

# Efficient Malware Detection with Optimized Learning on High-Dimensional Features

Aditya Choudhary

Department of CSE,  
School of Computational Sciences  
COEP Technological University  
choudharyap21.comp@coeptech.ac.in

Sarthak Pawar

Department of CSE,  
School of Computational Sciences  
COEP Technological University  
sarthaksp22.comp@coeptech.ac.in

Yashodhara Haribhakta

Department of CSE,  
School of Computational Sciences  
COEP Technological University  
ybl.comp@coeptech.ac.in

**Abstract**—Malware detection using machine learning requires feature extraction from binary files, as models cannot process raw binaries directly. A common approach involves using LIEF for raw feature extraction and the EMBER vectorizer to generate 2381-dimensional feature vectors. However, the high dimensionality of these features introduces significant computational challenges. This study addresses these challenges by applying two dimensionality reduction techniques: XGBoost-based feature selection and Principal Component Analysis (PCA). We evaluate three reduced feature dimensions (128, 256, and 384), which correspond to approximately 5.4%, 10.8%, and 16.1% of the original 2381 features, across four models—XGBoost, LightGBM, Extra Trees, and Random Forest—using a unified training, validation, and testing split formed from the EMBER-2018, ERMDs, and BODMAS datasets. This approach ensures generalization and avoids dataset bias. Experimental results show that LightGBM trained on the 384-dimensional feature set after XGBoost feature selection achieves the highest accuracy of 97.52% on the unified dataset, providing an optimal balance between computational efficiency and detection performance. The best model, trained in 61 minutes using 30 GB of RAM and 19.5 GB of disk space, generalizes effectively to completely unseen datasets, maintaining 95.31% accuracy on TRITIUM and 93.98% accuracy on INFERNO. These findings present a scalable, compute-efficient approach for malware detection without compromising accuracy.

**Index Terms**—Malware detection, Feature selection, XGBoost, Dimensionality reduction, LightGBM.

## I. INTRODUCTION

In an increasingly digital landscape, malware remains one of the most critical threats to cybersecurity. Malware—short for malicious software—exploits vulnerabilities in computer systems, often resulting in data breaches, system compromises, and other forms of cyberattacks. As both the complexity and volume of malware continue to grow, traditional detection techniques, such as signature-based methods, are proving insufficient. In response, machine learning (ML)-based malware detection has emerged as a promising alternative due to its ability to identify previously unseen or zero-day threats by learning patterns from data.

However, ML models cannot process raw malware binaries directly; these must first be converted into structured feature vectors. A widely used tool for feature extraction is LIEF [1], which extracts static features from binaries. These features are then processed by the EMBER vectorizer to produce 2,381-dimensional feature vectors used as input for machine learning

models. While this transformation makes the data suitable for training machine learning models, the high dimensionality leads to heavy computational costs, especially when dealing with large datasets.

To mitigate this, previous studies have primarily followed two approaches. The first involves training models on smaller subsets of the data to reduce computational load. However, this often limits the model’s exposure to diverse patterns, potentially causing it to miss critical features and newly emerging threats. The second approach utilizes high-end computational resources—such as GPUs or cloud infrastructure—to manage large data and feature spaces. While effective, this method is not scalable for continuous learning, where frequent updates are required to counter the evolving malware landscape. The dependency on costly infrastructure limits its practicality in many real-world scenarios.

This paper explores an alternative strategy using dimensionality reduction to enhance computational efficiency without compromising detection performance. We apply two dimensionality reduction techniques—feature selection using XGBoost and Principal Component Analysis (PCA)—to reduce the size of the original 2381-dimensional vectors. Specifically, we generate reduced versions with 128, 256, and 384 dimensions, allowing us to study how varying levels of compression affect model performance and resource usage. We then evaluate four machine learning models—XGBoost, LightGBM, Extra Trees, and Random Forest—on these reduced feature sets. Our results show that dimensionality reduction offers a practical trade-off—maintaining high detection accuracy while significantly lowering training costs—and enables efficient model updates as new malware variants emerge, making it suitable for scalable and continuous malware detection systems.

