1) **Low rate TCP:** Throttling the TCP SYN attack by sending packets in bursts of 1 second and idle time of 1s. Simulating a well-established Low-Rate TCP Targeted attack [46] targeting TCP Retransmission Timeout (RTO) mechanism.

2) **Nkiller2:** Setting TCP window size of TCP packets to 0 and sending packets at 2.2x rate as the benign traffic. Simulating the vulnerability [16], [18], [75].

3) **HTTP DoS (HTTP Overflow):** Setting HTTP Accept-Language to overflow of a long valid string, causing the server to spend excessive time processing (as described in [76]). The packets are sent at 2.2x rate as benign traffic.

4) **HTTP DoS (Mal Header):** Setting HTTP Accept-Encoding to manifold of invalid values which can cause kernel to crash, simulating vulnerability [77]. The packets are sent at 2.2x rate as benign traffic.

5) **HTTP Request Smuggling (Bypass Char 1):** Injecting malicious commands into HTTP packet headers using bypass space technique (for details refer to [78]).

6) **HTTP Request Smuggling n (Bypass Space 2):** Injecting malicious commands into HTTP packet headers using bypass space technique (for details refer to [78]).

# Appendix B.
# Captured Attacks

The CAP dataset was collected on two hosts: a MacBook running Sonoma 14.1 and a Dell PC running Ubuntu 14. The Dell PC was configured with an nginx http service when http dos was captures or with ssh service when brute force was launched. The Dell laptop acted as the victim. All traffic was captured using Wireshark on the victim. The MacBook machine acted as the attacker and launched attacks by replaying 5 different HTTP DoS attack strategies using Fiberfox [47]. Each of the following attack strategies ran for 20 minutes (from `README.md` for the repository Fiberfox [47]) or SSH Patator (the first 5 are fiberfox attacks):

1) `--AVB`: Issues HTTP GET packets into an open connection with long delays between send operations.

2) `--GET`: Sends randomly generated HTTP GET requests over an open TCP connection.

3) `--STRESS`: Sends a sequence of HTTP requests with a large body over a single open TCP connection.

4) `--BYPASS`: Sends HTTP GET requests over an open TCP connection, reads response back.

5) `--SLOW`: Similarly to STRESS, issues HTTP requests and attempts to keep connection utilized by reading back a single byte and sending additional payload with time delays between send operations.

6) `--SSH_PERSISTENT`: If set to 0 then it create a new connection for each attempt otherwise it will send multiple brute force attempt through the same connection until.

# Appendix C.
# Additional Plots, Figures, and Appendices

---

**Algorithm 1** Register Network Statistics at Time $t$

---

1: **function** UPDATE($CS_{t-x}, W_{t-x}, SSR_{t-x}, \sigma_{t-x}, x_i$)
2: $\quad CS_t \leftarrow CS_{t-x} + x_i$
3: $\quad W_t \leftarrow W_{t-x} + 1$
4: $\quad r_t \leftarrow \frac{CS_t}{W_t}$
5: $\quad SSR_t \leftarrow SSR_{t-x} + (x_i - r_t)^2$
6: $\quad \sigma_t \leftarrow \sqrt{\frac{SSR_t}{W_t}}$
7: $\quad$ **return** $CS_t, W_t, r_t, SSR_t, \sigma_t$
8: **end function**

---

TABLE 6: Grid search parameters and number of combinations for each model.

| Model | Parameters | Combinations |
|---|---|---|
| GMM | `n_components`: [3, 5, 7, 10, 20, 30, 40]<br>`covariance_type`: [diag, spherical] | 14 |
| DA | `hidden_ratio`: [0.1, 0.3, 0.5]<br>`learning_rate`: [0.0001, 0.001]<br>`corruption_level`: [0.05, 0.1, 0.2]<br>`activation`: [relu, tanh]<br>`l2_reg`: [0.0001, 0.001] | 72 |
| KIT-ML (Ensemble of DA) | `maxAE`: [1, 10]<br>`learning_rate`: [0.001, 0.01]<br>`hidden_ratio`: [0.25, 0.5] | 8 |

