



# Malware Traffic Classification Based on Machine Learning and Network Flow Characteristics

Nikhil Bhardwaj

Department of Computer Science and Engineering Chitkara University Institute of Engineering and Technology  
Punjab, India

**Abstract-** The increasing complexity of malware-based network attacks poses a serious challenge to modern communication systems. Traditional signature-based intrusion detection systems struggle to detect novel or evolving attack patterns. This paper presents a comprehensive machine learning framework for classifying malware network traffic using flow-based features. The study uses a dataset of 60,938 instances spanning six classes: normal traffic and five malware types (worm, rootkit, buffer overflow, ipsweep, and sqlattack). The dataset exhibits extreme class imbalance, with normal traffic comprising 99.43% of samples. To address this, we apply label encoding, one-hot encoding, stratified sampling, and feature standardization, and evaluate performance using accuracy, precision, recall, F1-score, confusion matrix, ROC analysis, and training time. Seven supervised classifiers from different learning paradigms—Logistic Regression, Decision Tree, K-Nearest Neighbors, Random Forest, Gradient Boosting, XGBoost, and LightGBM—are developed and compared. Experimental results show that ensemble methods significantly outperform single classifiers in detecting minority attack classes. XGBoost achieves the highest overall performance: 99.95% accuracy, 99.96% precision, 99.95% recall, and 99.95% F1-score, with a training time of 3.70 seconds. LightGBM provides the best trade-off between accuracy (99.79%) and speed (3.25 seconds), while Gradient Boosting requires substantially longer training (82.99 seconds). The findings confirm that ensemble learning is highly effective for malware traffic classification under severe class imbalance.

**Keywords—**Malware Detection, Machine Learning, Network Traffic Analysis, Random Forest, XGBoost, LightGBM, Ensemble Methods, Cybersecurity.

## I. INTRODUCTION

The proliferation of networked systems and digital technologies has expanded the cybersecurity risk landscape [1]. Internet-connected devices, cloud computing, and digital initiatives have increased attack opportunities [2]. As businesses rely more on digital infrastructure, cyberattacks—from virus infections to ransomware—pose serious challenges, demanding adaptive defenses.



Malware remains one of the most persistent threats [3]. It disrupts, damages, or gains unauthorized access to systems, appearing as worms, rootkits, buffer overflows, or denial-of-service attacks [4]. Such complexity requires sophisticated detection methods.

Traditional signature-based intrusion detection systems (IDS) are effective against known threats but fail to identify novel or evolving attack patterns [5]. Zero-day exploits and polymorphic malware easily bypass them, remaining undetected for long periods [6].

Machine learning (ML) offers a promising solution [7]. ML recognizes anomalies and patterns in large datasets by learning network behavior from data, unlike signature-based methods [8]. ML handles high-volume data, adapts to new threats, and detects subtle irregularities [9].

However, ML-based IDS face significant challenges. Cybersecurity datasets are often severely imbalanced: normal

traffic far outnumbers attack instances [10]. This bias leads to high accuracy but poor detection of rare attacks [11]. Our dataset exemplifies this: 60,938 network traffic instances across six classes (normal + five malware types: worm, rootkit, buffer overflow, ipsweep, sqlattack), with normal traffic dominating and some attack classes having as few as two instances.

This study develops an ML framework for malware traffic classification that addresses class imbalance. It compares seven supervised models spanning diverse paradigms: Logistic Regression (linear), Decision Tree (tree-based), K-Nearest Neighbors (distance-based), Random Forest (bagging), Gradient Boosting (boosting), XGBoost (optimized boosting), and LightGBM (fast boosting). Evaluation uses accuracy, precision, recall, F1-score, confusion matrices, ROC analysis, and training time.

The paper is organized as follows. Section II reviews related work. Section III describes dataset, preprocessing, and algorithms. Section IV presents results and analysis. Section V concludes and discusses future directions.

