

# Robust Anomaly Detection in Network Traffic: Evaluating Machine Learning Models on CICIDS2017

<sup>1st</sup> Zhaoyang Xu

Department of Electrical and Computer Engineering  
University of Southern California  
Los Angeles, USA  
xuzhaoya@usc.edu

<sup>2nd</sup> Yunbo Liu \*

Department of Electrical and Computer Engineering  
Duke University  
Durham, USA  
yunbo.liu954@duke.edu \*

**Abstract**—Identifying suitable machine learning paradigms for intrusion detection remains critical for building effective and generalizable security solutions. In this study, we present a controlled comparison of four representative models—Multi-Layer Perceptron (MLP), 1D Convolutional Neural Network (CNN), One-Class Support Vector Machine (OCSVM) and Local Outlier Factor (LOF)—on the CICIDS2017 dataset under two scenarios: detecting known attack types and generalizing to previously unseen threats. Our results show that supervised MLP and CNN achieve near-perfect accuracy on familiar attacks but suffer drastic recall drops on novel attacks. Unsupervised LOF attains moderate overall accuracy and high recall on unknown threats at the cost of elevated false alarms, while boundary-based OCSVM balances precision and recall best, demonstrating robust detection across both scenarios. These findings offer practical guidance for selecting IDS models in dynamic network environments.

**Keywords**—Intrusion Detection System (IDS), Anomaly Detection, unsupervised Learning, CICIDS2017, Network Security, Machine Learning

## I. INTRODUCTION

Following the global COVID-19 pandemic, people have become increasingly dependent on network infrastructure. With the widespread adoption of smart devices and the Internet of Things (IoT), a great volume of data are generated every day. Despite the significant convenience brought about by IoT devices, their widespread adoption inevitably leads to the transmission of vast amounts of sensitive information across networks. A substantial portion of this data—including personal, financial [1], and operational records—is now stored directly on cloud platforms such as AWS, Azure, and Google Cloud. This growing reliance on cloud-based infrastructure raises a critical question: how can we ensure the security and privacy of this data as it flows through increasingly complex and distributed networks [2], [26].

To address these security concerns, Intrusion Detection Systems (IDS) have been widely adopted as a critical component of modern network defense strategies [3]. IDS can be categorized to 2 types: Anomaly-based Intrusion Detection Systems (A-IDS) and Signature-based Intrusion Detection Systems (S-IDS).

In signature-based intrusion detection (S-IDS), the system monitors user or network activity and compares it against a database of known attack patterns, also referred to as signatures. If a match is found between the observed behavior and any of the stored signatures, the activity is flagged as an attack. However, because it can only detect attacks that are already known and recorded, if there is a novel attack, S-IDS unable to identify it. In contrast, Anomaly-based Intrusion Detection Systems (A-IDS) operate by establishing a baseline of normal behavior within a network or system, this capability to detect novel threats is well-supported in the anomaly-detection literature [4]. Any significant deviation from this baseline is treated as a potential intrusion. This makes the system has the ability to identify novel attacks.

In this work, we implement an anomaly-based intrusion detection pipeline that integrates multiple learning-based models, including MLP, OCSVM, and LOF. The system is evaluated using the CICIDS2017 dataset to examine its effectiveness in identifying malicious traffic under realistic network conditions [5].

The remainder of this paper is organized as follows. Section II reviews related works. Section III describes the explanation of the method adopted for this study. Section IV presents the evaluation of results. Section V is discussion and future work.

## II. LITERATURE REVIEW

Mashaly et al [5] evaluated the CICIDS2017 dataset using KNN, enhanced KNN, and LOF, showing that semi-supervised approaches trained only on normal data could effectively detect anomalies, especially with PCA-based feature selection.

Wang and Yang [6] proposed an intelligent cloud security framework combining deep learning models like CNN and LSTM with reinforcement learning, achieving 97.3% detection accuracy on real-time traffic.

Kale et al [7] developed a hybrid deep learning anomaly detection framework integrating spatial-temporal analysis, validated on benchmark datasets. Similarly, Tossou et al [8] applied deep learning to anomaly-based intrusion detection and reported improved generalization to unknown attack types.

Abrar et al [9] used classical ML classifiers on the NSL-KDD dataset, achieving robust performance through optimized feature selection, confirming the continued relevance of traditional models.

Tavallae et al. [10] presented a comprehensive analysis of the KDD Cup 99 dataset, identifying key preprocessing steps and attack categorization strategies that have become standard in IDS research.

Hizal et al [11] proposed a deep learning-based IDS specifically designed for cloud environments, incorporating real-time analytics and demonstrating strong results on CICIDS2017.

In a study addressing data imbalance in security applications, Wang et al [12] compared classical and deep architectures, highlighting the risk of overfitting in purely supervised models under skewed data distributions.