## II. RELATED WORK

A summary of recent malware detection studies is provided in Table I. The EMBER-2018 dataset has been widely used to evaluate machine learning and deep learning models for static malware detection. Anderson and Roth [2] achieved strong results using LightGBM, reporting an AUC above 0.9911, outperforming early models like MalConv [3]. Wu et al. [4] improved classification accuracy from 15.75% to 93.5% by

TABLE I  
SUMMARY OF PRIOR MALWARE DETECTION STUDIES

Paper	Dataset	Approach	Compute Environment	Detection Metric
Anderson & Roth [2]	EMBER-2018	LightGBM	-	AUC >0.9911
Wu et al. [4]	EMBER-2018	Reinforcement Learning with Gym-Plus	-	Accuracy ↑ from 15.75% to 93.5%
Oyama et al. [5]	EMBER-2018	File Metadata + Imported Functions + LightGBM	Intel Xeon E5-2620, 128GB RAM	Not specified
Vinayakumar et al. [6]	EMBER-2018	WSBD: ML + MalConv	NVIDIA India GPU grant	Accuracy: 98.9%
Pramanik & Teja [7]	EMBER-2018	CNN vs FFNN	-	Precision, Recall, F1: 0.97
Galen & Steele [8]	EMBER-2018	LightGBM on time-sequenced subset	-	Accuracy: 94.80%
Loi et al. [9]	EMBER-2018	ML pipeline with static features	-	Accuracy: 96.9%
Kundu et al. [10]	EMBER-2018	AutoML tuning LightGBM hyperparameters	AWS, 96/72 cores, 192GB RAM	TPR ↑ from 86.8% to 90%
Thosar et al. [13]	EMBER-2018	Gradient Boosting + CNN	Acer Aspire 7, Intel i5 9th Gen, 8GB RAM	Accuracy: 96%
Lad & Adamuthe [14]	EMBER-2018	Deep Learning (Static)	Intel Core i5-4500 CPU, 8GB RAM, GeForce 940M GPU	Accuracy: 94.09%
Vo et al. [15]	EMBER-2018	PEMA (XGBoost, CatBoost, LightGBM)	2x Intel Xeon Platinum, 384GB RAM, 6TB SSD	Accuracy: 97.65%
Shinde et al. [16]	EMBER-2018	Random Forest + Dimensionality Reduction	-	Accuracy: 90%
Dener & Gulburun [17]	EMBER-2018, BODMAS	Supervised + Unsupervised (k-means + feature selection)	Google Colaboratory	Accuracy: 96.77% (EMBER), 99.74% (BODMAS)
Connors & Sarkar [18]	EMBER-2018	Neural Network Model	-	Accuracy: 95.22%
Manikandaraja et al. [19]	TRITIUM, INFERNO	Concept Drift Framework + Adversarial Samples	-	Not specified
Rayankula [20]	BODMAS	K-Nearest Neighbors	MacOS M1, 8GB RAM	Accuracy: 94.9% (Multi), 94.8% (Binary)
Brown et al. [22]	SOREL-20M, EMBER-2018	AutoML (Static and Online)	92 vCPUs, 448GB RAM, 8x Tesla V100 GPUs	Accuracy: 95.8%
Shashank et al. [24]	EMBER-2018	Ensemble Learning (Bagging)	Nvidia DGX Station A100, 4x Tesla A100 GPUs	Accuracy: 96.56%
Maryam et al. [25]	EMBER-2018	SVM (Linear SVC)	-	Accuracy: 98.9% (14.7K samples), dropped to 92.6% (132K samples)
Bhardwaj et al. [26]	BODMAS	MD-ADA: Adversarial Domain Adaptation	-	Accuracy: 99.29%, F1-score: 99.13%
Buriro et al. [27]	BODMAS	Anomaly Detection + Random Forest	-	Accuracy: 99.73%, FPR: 0%
Farfoura et al. [28]	BODMAS	Dimensionality Reduction (MBMD) + Random Forest	-	Accuracy: 99%

incorporating reinforcement learning through Gym-Plus with enhanced PE modification strategies.

Later works explored combining static features such as file metadata and imported functions with ensemble models (Oyama et al. [5]), deep learning frameworks blending CNNs and MalConv (Vinayakumar et al. [6]), and AutoML-based hyperparameter tuning on LightGBM (Kundu et al. [10]), showing improvements in true positive rates at low false positive levels.

Several studies examined the trade-offs between accuracy, training time, and hardware requirements. Thosar et al. [13] and Lad & Adamuthe [14] reported over 94% accuracy using modest hardware. Vo et al. [15] and Shashank et al. [24] used high compute resources to reach nearly 97%, demonstrating the impact of increased computational capacity.