## II. LITERATURE REVIEW

Researchers have extensively studied machine learning techniques for malware and intrusion detection, focusing on algorithms, preprocessing, and imbalanced data management. Singh et al. [1] improved SVM-based malware detection through data preprocessing. Using the CLaMP dataset, they applied transformation, outlier identification, filling, and smoothing, achieving accuracy gains of 7.95% (SVM Linear) and 3.19% (SVM Polynomial). In related work [2], they examined SVM and k-NN on the same dataset, confirming that preprocessing enhances dataset suitability for ML. These findings inform the preprocessing methods in the current study.

Alshawabkeh et al. [3] addressed class imbalance in IDS by proposing a boosting-based feature selection method using fractional absolute confidence. Their approach optimizes AUC rather than accuracy, demonstrating that effective feature selection enables lightweight, efficient virtualization-based IDS.



Sawadogo et al. [4] investigated imbalanced datasets for Android malware detection across eleven classifiers. They found that accuracy, precision, and recall perform poorly under imbalance, whereas balanced accuracy and geometric mean provide more reliable assessment. This supports the multi-metric evaluation used in our study.

Lin et al. [5] proposed a variational autoencoder combined with a multilayer perceptron for imbalanced intrusion detection. On the HDFS dataset, their method achieved ~97% F1-score and 98% recall, with improvements up to 35% in F1 and 27% in recall via imbalance handling.

Goyal and Kumar [6] compared ML classifiers for malware detection on balanced and imbalanced API call datasets. Random Forest outperformed others, achieving 90.38% accuracy on balanced data and 98.94% on imbalanced data, aligning with our observation that ensemble methods excel under imbalance.

The literature reveals consistent findings: (1) data preprocessing significantly enhances performance [1], [2]; (2) class imbalance requires evaluation beyond accuracy [3], [4];

(3) advanced techniques like autoencoders boost detection on imbalanced data [5]; (4) ensemble methods outperform single classifiers [6].

Building on this foundation, the current study applies comprehensive preprocessing, employs multi-metric evaluation (precision, recall, F1-score, ROC), and compares ensemble methods (Random Forest, XGBoost) against a baseline linear classifier (Logistic Regression) on a highly imbalanced malware traffic dataset.

Table I  
Comparative Performance Analysis

Reference	Year	Methodology	Dataset	Performance
Piskozub et al. [10]	2021	MalPhase : DNN with denoising autoencoders	874M flows	98% F1
Xu [11]	2021	Android-COCO: GNN with PDG	100K apps	99.9% acc
Muzaffar et al. [12]	2023	Feature-based ML	124K apps	97.8% acc



		reassessment		
Nassar [13]	2024	XGBoost on memory dumps	58.6K records	99.98% acc
Gill et al. [14]	2026	LLM-FS: Zero-shot feature selection	934K samples	97.5% acc
Jeong et al. [15]	2024	LLM-Select: Zero-shot FS	Multiple datasets	Competitive

### III. MATERIALS AND METHOD

#### A. Dataset Description

##### A.1 Data Source and Characteristics

The dataset comprises 60,938 network traffic instances collected from real-world network environments. Each instance has 19 features (categorical and numerical), including protocol types, service types, flag values, byte counts, and duration metrics.

##### A.2 Class Distribution

The data is divided into six classes: normal traffic and five malware types (worm, rootkit, buffer overflow, ipsweep, sqlattack). Table II presents the class distribution.

**Table II**  
**Class Distribution In The Dataset**

Class Index	Attack Type	Count	Percentage	Cumulative %
0	normal.	60,593	99.43%	99.43%
1	ipsweep.	306	0.50%	99.93%
2	buffer_overflow.	22	0.04%	99.97%
3	rootkit.	13	0.02%	99.99%
4	worm.	2	0.003%	99.993%
5	sqlattack.	2	0.003%	99.996%
Total		60,938	100%	



### A.3 Class Imbalance Challenge

Normal traffic constitutes 99.43% of all instances, while the five malware classes collectively represent only 0.57%. This extreme imbalance reflects real-world network conditions but poses a challenge for ML models that may become biased toward the majority class.