Additionally, Moustafa and Slay [2] and Chandola et al [4] laid the groundwork for anomaly detection benchmarks, emphasizing the importance of statistical diversity and outlier behavior in datasets like UNSW-NB15 and KDD99.

Together, these studies demonstrate the evolution from signature-based IDS toward hybrid and anomaly-based methods. However, the trade-offs between generalization, accuracy, and interpretability remain open challenges—particularly in the context of unseen or evolving attack patterns. Similar challenges have also been observed in risk modeling tasks beyond cybersecurity, such as credit risk assessment using structural graph-based models [21].

### III. METHODOLOGY

In this paper, we used CICIDS2017 data set. This dataset was created by the Canadian Institute for Cybersecurity and provides a comprehensive collection of network traffic that reflects real-world scenarios [5], including both benign and malicious activities. It contains over 80 features extracted from packet flows and covers various attack types such as DoS, DDoS, brute force, infiltration, botnet, and web-based attacks. The dataset is widely used for evaluating intrusion detection systems due to its diversity and realistic traffic patterns.

We selected four representative models: Multi-Layer Perceptron (MLP), 1D Convolutional Neural Network (CNN), One-Class SVM (OCSVM) and Local Outlier Factor (LOF). We compared their performance to evaluate their effectiveness in anomaly-based intrusion detection.

#### A. Data Processing

CICIDS2017 contains 2,830,743 records in total, consisting of both benign and malicious traffic across 15 labeled categories. The number of samples in each category is shown in Table 1. To evaluate model performance under both known and unknown attack scenarios, we divided the dataset into three subsets: Training Set, Overall Test Set, and Unknown Attack Test Set.

In our design, we selected three specific attack types—DoS slowloris, DoS Slowhttptest, and Bot—as unknown attacks. These attacks were completely excluded from all training sets

TABLE I: Attack type count

Label	Count
BENIGN	2 270 397
DoS Hulk	231 073
PortScan	158 930
DDoS	128 027
DoS GoldenEye	10 293
FTP-Patator	7 938
SSH-Patator	5 897
DoS slowloris	5 796
DoS Slowhttptest	5 499
Bot	1 966
Web Attack – Brute Force	1 507
Web Attack – XSS	652
Infiltration	36
Web Attack – SQL Injection	21
Heartbleed	11

to ensure that their patterns were not learned in advance. The selection was guided by two main considerations:

- To evaluate how well the models generalize to unknown attacks, simulating realistic deployment scenarios where not all attack types are known beforehand.
- The sample sizes of these attacks are moderate—not too large to compromise the training set quality, and not too small to produce unstable or unreliable test results.

The following is how we made the training set and test set. Training set:

- For MLP and CNN, the training set includes 80% of benign samples and 11 attack types except for three unknown attacks: DoS slowloris, DoS Slowhttptest, and Bot.
- For OCSVM and LOF, the training set includes only 80% of benign traffic [5].

Test set:

- Overall-Test-Set: This set includes 20% of benign samples and 11 attack types plus all three known attack types. It is used to evaluate general detection performance across both seen and unseen attacks.
- Unknown Attack Test Set: This set includes only three excluded unknown attack types and equal number of benign traffic.

All random sampling operations were performed with `random_state=42` to ensure consistent and reproducible training and test splits.

To ensure fair comparison and stable training, we applied `scikit-learn's StandardScaler`: we fit and transformed the training set, then used the same fitted scaler to transform both test splits. This procedure centers each feature to zero mean and scales it to unit variance, improving model convergence and comparability across different feature distributions.

#### B. Implementation Environment

All experiments were conducted on AWS EC2:

- **Instance type:** `g4dn.xlarge` (1× NVIDIA T4 GPU, 4 vCPUs, 16 GiB RAM)
- **OS:** Ubuntu 20.04 LTS

- **Python:** 3.8.10
- **scikit-learn:** 1.0.2
- **TensorFlow:** 2.6.0
- **pandas:** 1.3.4, **numpy:** 1.21.2
- **Other libraries:** matplotlib 3.4.3

### C. Classifiers

1) *Multi-Layer Perceptron (MLP)*: A Multi-Layer Perceptron (MLP) is a fully connected feedforward neural network composed of an input layer, one or more hidden layers, and a single-node output layer. Each hidden unit computes a weighted sum of its inputs followed by a nonlinear activation:

$$h^{(l)} = \sigma \left( W^{(l)} h^{(l-1)} + b^{(l)} \right)$$

where  $h^{(l)}$  is the activation of the  $l$ -th layer,  $W^{(l)}$  and  $b^{(l)}$  are the weight matrix and bias vector, and  $\sigma(\cdot)$  denotes a nonlinear activation function such as ReLU.