Dimensionality reduction and feature selection were also crucial for handling the high-dimensional EMBER-2018 dataset, as shown by Shinde et al. [16] and Dener & Gulburun [17]. Newer datasets like TRITIUM and INFERNO, introduced by Manikandaraja et al. [19], include adversarial samples to study concept drift and robustness.

On related datasets like BODMAS, highly accurate models (> 99%) have been developed using domain adaptation [26], anomaly detection [27], and novel dimensionality reduction methods [28]. Together, these works highlight the evolving landscape of static malware detection, balancing accuracy, scalability, and robustness across datasets and hardware settings.

### III. PROPOSED METHODOLOGY

#### A. Dataset

For this study, we utilized five datasets to develop and evaluate a generalized and robust malware detection model. These datasets are EMBER-2018, BODMAS, ERMDs, INFERNO, and TRITIUM, summarized in Table II. EMBER-2018 and BODMAS contain malware samples from 2018 and 2019–2020, respectively, providing temporal diversity. ERMDs, INFERNO, and TRITIUM consist of obfuscated and adversarial malware samples exhibiting evasive behaviors, thereby increasing the model’s exposure to real-world evasion tactics and enhancing its generalization.

TABLE II  
DATASET SUMMARY WITH CLASS-WISE SAMPLE DISTRIBUTION

Dataset	Description	Data Distribution
EMBER-2018	EMBER-2018 is a large-scale dataset containing features from 1 million Windows PE files scanned before 2018, designed for challenging malware detection tasks.	400K benign 400K malicious 300K unlabelled
BODMAS	BODMAS is a dataset collected between August 2019 and September 2020.	77,142 benign 57,293 malicious
ERMDS	ERMDS is a dataset created to assess the effectiveness of learning-based malware detection systems against obfuscated malware samples.	30,455 benign 86,685 malicious
TRITIUM	TRITIUM was introduced to analyze concept drift by evaluating threats from a distinct source and threat profile.	10,785 benign 12,471 malicious
INFERNO	INFERNO is a dataset developed to assess classifier robustness against adversarial malware.	1430 benign 1430 malicious

A unified train-validation-test split with stratified sampling was established using EMBER-2018, BODMAS, and ERMDS to encourage generalization and reduce dataset-specific bias. Table III shows the distribution of this unified dataset. Along with the unified test split, the INFERNO and TRITIUM datasets were reserved exclusively for evaluation to test the model’s performance on unseen and adversarial samples.

#### B. Data Preprocessing and Dimensionality Reduction

The datasets used in this study already contained 2,381-dimensional feature vectors extracted using the LIEF tool and EMBER vectorizer methodology described in [2]. As shown in Figure 1, preprocessing began with the raw dataset by removing samples with missing features, followed by the elimination of duplicate entries. These steps ensured data quality and consistency before further processing.

Next, the cleaned feature vectors underwent a two-stage normalization process. First, RobustScaler was applied to reduce the influence of outliers by scaling features based on the median and interquartile range, effectively addressing the skewed distributions common in malware data. Then, MinMaxScaler scaled the features to a uniform range of [0, 1], promoting numerical stability and consistency across all feature dimensions. This combined scaling approach produced input data that was both robust to extreme values and well-normalized, facilitating efficient and stable model training.

TABLE III  
TRAIN, VALIDATION, AND TEST SPLITS FOR UNIFIED DATASET

Dataset	Number of Train Samples	Number of Validation Samples	Number of Test Samples
EMBER-2018	479936	119984	199956
BODMAS	36662	9166	11458
ERMDS	62748	15687	19609
<b>TOTAL</b>	<b>579346</b>	<b>144837</b>	<b>231023</b>