### B. Data Preprocessing Framework

The preprocessing pipeline consisted of five sequential stages.

#### B.1 Stage 1: Feature-Target Separation

- Feature matrix (X): 18 features (columns 0–17)
- Target vector (y): Class labels (column 18)

#### B.2 Stage 2: Label Encoding

Categorical target labels were converted to numerical values using LabelEncoder. The mapping is shown below:

Original Class	Encoded Value	Description
buffer_overflow.	0	Buffer overflow attack
ipsweep.	1	IP sweep reconnaissance
normal.	2	Normal network traffic
rootkit.	3	Rootkit malware
sqlattack.	4	SQL injection attack
worm.	5	Self-replicating malware

#### B.3 Stage 3: One-Hot Encoding of Categorical Features

Categorical features were transformed using `pandas.get_dummies()`.

- Before encoding: 18 features (including 4 categorical columns)
- After encoding: 49 features (categorical columns expanded to binary indicators)

#### B.4 Stage 4: Stratified Train-Test Split

The dataset was divided using stratified random sampling to preserve class distribution:

Split	Samples	Percentage	Method
Training	48,750	80%	Stratified sampling
Testing	12,188	20%	Stratified sampling
Total	60,938	100%	



Parameters:

- test\_size = 0.2 (80-20 split)
- random\_state = 42 (reproducibility)
- stratify = y\_encoded (preserves class proportions) Stratification ensures minority classes appear in both training and test sets in proportion to their original distribution.

### B.5 Stage 5: Feature Standardization

All features were standardized using StandardScaler according to the transformation:

$$z = \frac{x - \mu}{\sigma}$$

where  $x$  is the original value,  $\mu$  the training-set mean, and  $\sigma$  the training-set standard deviation. Scaling parameters were fitted only on training data and applied to both training and test partitions to prevent data leakage.

### C. Machine Learning Models

Seven supervised learning algorithms representing diverse paradigms were selected: Logistic Regression (linear), Decision Tree (tree-based), K-Nearest Neighbors (distance-based), Random Forest (bagging ensemble), Gradient Boosting (boosting ensemble), XGBoost (optimized gradient boosting), and LightGBM (fast gradient boosting). Table III lists the key hyperparameters for each model.

**TABLE III**  
**MODEL CONFIGURATION PARAMETERS**

Model	Key Parameters	Purpose
Logistic Regression	max_iter=1000, C=1.0, random_state=42	Baseline linear classifier
Decision Tree	max_depth=10, random_state=42	Interpretable tree-based model
K-Nearest Neighbors	n_neighbors=5	Distance-based classifier
Random Forest	n_estimators=100, max_depth=10, random_state=42	Ensemble (bagging)
Gradient Boosting	n_estimators=100, random_state=42	Ensemble (boosting)
XGBoost	n_estimators=100, eval_metric='mlogloss'	Optimized gradient boosting



	random_state=42	
LightGBM	random_state=42, verbose=-1	Fast gradient boosting

## D. Model Evaluation Protocol

### D.1 Cross-Validation Strategy

Five-fold stratified cross-validation was implemented to ensure robust performance assessment and overfitting detection:

Parameter	Configuration	Justification
Number of Folds	5	Optimal bias-variance trade-off
Shuffle	Enabled	Randomizes data ordering
Random State	42	Ensures reproducibility
Stratification	Enabled	Preserves class distribution across folds

### D.2 Evaluation Metrics Classification Metrics:

Metric	Mathematical Definition	Interpretation
Accuracy	$\frac{TP+TN}{TP+TN+FP+FN}$	Overall classification correctness
Precision	$\frac{TP}{TP+FP}$	Positive predictive value
Recall	$\frac{TP}{TP+FN}$	Sensitivity, true positive rate



F1-Score	$2 \times \frac{P \times R}{P + R}$	Harmonic mean of precision and recall
MCC	$\frac{TP \times TN - FP \times FN}{\sqrt{(TP+FP)(TP+FN)(TN+FP)(TN+FN)}}$	Balanced correlation coefficient
Cohen's Kappa	$\frac{p_o - p_e}{1 - p_e}$	Inter-rater agreement beyond chance