| Category | Example Statistics | # Features |
|---|---|---|
| General Flow Metrics | Duration, packet counts, rates, throughput, transmission rate | 13 |
| Packet Size Statistics | Total, mean, std deviation, SSR | 4 |
| Timing Metrics | Inter-arrival time (IAT), Time-To-Live (TTL) | 20 |
| TCP Metrics | TCP flags, window size stats, payload stats, TCP state values | 43 |
| HTTP Metrics | Request methods, payload/header stats, response codes | 66 |
| SSH Metrics | Packet size stats, authentication counts, key exchange | 19 |
| ICMP Metrics | Ping requests, packet count, error count | 6 |
| Connection/State Metrics | Idle connections, establishment/termination attempts, request/response intervals | 10 |
| Error/Retransmission | Invalid checksums, retransmissions, packet loss | 15 |
| Keep-alive | Keep-alive message counts | 2 |
| Destination-specific Metrics | Duplicate metrics tracked separately for destination host | 99 |

TABLE 7: Feature extracted for the prototype.

TABLE 8: Attack keywords used in search and number of repositories found

| Attack Name | Base Keywords | Variations | Repos |
|---|---|---|---|
| TCP DoS | TCP | Flood, DoS, Denial of Service, Attack | 2333 |
| HTTP DoS | HTTP, HTTP GET, HTTP POST, HTTP Request | Flood, DoS, Denial of Service, Attack, Bombardment, Overload | 2935 |
| SSH Brute Force | SSH Dictionary, SSH Brute Force, SSH Brute-Force, SSH Password Guessing, SSH Password Cracking, SSH, SSH Login Guessing, SSH Password Spraying | Attack, Technique, Method, Tool, Hack, Crack, Attempt, Threat, Vulnerability, Exploit, Breach, Brute Force, Dictionary Attack | 2585 |

**Algorithm 2** Online Algorithm to Extract Relevant Feature Values for Packet at Time $t$

1: **function** EXTRACTFEATURES(packet $p$)
2:     $H_1, H_2 \leftarrow p.identifiers$    ▷ $H_1$: source identifier, $H_2$: destination identifier
3:     $V_{H_1,H_2} \leftarrow V_{channel}(H_1, H_2)$
4:     $V_{H_2} \leftarrow V_{destination}(H_2)$
5:     $\Delta V \leftarrow [\,]$
6:     **for all** $v$ in $V_{H_1,H_2}$ **do**
7:         $x_i \leftarrow p.f_i$ ▷ Get corresponding value to update
8:         $W_{t-x}, CS_{t-x}, r_{t-x}, SSR_{t-x}, \sigma_{t-x} \leftarrow v$
9:         $\Delta CS \leftarrow CS_{t-x}$        ▷ Store value at previous timestep
10:         $W_t, CS_t, r_t, SSR_t, \sigma_t$           $\leftarrow$        UP-DATE$(CS_{t-x}, W_{t-x}, SSR_{t-x}, \sigma_{t-x}, x_i)$
11:         $v \leftarrow W_t, CS_t, r_t, SSR_t, \sigma_t$
12:         $\Delta CS \leftarrow |CS_t - \Delta CS|$
13:         $\Delta V \leftarrow \Delta V \cup \{\Delta CS\}$
14:     **end for**
15:     **for all** $(i, v)$ in ENUMERATE$(V_{H_2})$ **do**
16:         $x_i \leftarrow \Delta V[i]$              ▷ Get change in value
17:         $W_{t-x}, CS_{t-x}, r_{t-x}, SSR_{t-x}, \sigma_{t-x} \leftarrow v$
18:         $W_t, CS_t, r_t, SSR_t, \sigma_t$           $\leftarrow$        UP-DATE$(CS_{t-x}, W_{t-x}, SSR_{t-x}, \sigma_{t-x}, x_i)$
19:         $v \leftarrow W_t, CS_t, r_t, SSR_t, \sigma_t$    ▷ Store updated statistics
20:     **end for**
21:     **return** $V_{H_1,H_2} \| V_{H_2}$      ▷ Return updated feature values
22: **end function**