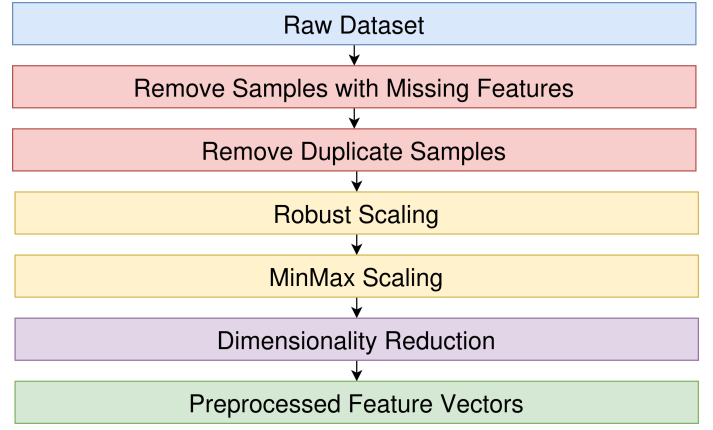


Fig. 1. Preprocessing Pipeline Overview

Finally, to improve computational efficiency without compromising classification performance, dimensionality reduction was applied using two separate techniques: feature selection via XGBoost and Principal Component Analysis (PCA). The original 2,381-dimensional feature space was independently reduced to lower dimensions of 128, 256, and 384 for each technique. For every reduced dimension and method, classification models were trained and evaluated to identify the optimal balance between dimensionality and model accuracy. This approach ensured that the reduced feature sets retained sufficient informative value while significantly lowering computational overhead during training.

#### C. Model Training and Evaluation on Reduced Features

After completing the preprocessing steps (Figure 1), the processed feature vectors were used in the training and evaluation framework, as illustrated in Figure 2. To manage computational constraints, the training and validation data are divided into two equal stratified partitions.

In this study, we evaluate four machine learning classifiers—XGBoost, LightGBM, Extra Trees, and Random Forest—training each independently on both partitions, resulting in two instances per classifier. Hyperparameter tuning for each instance of model is conducted using FLAML [30] to optimize performance.

The outputs of the two instances of each model are combined using a weighted soft voting mechanism. Let  $p_1$  and  $p_2$  denote the probability scores from the first and second model instances, respectively. The final prediction  $\hat{p}$  is computed as:

$$\hat{p} = w_1 \cdot p_1 + w_2 \cdot p_2$$

where  $w_1 + w_2 = 1$ . Optimal weights are identified by iterating  $w_1$  from 0 to 1 in increments of 0.1 during training. The weight pair  $(w_1, w_2)$  yielding the highest performance is selected and consistently applied during validation and testing. Model effectiveness is assessed using accuracy, precision, recall, F1-score, and area under the ROC curve (AUC).

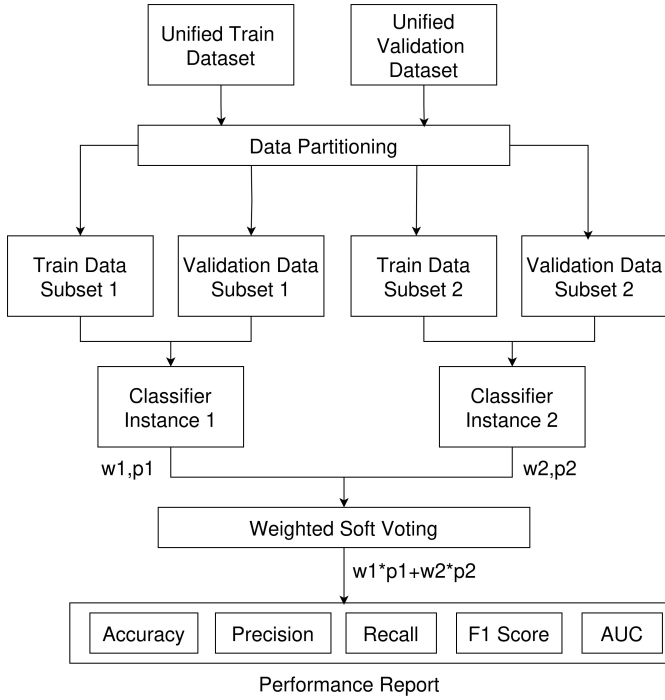


Fig. 2. Proposed system architecture for training and evaluating models on reduced feature sets

## IV. RESULTS AND DISCUSSION

### A. Experimental Setup

All experiments were conducted using the free resources provided by Kaggle<sup>1</sup>. The computational setup included an Intel Xeon 2.20 GHz CPU, 30 GB of RAM, and 19.5 GB of disk space. Several Python libraries were employed throughout the implementation: pandas for dataset handling and operations, numpy for numerical computations, scikit-learn for preprocessing (including scalers and PCA), baseline classifier implementations, and evaluation metrics, and flaml for automated hyperparameter tuning.