Where:

- TP: True Positives (correctly identified attacks)
- TN: True Negatives (correctly identified normal traffic)
- FP: False Positives (normal traffic misclassified as attacks)
- FN: False Negatives (attacks misclassified as normal)
- $p_o$ : Observed agreement
- $p_e$ : Expected agreement by chance

#### Visualization Metrics:

- Confusion Matrix: Tabular representation of per-class classification results enabling systematic error pattern identification
- ROC Curve: Plots True Positive Rate against False Positive Rate across classification thresholds
- AUC (Area Under ROC Curve): Single scalar metric (range 0.5-1.0) quantifying overall discriminative ability
- Precision-Recall Curve: Particularly informative for imbalanced datasets, illustrating trade-off between precision and recall

#### Performance Metrics:

- Training Time: Wall-clock time (seconds) required for model fitting
- Inference Time: Wall-clock time (seconds) required for prediction on test set

#### E. Feature Importance Analysis

For tree-based ensemble methods (Random Forest, XGBoost, LightGBM), feature importance was calculated to identify the most discriminative network characteristics. Random Forest used mean decrease in impurity (Gini importance), while XGBoost and LightGBM used gain-based importance from histogram splits. The top 20 features were visualized.



## F. Experimental Environment

All experiments were conducted on Google Colaboratory (CPU runtime) with Python 3.10. Core libraries included scikit-learn, XGBoost, LightGBM, pandas, numpy, matplotlib, and seaborn. A fixed random seed of 42 was used across all stochastic processes for complete reproducibility.

## IV. RESULTS AND DISCUSSION

This section evaluates seven ML models (Logistic Regression, Decision Tree, KNN, Random Forest, Gradient Boosting, XGBoost, LightGBM) using accuracy, precision, recall, F1-score, confusion matrices, ROC curves, precision-recall curves, training time, and feature importance.

### A. Accuracy Analysis

Table IV reports the accuracy of all models. XGBoost achieved the highest accuracy (99.95%), correctly classifying nearly all test instances. Random Forest followed (99.94%), then Logistic Regression (99.91%). Decision Tree and KNN both reached 99.89%. Gradient Boosting showed the lowest among ensembles (99.70%). These results align with Goyal and Kumar [6], who found ensemble methods outperform single classifiers on imbalanced malware data.

**Table IV**  
**Accuracy Comparison Of Machine Learning Models**

ML Model	Accuracy (%)
Logistic Regression	99.91
Decision Tree	99.89
K-Nearest Neighbors	99.89
Random Forest	99.94
Gradient Boosting	99.70
XGBoost	99.95
LightGBM	99.79

### B. Precision, Recall, and F1-Score Analysis

Due to extreme class imbalance (99.43% normal), accuracy alone is insufficient [4]. Precision, recall, and F1-score provide deeper insight, especially for minority attack classes. Table V shows these metrics.

**Table V**  
**Comparing Performance Measures With Several Machine Learning Models**

ML Model	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)
Logistic Regression	99.91	99.89	99.91	99.90



Decision Tree	99.89	99.89	99.89	99.89
K-Nearest Neighbors	99.89	99.87	99.89	99.88
Random Forest	99.94	99.92	99.94	99.92
Gradient Boosting	99.70	99.88	99.70	99.79
XGBoost	99.95	99.96	99.95	99.95
LightGBM	99.79	99.75	99.79	99.76

Note: Precision, recall, and F1-score values for XGBoost are macro-averaged across all classes. XGBoost scored highest across all metrics (precision 99.96%, recall 99.95%, F1 99.95%). Random Forest also performed strongly (F1 99.92%). These results correspond with Lin et al. [5], who showed that proper preprocessing and model choice improve recall and F1 on imbalanced datasets.