For binary classification, the output layer uses the sigmoid function:

$$\hat{y} = \frac{1}{1 + e^{-z}}$$

where  $z$  is the output of the final layer before activation. The network is trained using the binary cross-entropy loss:

$$\mathcal{L}_{\text{BCE}} = -[y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})]$$

where  $y$  is the true label (0 for benign, 1 for malicious), and  $\hat{y}$  is the predicted probability.

To train this network effectively, we configured the following hyperparameters for the MLP:

- Hidden layers: (100, 50) neurons
- Activation: ReLU
- Optimizer: Adam, learning rate = 0.001
- Batch size: 256
- Max epochs: 100 with early stopping (patience = 5)
- Random seed: 42

2) *1D Convolutional Neural Network (CNN)*: benign temporal and spatial patterns within the network flow features, we employ a one-dimensional convolutional neural network (1D-CNN) composed of multiple convolutional and pooling stages [20] followed by a fully connected classifier. Concretely, each convolution layer computes

$$(x * w)[t] = \sum_{k=0}^{K-1} x[t+k] w[k] + b,$$

where  $x$  is the input sequence,  $w$  is the convolution kernel of size  $K$ , and  $b$  is a bias term. The convolution output is passed through a nonlinear activation function, here chosen as ReLU:

$$\text{ReLU}(z) = \max(0, z).$$

After each convolution we apply max-pooling over a fixed window to reduce resolution and introduce translation invariance:

$$y[t] = \max_{0 \leq k < P} z[t \cdot S + k],$$

where  $P$  is the pooling size and  $S$  the stride. We also include dropout layers to mitigate overfitting by randomly zeroing a fraction of activations.

Finally, the flattened feature map is fed into one or more dense layers terminating in a sigmoid output for binary classification. To train this network effectively, we used the following configuration:

- Optimizer: Adam, learning rate = 0.001
- Loss: Focal Loss ( $\gamma = 2.0$ ,  $\alpha = 0.25$ )
- Batch size: 512, validation split = 0.2
- Epochs: 50 with early stopping (patience = 3)
- Class weights: benign=1.0, malicious=5.0
- Random seed: 42

3) *One-Class Support Vector Machine (OCSVM)*: OCSVM method is an extension of Support Vector Machine (SVM) method [7]. It is a method suitable for handling unlabeled data. It attempts to learn a decision boundary that encloses the majority of the normal data points, treating anything outside the boundary as anomalous. This property makes OCSVM particularly suitable for real-world intrusion detection scenarios where attack data is scarce or unknown.

$$\min_{w, \rho, \xi} \frac{1}{2} \|w\|^2 + \frac{1}{\nu n} \sum_{i=1}^n \xi_i - \rho$$

subject to:

$$(w \cdot \phi(x_i)) \geq \rho - \xi_i, \quad \xi_i \geq 0$$

where

- $\phi(x)$  maps the input into a high-dimensional kernel space,
- $\nu \in (0, 1]$  controls the trade-off between margin size and the number of training errors,
- $\rho$  defines the decision boundary,
- $\xi_i$  are slack variables for soft margins.

The decision function is the following.

$$f(x) = \text{sign} \left( \sum_{i=1}^n \alpha_i K(x_i, x) - \rho \right)$$

where  $K(x_i, x) = \langle \phi(x_i), \phi(x) \rangle$  (commonly RBF kernel). Samples with  $f(x) < 0$  are considered anomalies. To configure the OCSVM for our experiments, we set:

- Kernel: RBF
- $\nu = 0.05$
- Gamma: scale
- Random seed: 42

4) *Local Outlier Factor (LOF)*: The Local Outlier Factor algorithm identifies points with significantly lower density than their neighbors. It compares the local density of a data point to those of its neighbors to determine how isolated the point is. LOF is especially effective for detecting local anomalies in datasets with varying density, making it a useful baseline

in IDS settings where malicious traffic often deviates subtly from normal patterns.

$$\text{lrd}_k(x) = \left( \frac{1}{|N_k(x)|} \sum_{y \in N_k(x)} \max\{\text{dist}(x, y), \text{reach-dist}_k(y)\} \right)^{-1}$$

where  $N_k(x)$  is the set of  $k$ -nearest neighbors of  $x$ , and  $\text{reach-dist}_k(y) = \max\{k\text{-dist}(y), \text{dist}(x, y)\}$ .

The LOF score is given by:

$$\text{LOF}_k(x) = \frac{1}{|N_k(x)|} \sum_{y \in N_k(x)} \frac{\text{lrd}_k(y)}{\text{lrd}_k(x)}$$

A LOF score close to 1 indicates that  $x$  has similar density to its neighbors, while a substantially larger value suggests an outlier. For our experiments, we configured LOF as follows:

- Number of neighbors ( $k$ ): 80
- Novelty detection: True
- Distance metric: Euclidean ( $\text{metric} = \text{minkowski}, p = 2$ )
- Leaf size: 80
- Random seed: 42

#### D. Model Selection Justification

The selected models represent four distinct methodological approaches:

- MLP: Standard supervised deep learning
- CNN: Supervised deep learning with convolutional feature extraction
- OCSVM: Classical boundary-based unsupervised learning
- LOF: Established density-based unsupervised technique

This controlled comparison isolates the effects of supervision paradigm, following the experimental design principles in [13].