### B. Performance on Unified Test Data

Table IV summarizes the performance of four classifiers—XGBoost, LightGBM, Extra Trees, and Random Forest—trained on feature sets reduced to 128, 256, and 384 dimensions using two techniques: supervised feature selection with XGBoost and unsupervised PCA.

Among the models using XGBoost-based feature selection, LightGBM demonstrated the highest performance with 384 features, achieving 97.52% accuracy, 97.73% F1-score, and 99.58% AUC. XGBoost followed closely with 97.25% accuracy and 97.48% F1-score. Extra Trees and Random Forest also performed competitively, particularly at higher feature counts. The steady improvement from 128 to 384 dimensions indicates that additional features contributed discriminative

value, enhancing classification performance. These results confirm that supervised selection effectively retains task-relevant features, benefiting tree-based classifiers. The complete training workflow for the best-performing LightGBM model with 384 features was executed in approximately 61 minutes.

In contrast, PCA-based models showed comparatively lower performance. The best PCA result—XGBoost with 384 components—reached 95.17% accuracy and 95.55% F1-score. Other models exhibited similar trends, with performance increasing from 128 to 384 dimensions but remaining below the XGBoost selection baseline. Random Forest was most affected, with its accuracy dropping to 91.78%. Since PCA preserves variance without considering class labels, it may discard critical features necessary for malware discrimination, weakening downstream classifier performance.

Overall, XGBoost-based feature selection consistently outperformed PCA, confirming the importance of label-aware dimensionality reduction in malware detection. While increasing dimensions led to better results in both methods, the performance gains tended to saturate beyond 256 features, especially for PCA. This suggests a diminishing return on including more components, and underscores the benefit of selecting fewer, yet more informative, features. The combination of supervised selection and tree-based classifiers offers a robust, efficient, and scalable solution well-suited for real-time detection under limited computational resources.

### C. Performance on TRITIUM and INFERNO Dataset

As shown in Table V, we further assessed the robustness and generalization capability of our model on the TRITIUM and INFERNO datasets, which were not used during training or validation. These datasets contain obfuscated and adversarial malware samples, making them ideal benchmarks for evaluating real-world detection effectiveness.

We employed our best-performing pipeline—XGBoost-based feature selection reducing the input to 384 dimensions, followed by LightGBM classification—for this evaluation. On the TRITIUM dataset, the model achieved 95.31% accuracy, with a precision of 99.56%, recall of 91.75%, and F1-score of 95.49%. An AUC of 99.79% further demonstrated strong separability between benign and malicious samples. On the INFERNO dataset, the model maintained robust performance, achieving 93.98% accuracy, 91.99% precision, 96.04% recall, 93.97% F1-score, and an AUC of 98.29%. These results highlight the pipeline’s strong generalization across unseen, evasive threats and reinforce its practical applicability for real-time malware detection in dynamic and adversarial cybersecurity environments.

### D. Comparison with Existing Methods

A detailed comparison of our approach with existing methods is presented in Table VI. Since many prior studies benchmarked their models using the EMBER-2018 dataset, we evaluated our best-performing model (LightGBM trained on 384 features selected via XGBoost-based feature selection) on the same dataset for a fair comparison. Our model achieves

<sup>1</sup><https://www.kaggle.com>

TABLE IV  
PERFORMANCE COMPARISON OF MODELS USING DIFFERENT DIMENSIONALITY REDUCTION TECHNIQUES AND FEATURE SIZES ON UNIFIED TEST DATA