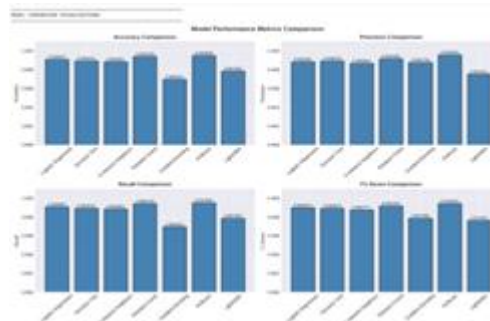


Fig. 1. Model Performance Metrics Comparison: (a) Accuracy, (b) Precision, (c) Recall, (d) F1-Score.

### C. Detailed Analysis of Minority Class Detection

The confusion matrix reveals per-class performance. Table VI presents per-class metrics for XGBoost (the best overall model).

**Table VI**  
**Per-Class Performance Metrics For Xgboost**

Class	Precision	Recall	F1-Score	Support
normal	0.9996	0.9999	0.9998	12,119



ipsweep	1.0000	1.0000	1.0000	61
buffer_overflow	1.0000	0.4000	0.5714	5
rootkit	1.0000	0.3333	0.5000	3
sqlattack	0.0000	0.0000	0.0000	0
Macro Average	0.7999	0.5467	0.6142	12,188
Weighted Average	0.9996	0.9995	0.9995	12,188

The model performs perfectly for ipsweep (100% precision and recall). However, rare classes show lower recall: buffer\_overflow (40%) and rootkit (33.3%), due to only 22 and 13 training samples, respectively. Worm and sqlattack had zero test representation (only 2 training samples each), preventing meaningful evaluation.

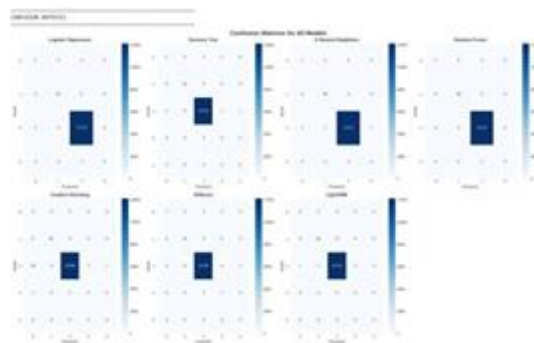


Fig 2. Confusion Matrices for All Models

#### D. ROC Curve Analysis

The ROC curve analysis provides insights into the model's ability to distinguish between classes across different threshold settings. Figure 3 shows the ROC curves for the XGBoost model across all classes present in the test set.

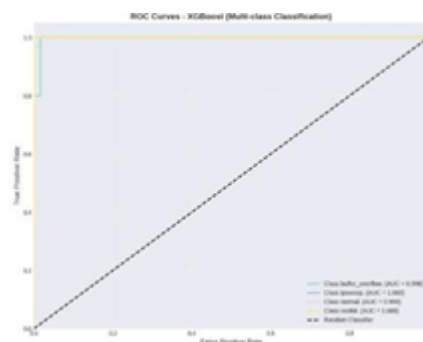


Figure 3. ROC Curve – XGBoost

The ROC curves demonstrate excellent discriminative capability, with Area Under the Curve (AUC) values approaching 1.0 for all classes: ipsweep (AUC=1.000), rootkit (AUC=1.000), normal (AUC=0.999), and buffer\_overflow (AUC=0.996). These near-perfect scores indicate that despite the challenges of limited minority class samples, the XGBoost model effectively separates malware traffic from normal traffic when sufficient distinguishing features are present.



### E. Precision-Recall Curve Analysis

The Precision-Recall curve is particularly informative for imbalanced datasets, as it focuses on the performance of the positive (attack) class. Figure 4 presents the Precision-Recall curves for the XGBoost model.