#### E. Analysis

To evaluate the performance of each classifier, standard binary classification metrics derived from the confusion matrix (True Positives: TP, False Positives: FP, True Negatives: TN, and False Negatives: FN) are used [10]. The following metrics are computed:

- **Accuracy:**

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

- **Precision:**

$$\text{Precision} = \frac{TP}{TP + FP}$$

- **Recall:**

$$\text{Recall} = \frac{TP}{TP + FN}$$

- **F1-Score:**

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

In our analysis, class ‘0’ refers to benign traffic and class ‘1’ refers to malicious traffic. These metrics provide a clear

view of each model’s ability to distinguish attacks from normal traffic under binary classification.

## IV. RESULTS

This section compares the detection capabilities of MLP, CNN, LOF and OCSVM under two scenarios:

- (1) Overall Test Set, which includes all traffic types
- (2) Unknown Attack Test Set, containing only truly novel attack types.

Tables II and III present the key evaluation metrics. We then delve into confusion matrices to highlight each model’s behavior differences, especially their handling of unseen attacks.

TABLE II: Performance on Overall Test Set

Model	Accuracy	Precision	Recall	F1-score
MLP	0.9775	0.9894	0.9036	0.9446
CNN	0.9650	0.9316	0.9010	0.9160
LOF	0.8046	0.9106	0.8341	0.8706
OCSVM	0.8356	0.6525	0.4784	0.5520

TABLE III: Performance on Unknown Attack Test Set

Model	Accuracy	Precision	Recall	F1-score
MLP	0.5863	0.9860	0.1750	0.2973
CNN	0.5882	0.9110	0.1954	0.3218
LOF	0.6087	0.5746	0.8370	0.6814
OCSVM	0.7919	0.9072	0.6503	0.7575

TABLE IV: Per-class Accuracy on the Overall Test Set

Attack Type	CNN	OCSVM	MLP	LOF
DoS Hulk	0.9978	0.6886	0.9921	0.7633
BENIGN	0.9822	0.9316	0.9974	0.8341
PortScan	0.9713	0.0095	0.9995	0.5280
DDoS	0.9996	0.6299	0.9993	0.9297
Web Attack – Brute Force	0.1860	0.0100	0.1395	0.1894
FTP-Patator	0.9987	0.0132	0.9987	0.6688
DoS GoldenEye	0.9995	0.7358	0.9932	0.8193
SSH-Patator	0.9864	0.0008	0.9839	0.9915
Infiltration	0.5714	0.8571	0.0000	0.8571
Web Attack – XSS	0.0231	0.0308	0.0308	0.0385
Web Attack – Sql Injection	0.0000	0.0000	0.2500	0.0000
Heartbleed	0.5000	1.0000	1.0000	1.0000
Bot	0.0000	0.0443	0.0000	0.4680
DoS slowloris	0.3675	0.5733	0.3675	0.3176
DoS Slowhttptest	0.0838	0.9480	0.0347	0.4150

Table II presents the performance of the four models on the *Overall Test Set*:

- **MLP** achieves the best results across all metrics on the Overall Test Set. Its high precision (0.9894) and F1-score (0.9446) for the malicious class demonstrate strong capability in detecting known attacks.
- **CNN** closely follows MLP with an overall accuracy of 0.9650, precision of 0.9316, and F1-score of 0.9160, indicating that convolutional feature extraction also excels at capturing attack patterns.
- **LOF** and **OCSVM** obtain lower overall accuracy (0.8046 and 0.8356, respectively) but still exhibit robust anomaly



detection performance, particularly OCSVM's ability to generalize to a variety of normal traffic profiles.

Table III shows the performance on the *Unknown Attack Test Set*, which consists only of the three designated unknown attack types:

- **MLP** shows a highly imbalanced detection behavior: although its precision is high (0.9860), its recall is very low (0.1750), yielding a poor F1-score of 0.2973. This indicates a strong tendency to under-report unknown attacks.
- **CNN** similarly struggles with novel threats, achieving precision of 0.9110 but recall of only 0.1954 ( $F_1 = 0.3218$ ), suggesting that convolutional feature extraction alone still overfits to familiar patterns.
- **LOF** demonstrates a more balanced profile on novel attacks (Accuracy = 0.6087, Precision = 0.5746, Recall = 0.8370,  $F_1$ -score = 0.6814), indicating that density-based anomaly detection can recover a large fraction of unseen intrusions at the cost of moderate false alarms.
- **OCSVM** achieves the best performance in this scenario (Accuracy = 0.7919,  $F_1 = 0.7575$ ), which shows that training exclusively on benign data supports better generalization to unseen malicious behaviors.