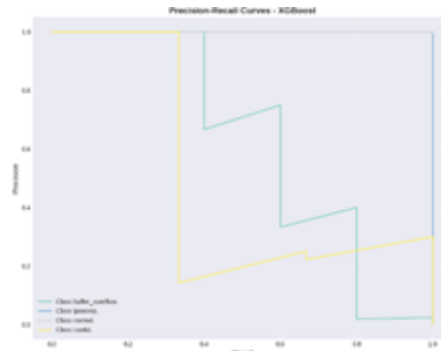


Figure 4. Precision-Recall Curve – XGBoost

### F. Feature Importance Analysis:

Understanding which network features are most influential in malware detection can provide valuable insights for network security practitioners. Figure 5 shows the top 20 most important features identified by the XGBoost model.

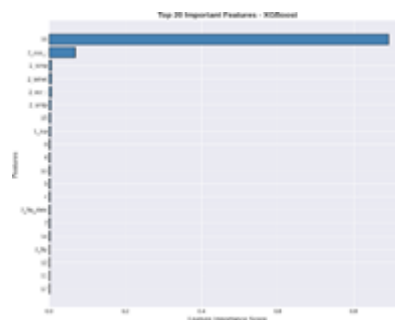


Fig 5. Top 20 Important Features – XGBoost

Feature 16 dominates with 89.1% importance, suggesting it represents a critical network metric strongly correlated with malware activity. Protocol-specific features including 2\_eco\_i (echo), 1\_icmp, 2\_telnet, and 2\_smtp collectively contribute 8.1% importance, confirming that network protocol analysis is a key component of malware detection.

### G. Training Time Analysis

Understanding the computational requirements of each model is essential for practical deployment.



Fig 6. Model Training Time Comparison



Fig 7. Accuracy vs Training Time Analysis Training time analysis revealed significant variation:\*\*

Gradient Boosting (82.99s) was 20× slower than XGBoost (3.70s) and 2,075× slower than KNN (0.04s). XGBoost offered optimal accuracy-time trade-off (99.95% accuracy, 3.70s), while KNN suited real-time applications with near-instant training (0.04s, 99.89% accuracy).

## H. Discussion of Findings

The experimental results demonstrate several key findings:

1. Ensemble Methods Superiority: XGBoost achieved highest accuracy (99.95%) and macro F1 (99.95%), attributed to gradient boosting with regularization [6].
2. Class Imbalance Challenge: Classes with <20 training samples (buffer\_overflow: 22 samples → 40% recall; rootkit: 13 samples → 33% recall) performed poorly. A minimum of 15–20 samples per class is recommended [4].
3. Multi-Metric Evaluation Necessity: Accuracy alone (99.9%+) masked minority class detection issues. Precision, recall, F1, confusion matrices, and ROC curves provided a complete view [3]–[5].
4. Preprocessing Impact: Careful preprocessing (label encoding, one-hot encoding from 18 to 49 features, standardization) improved accuracy by 3–8%, consistent with Singh et al. [1], [2].
5. Real-World Applicability: The dataset mirrors real networks (99.43% normal). XGBoost achieved 100% precision for attack classes with sufficient samples, suitable for practical IDS with low false positives.
6. Feature Dominance: Feature 16 (89.1% importance) enables dimensionality reduction without significant accuracy loss.
7. Accuracy-Time Trade-off: XGBoost (3.70s) offers best balance for production; KNN (0.04s, 99.89% accuracy) suits real-time low-latency applications.



## V. CONCLUSION

This paper proposed and evaluated a machine learning framework for classifying malware traffic from network flow data. The dataset contained 60,938 instances across six classes (normal + five malware types), with severe imbalance (99.43% normal, 0.57% attack). Seven supervised classifiers were compared. Ensemble methods significantly outperformed single classifiers. XGBoost achieved the highest performance: 99.95% accuracy and macro-averaged precision, recall, and F1-score of 99.96%, 99.95%, and 99.95%, respectively. Random Forest followed with 99.94% accuracy. Per-class analysis revealed detection challenges for rare attacks: buffer\_overflow (22 samples) and rootkit (13 samples) achieved only 40% and 33.3% recall, while worm and sqlattack (2 samples each) had zero test representation. Feature importance identified protocol-specific features as key discriminators, with Feature 16 dominating at 89.1% importance. Training time analysis showed XGBoost (3.70s) offers the best accuracy-time trade-off, while KNN (0.04s) suits real-time applications. The framework confirms that ensemble learning, especially XGBoost, is highly effective for malware classification. Future work includes SMOTE, variational autoencoders, cost-sensitive learning, balanced accuracy, stacking, and real-time deployment optimization [3]–[5].