Table IV breaks down each model's accuracy by attack type. Several patterns emerge:

- **Supervised models (MLP & CNN)** excel on high-volume attacks and the benign class but fail on rare or subtle threats. Both achieve near-perfect accuracy on PortScan (MLP: 0.9995; CNN: 0.9713), DDoS (MLP: 0.9993; CNN: 0.9996), FTP-Patator (0.9987; 0.9987), and BENIGN (MLP: 0.9974; CNN: 0.9822), reflecting strong learning when ample labeled data is available. However, they both struggle on low-volume Web Attack – Brute Force (MLP: 0.1395; CNN: 0.1860), Web Attack – XSS (MLP: 0.0308; CNN: 0.0231), and Infiltration (MLP: 0.0000; CNN: 0.5714)—indicating overfitting to frequent classes and difficulty generalizing from limited examples. On the three truly novel attacks (DoS slowloris, Slowhttptest, Bot), their accuracy further collapses ( $\leq 0.37$  on slowloris;  $\leq 0.08$  on Slowhttptest; 0.00 on Bot), underscoring their reliance on seen patterns.
- **Density-based LOF** shows notably higher per-class accuracy on certain novel attack types compared with supervised models. On the Unknown Attack Set, LOF attains accuracy of 0.3176 on DoS slowloris (versus  $\leq 0.37$  for MLP/CNN), 0.4150 on Slowhttptest (versus  $\leq 0.08$ ), and 0.4680 on Bot—demonstrating its ability to flag these outliers more effectively.
- **Boundary-based OCSVM** stands out on truly novel attacks, achieving moderate accuracy on DoS slowloris (0.57), Slowhttptest (0.94), while maintaining balanced detection across both known and unknown threats (Overall: Accuracy = 0.8356,  $F_1 = 0.5520$ ; Unknown: Accuracy = 0.7919,  $F_1 = 0.7575$ ). This highlights OCSVM's strength in learning a compact boundary around benign

traffic and flagging deviations—even for completely unseen malicious behaviors.

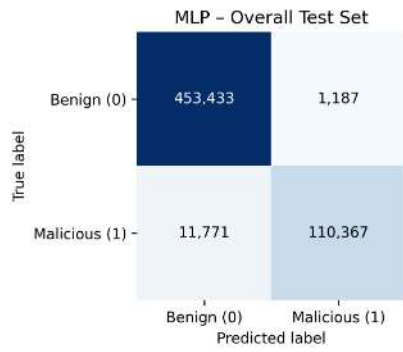
**Confusion Matrix Analysis.** Figure 1 illustrates classifier performance under known attack conditions:

- **MLP** achieves nearly perfect separation, with 453,433 true negatives and 110,367 true positives, suffering only 1,187 false positives and 11,771 false negatives. Such results reflect its high discriminative power when trained on the full range of attack types.
- **CNN** likewise performs strongly, correctly classifying 446,537 benign instances and 110,044 attacks, with 8,083 false positives and 12,094 false negatives—indicating that convolutional feature extraction yields comparable separation to MLP on known threats.
- **LOF** maintains a solid detection ability on known attacks, correctly classifying 379,190 benign instances but incurring 75,430 false positives and missing 37,244 malicious samples—suggesting that density-based anomaly detection still struggles with perfectly separating benign from malicious traffic in familiar settings.
- **OCSVM** balances these trade-offs more effectively: it significantly reduces false positives to 31,113, albeit with 63,711 false negatives, indicating that unsupervised learning can better delineate benign behaviour with fewer mislabeled benign instances.

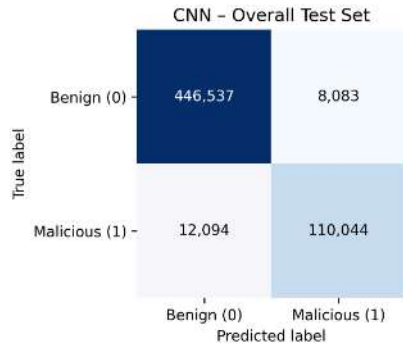
Figure 2 shifts focus to the Unknown Attack Test Set:

- **MLP** performance collapses: though it correctly identifies 13,228 benign flows, it mislabels 10,940 malicious samples and only detects 2,321. This stark contrast highlights its inability to generalize beyond the training distribution.
- **CNN** shows a similar collapse: it correctly labels 13,008 benign samples but misclassifies 10,670 attacks, detecting only 2,591, which underscores its limited generalization to truly novel patterns.
- **LOF** shows marked improvement with novel attacks, correctly excluding 11,100 benign flows and flagging 5,043 malicious samples, while incurring 2,161 false positives and missing 8,218 attacks—demonstrating its better adaptability to unseen threats at the cost of moderate misclassification.
- **OCSVM** delivers the most balanced performance in novel scenarios: with only 882 benign samples misclassified and 4,638 attacks missed, it achieves 8,623 true positives. This underscores its superior recall and precision balance in detecting unseen attack types.

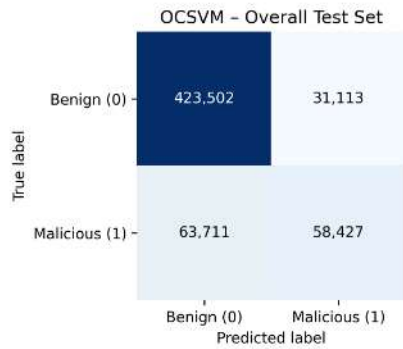
**Summary of Observed Trends:** The confusion matrix results underscore a clear pattern: **MLP and CNN both excel in familiar contexts but severely degrade on new threats**, with CNN showing marginally lower separation power than MLP on known attacks and equally poor generalization to unknown patterns. In contrast, **LOF offers greater resilience albeit with higher error rates**, and **OCSVM combines the strengths of both**, maintaining strong detection rates while reducing both false alarms and misses even when facing



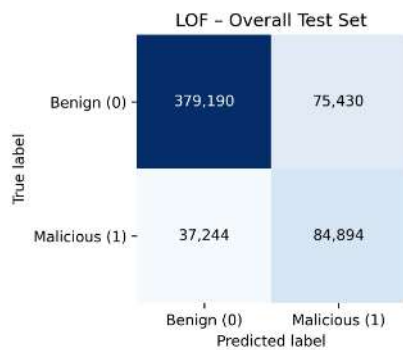
(a) MLP – Overall



(b) CNN – Overall

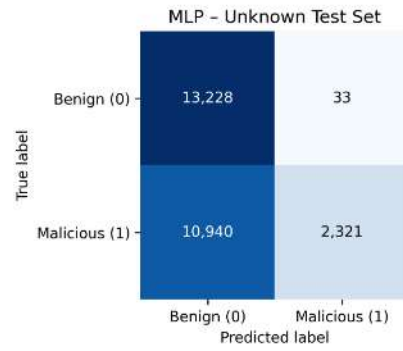


(c) OCSVM – Overall

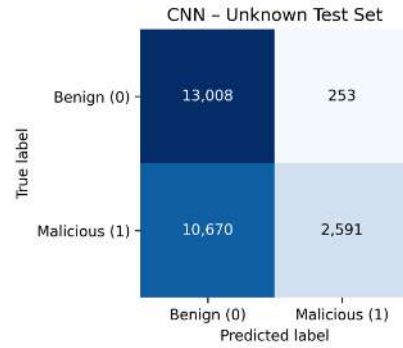


(d) LOF – Overall

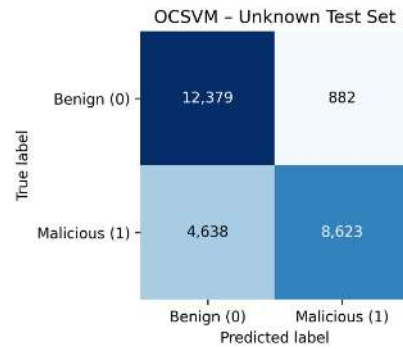
Fig. 1: Confusion matrices on the Overall Test Set. (a)–(d) correspond to MLP, CNN, OCSVM, and LOF, respectively.



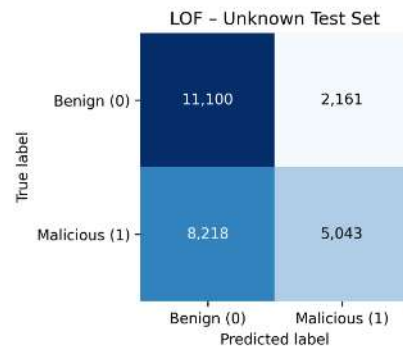
(a) MLP – Unknown



(b) CNN – Unknown



(c) OCSVM – Unknown



(d) LOF – Unknown

Fig. 2: Confusion matrices on the Unknown Attack Test Set. (a)–(d) correspond to MLP, CNN, OCSVM, and LOF, respectively.

unknown attacks—demonstrating why it emerges as the most practical approach for real-world intrusion detection.