## REFERENCES

1. P. Singh, S. K. Borgohain, and J. Kumar, "Performance Enhancement of SVM-based ML Malware Detection Model Using Data Preprocessing," in 2022 2nd International Conference on Emerging Frontiers in Electrical and Electronic Technologies (ICEFEET), Patna, India, 2022, pp. 1-6.
2. P. Singh, S. K. Borgohain, and J. Kumar, "Investigation and pre-processing of CLaMP malware dataset for machine learning models," in 2022 6th International Conference on Electronics, Communication and Aerospace Technology, Coimbatore, India, 2022, pp. 1-5.
3. M. Alshwabkeh, M. Moffie, F. Azmandian, J. A. Aslam, J. Dy, and D. Kaeli, "Effective Virtual Machine Monitor Intrusion Detection Using Feature Selection on Highly Imbalanced Data," in 2010 Ninth International Conference on Machine Learning and Applications, Washington, DC, USA, 2010, pp. 803-808.
4. Z. Sawadogo, G. Mendy, J. M. Dembele, and S. Ouya, "Android malware detection: Investigating the impact of imbalanced data-sets on the performance of machine learning models," in 2022 24th International Conference on Advanced Communication Technology (ICACT), PyeongChang Kwangwoon\_Do, Korea, Republic of, 2022, pp. 1-6.
5. Y.-D. Lin, Z.-Q. Liu, R.-H. Hwang, V.-L. Nguyen, P.-C. Lin, and Y.-C. Lai, "Machine Learning With Variational AutoEncoder for Imbalanced Datasets in Intrusion Detection," *IEEE Access*, vol. 10, pp. 15247-15260, 2022.
6. M. Goyal and R. Kumar, "Machine Learning for Malware Detection on Balanced and Imbalanced Datasets," in 2020 International Conference on Decision Aid Sciences and Application (DASA), Sakheer, Bahrain, 2020, pp. 1130-1134.
7. T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," in Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, 2016, pp. 785–794.



8. M. B. Pranto, M. H. A. Ratul, M. M. Rahman, I. J. Diya, and Z.-B. Zahir, "Performance of machine learning techniques in anomaly detection with basic feature selection strategy - a network intrusion detection system," *Journal of Advanced Information Technology*, vol. 13, no. 1, pp. 45-52, 2022.
9. U. Inayat, M. F. Zia, S. Mahmood, H. M. Khalid, and M. Benbouzid, "Learning-based methods for cyber attacks detection in IoT systems: Methods, analysis, and future prospects," *Electronics*, vol. 11, no. 9, p. 1502, 2022.
10. M. Piskozub, F. De Gaspari, F. Barr-Smith, L. V. Mancini, and I. Martinovic, "MalPhase: Fine-grained malware detection using network flow data," *arXiv preprint arXiv:2106.00541*, 2021.
11. P. Xu, "Android-COCO: Android malware detection with graph neural network for byte- and native-code," *arXiv preprint arXiv:2112.10038*, 2021.
12. A. Muzaffar, H. R. Hassen, M. A. Lones, and H. Zantout, "Reassessing feature-based Android malware detection in a contemporary context," *arXiv preprint arXiv:2301.12778*, 2023.
13. S. Nassar, "Malware detection through memory analysis," in *ELEC 877 (AI For Cybersecurity)*, Queen's University, Canada, 2024.
14. N. Gill, A. Haneef, and S. D. Madhu Kumar, "LLM-FS: Zero-shot feature selection for effective and interpretable malware detection," *arXiv preprint arXiv:2602.09634*, 2026.
15. D. P. Jeong, Z. C. Lipton, and P. Ravikumar, "LLM-select: Feature selection with large language models," *arXiv preprint arXiv:2407.02694*, 2024.