## V. DISCUSSION

The paradigm-level differences manifest in three key aspects:

- **MLP** relies on labeled samples of both benign and known attack types during training. Its excellent performance on the Overall Test Set stems from its ability to learn discriminative patterns directly tied to these labels. However, it lacks exposure to unseen attack behaviors and therefore fails to generalize, which explains its drastic drop in recall when tested on novel attacks. Similar generalization challenges have been noted in other systems [14], including concept-level backdoor vulnerabilities in interpretable models [23].
- **CNN** extends supervised learning with convolutional feature extraction, capturing local temporal and spatial correlations in the data. This leads to performance comparable to MLP on known attacks, but like MLP, it suffers from overfitting to familiar patterns and shows similarly poor generalization to truly novel threats.
- **LOF** models normal behavior density without requiring attack labels. It detects anomalies purely based on deviation from established distributions. This grants higher recall on unseen attacks but also induces elevated false positives and false negatives, even for known attacks, reflecting the sensitivity of density thresholds and imperfect boundary estimation.
- **OCSVM** trains solely on benign samples, defining a boundary around normal traffic. It flags out-of-distribution patterns as anomalies. This method strikes an optimal balance: it avoids overfitting to known attacks while maintaining precise decision boundaries, resulting in robust detection across both known and novel attack scenarios [7], [19], [22].

**Conclusions from Model Behavior:** These observations suggest that anomaly detection models defining decision boundaries based only on benign data can be more robust to unseen attacks than supervised models relying on known attack patterns. Among the supervised methods, both **MLP** and **CNN** achieve strong performance on familiar threats but suffer significant drops in recall on novel attacks. Within the unsupervised paradigm, our results show that boundary-based **OCSVM** consistently outperforms density-based **LOF** for IDS tasks: OCSVM’s global decision surface yields fewer false alarms in mixed-density traffic and more stable recall on novel attacks, making it especially well-suited for real-world intrusion detection, echoing broader efforts to build robust and interpretable AI systems in safety-critical domains such as autonomous driving [24].

**Future Work:** Future research could explore hybrid architectures that combine supervised and unsupervised components, leveraging the sequence modeling capability of frameworks like SETransformer [15], [18]. Another promising direction

involves post-training model adaptation [16], where attribute unlearning techniques [14], [17], [25] allow intrusion detectors to incorporate new patterns without retraining from scratch.

## REFERENCES

- [1] Y. Chen, C. Zhao, Y. Xu, and C. Nie, “Year-over-year developments in financial fraud detection via deep learning: A systematic literature review,” 2025. [Online]. Available: <https://arxiv.org/abs/2502.00201>.
- [2] N. Moustafa and J. Slay, “The evaluation of network anomaly detection systems: Statistical analysis of the UNSW-NB15 data set and the comparison with the KDD99 data set,” *Information Security Journal: A Global Perspective*, vol. 25, no. 1-3, pp. 18–31, 2016. [Online]. Available: <https://doi.org/10.1080/19393555.2015.1125974>.
- [3] R. Sommer and V. Paxson, “Outside the closed world: On using machine learning for network intrusion detection,” in *Proceedings of the IEEE Symposium on Security and Privacy*, 2010, pp. 305–316.
- [4] V. Chandola, A. Banerjee, and V. Kumar, “Anomaly detection: A survey,” *ACM Computing Surveys*, vol. 41, no. 3, p. Article 15, 2009. [Online]. Available: <https://doi.org/10.1145/1541880.1541882>.
- [5] T. Elmasri, N. Samir, M. Mashaly, and Y. Atef, “Evaluation of CI-CIDS2017 with qualitative comparison of machine learning algorithm,” in *2020 IEEE Cloud Summit*, 2020, pp. 46–51.
- [6] Y. Wang and X. Yang, “Research on enhancing cloud computing network security using artificial intelligence algorithms,” 2025. [Online]. Available: <https://arxiv.org/abs/2502.17801>.
- [7] R. Kale, Z. Lu, K. W. Fok, and V. L. L. Thing, “A hybrid deep learning anomaly detection framework for intrusion detection,” in *2022 IEEE 8th Intl Conference on Big Data Security on Cloud (BigDataSecurity), IEEE Intl Conference on High Performance and Smart Computing (HPSC), and IEEE Intl Conference on Intelligent Data and Security (IDS)*, 2022, pp. 137–142.
- [8] S. Tossou, M. Qorib, and T. Kacem, “Anomaly based intrusion detection system: A deep learning approach,” in *2023 International Symposium on Networks, Computers and Communications (ISNCC)*, 2023, pp. 1–6.
- [9] I. Abrar, Z. Ayub, F. Masoodi, and A. M. Bamhdi, “A machine learning approach for intrusion detection system on NSL-KDD dataset,” in *2020 International Conference on Smart Electronics and Communication (ICOSEC)*, 2020, pp. 919–924.
- [10] M. Tavallaei, E. Bagheri, W. Lu, and A. A. Ghorbani, “A detailed analysis of the KDD Cup 99 data set,” in *2009 IEEE Symposium on Computational Intelligence for Security and Defense Applications*, 2009, pp. 1–6.
- [11] S. Hizal, Ü. Çavuşoğlu, and D. Akgün, “A new deep learning based intrusion detection system for cloud security,” in *2021 3rd International Congress on Human-Computer Interaction, Optimization and Robotic Applications (HORA)*, 2021, pp. 1–4.
- [12] C. Wang, C. Nie, and Y. Liu, “Evaluating supervised learning models for fraud detection: A comparative study of classical and deep architectures on imbalanced transaction data,” 2025. [Online]. Available: <https://arxiv.org/abs/2505.22521>.
- [13] B. Schölkopf, J. C. Platt, J. Shawe-Taylor, A. J. Smola, and R. C. Williamson, “Estimating the support of a high-dimensional distribution,” *Neural Computation*, vol. 13, no. 7, pp. 1443–1471, 2001.
- [14] Y. Li, C. Chen, Y. Zhang, W. Liu, L. Lyu, X. Zheng, D. Meng, and J. Wang, “UltraRE: Enhancing RecEraser for recommendation unlearning via error decomposition,” in *Advances in Neural Information Processing Systems*, vol. 36, A. Oh et al., Eds. Curran Associates, Inc., 2023, pp. 12 611–12 625.
- [15] Y. Liu, X. Qin, Y. Gao, X. Li, and C. Feng, “SETransformer: A hybrid attention-based architecture for robust human activity recognition,” 2025. [Online]. Available: <https://arxiv.org/abs/2505.19369>.
- [16] Z. Li and Z. Ke, “Pruning visual concepts for efficient and interpretable transfer learning,” in *Second Workshop on Visual Concepts*, 2025. [Online]. Available: <https://openreview.net/forum?id=ADXwXt0lFt>.
- [17] Y. Li, C. Chen, X. Zheng, Y. Zhang, Z. Han, D. Meng, and J. Wang, “Making users indistinguishable: Attribute-wise unlearning in recommender systems,” in *Proceedings of the 31st ACM International Conference on Multimedia (MM ’23)*, ACM, Oct. 2023, pp. 984–994. [Online]. Available: <http://dx.doi.org/10.1145/3581783.3612418>.
- [18] C. Chen, Y. Zhang, Y. Li, J. Wang, L. Qi, X. Xu, X. Zheng, and J. Yin, “Post-training attribute unlearning in recommender systems,” 2024. [Online]. Available: <https://arxiv.org/abs/2403.06737>.

- [19] S. Sun, J. Liu, W. Liu, and T. Jian, "Robust detection of distributed targets based on Rao test and Wald test," *Signal Processing*, vol. 180, pp. 107801, 2021. [Online]. Available: <https://doi.org/10.1016/j.sigpro.2020.107801>
- [20] Q. Huang, Z. Chen, Z. Li, C. Wang, X. Song, Y. Hu, and L. Nie, "MEDIAN: Adaptive Intermediate-grained Aggregation Network for Composited Image Retrieval," in *Proc. IEEE ICASSP*, 2025, pp. 1–5.
- [21] Z. Zhang, Q. Shen, Z. Hu, Q. Liu, and H. Shen, "Credit risk analysis for SMEs using graph neural networks in supply chain," *arXiv preprint arXiv:2507.07854*, 2025. [Online]. Available: <https://arxiv.org/abs/2507.07854>
- [22] Y. He, X. Wang, and T. Shi, "DDPM-MoCo: Advancing industrial surface defect generation and detection with generative and contrastive learning," in *International Joint Conference on Artificial Intelligence (IJCAI)*, 2024, pp. 34–49.
- [23] S. Lai, Y. Huang, J. Yang, G. Huang, W. Chen, and Y. Yue, "Guarding the Gate: ConceptGuard Battles Concept-Level Backdoors in Concept Bottleneck Models," *arXiv preprint arXiv:2411.16512*, 2024. [Online]. Available: <https://arxiv.org/abs/2411.16512>
- [24] S. Lai, T. Xue, H. Xiao, L. Hu, J. Wu, N. Feng, R. Guan, H. Liao, Z. Li, and Y. Yue, "Drive: Dependable Robust Interpretable Visionary Ensemble Framework in Autonomous Driving," *arXiv preprint arXiv:2409.10330*, 2024. [Online]. Available: <https://arxiv.org/abs/2409.10330>
- [25] S. Sun, Y. Ren, C. Ma, and X. Zhang, "Large language models as topological structure enhancers for text-attributed graphs," *arXiv preprint arXiv:2311.14324*, 2023. [Online]. Available: <https://arxiv.org/abs/2311.14324>
- [26] J. Song, K. Ding, R. Cheng, X. Zhao, X. Luo, Y. Liu, Y. Tian, and Z. Duan, "Research on brand strategy of hotel enterprises—taking Hyatt Hotel Group as an example," *Open Journal of Business and Management*, vol. 13, no. 2, pp. 861–869, 2025.