Reduction Technique	Reduced Dimension	Model	Accuracy	Precision	Recall	F1 score	AUC score
Feature Selection using XGBoost	128	XGBoost	96.61%	96.51%	97.28%	96.89%	99.40%
		LightGBM	96.73%	96.78%	97.23%	97.00%	99.40%
		Extra Trees	96.42%	96.76%	96.65%	96.70%	99.34%
		Random Forest	95.44%	95.10%	96.60%	95.84%	99.10%
	256	XGBoost	97.05%	96.96%	97.64%	97.30%	99.51%
		LightGBM	97.34%	97.52%	97.59%	97.56%	99.54%
		Extra Trees	96.79%	97.05%	97.05%	97.05%	99.46%
		Random Forest	95.80%	95.59%	96.73%	96.16%	99.17%
	384	XGBoost	97.25%	97.34%	97.61%	97.48%	99.57%
		<b>LightGBM</b>	<b>97.52%</b>	<b>97.54%</b>	<b>97.91%</b>	<b>97.73%</b>	<b>99.58%</b>
		Extra Trees	97.01%	97.17%	97.33%	97.25%	99.46%
		Random Forest	94.96%	94.44%	96.41%	95.41%	98.86%
PCA	128	XGBoost	94.85%	95.38%	95.13%	95.25%	98.89%
		LightGBM	94.80%	95.39%	95.02%	95.21%	98.85%
		Extra Trees	94.16%	95.03%	94.19%	94.61%	98.71%
		Random Forest	92.35%	92.36%	93.68%	93.02%	97.88%
	256	XGBoost	95.14%	95.63%	95.43%	95.53%	98.98%
		LightGBM	95.00%	95.56%	95.23%	95.39%	98.92%
		Extra Trees	92.69%	92.99%	93.61%	93.30%	98.05%
		Random Forest	91.47%	91.57%	92.86%	92.21%	97.49%
	384	<b>XGBoost</b>	<b>95.17%</b>	<b>95.63%</b>	<b>95.48%</b>	<b>95.55%</b>	<b>98.95%</b>
		LightGBM	94.85%	95.64%	94.84%	95.24%	98.84%
		Extra Trees	94.01%	94.91%	94.04%	94.47%	98.59%
		Random Forest	91.78%	91.61%	93.44%	92.52%	97.73%

a favorable balance between detection accuracy and computational efficiency.

For instance, Vo et al. [15] attained 97.65% accuracy using 2,291 features on high-end hardware with 384 GB RAM and Xeon Platinum CPUs. In contrast, our approach achieves a comparable 97.25% accuracy using only 384 features and 30 GB of RAM. Similarly, Shashank et al. [22] reported 95.80% accuracy using 448 GB RAM and 8 Tesla V100 GPUs, whereas our method surpasses this performance without any GPU dependency. Furthermore, Dener and Gulburun [17] achieved 96.77% accuracy with 448 features, while our model achieves similar performance using fewer features, highlighting its efficiency.

Overall, our approach demonstrates strong predictive performance with minimal computational overhead, reduced feature complexity, and broader scalability—making it a viable alternative to more resource-intensive methods in real-world deployment scenarios.

## V. CONCLUSION AND FUTURE WORK

In conclusion, this study presents an optimized machine learning pipeline for static malware detection by applying dimensionality reduction techniques to high-dimensional feature vectors extracted from PE files. Specifically, XGBoost-based feature selection and Principal Component Analysis

(PCA) were used to reduce the original 2,381-dimensional vectors to 128, 256, and 384 dimensions. Among the evaluated classifiers, LightGBM trained on 384 features selected via XGBoost achieved the highest performance, reaching 97.52% accuracy and strong precision, recall, and F1-score values. The complete training workflow was executed in approximately 61 minutes, demonstrating the practicality of the proposed approach for efficient model development.

For future work, this system can be extended to support malware family classification and made more resilient through adversarial training. Incorporating deep learning-based dimensionality reduction techniques such as autoencoders and applying ensemble learning across heterogeneous classifiers may further enhance detection capabilities. Additionally, leveraging transfer learning could improve adaptability to evolving malware threats and support generalization across diverse datasets.

TABLE V  
EVALUATION OF LIGHTGBM WITH 384 FEATURES SELECTED VIA XGBOOST-BASED FEATURE SELECTION ON TRITIUM AND INFERNO DATASETS

Dataset	Accuracy	Precision	Recall	F1 Score	AUC Score
TRITIUM	95.31%	99.56%	91.75%	95.49%	99.79%
INFERNO	93.98%	91.99%	96.04%	93.97%	98.29%

TABLE VI  
COMPARISON OF VARIOUS APPROACHES ON THE EMBER-2018 DATASET

Model	Test Data Size	Feature Size	Experimental Setup	Test Accuracy
WSBD model	56,000	2351	Recieved GPU Grant from NVIDIA India	98.90%
ML pipeline with static features	2,00,000	2351	-	96.90%
Gradient Boosting + CNN	-	2351	8GB RAM, 512 SSD, Nvidia GeForce GTX 1650, Intel i5 9th gen processor	96.00%
Deep Learning (Static)	2,00,000	2381	Nvidia GeForce 940M 2GB GPU ,Intel Core i5-4500 processor, 8 GB RAM,Google Colab	94.09%
PEMA (XGBoost, CatBoost, LightGBM)	2,00,000	2291	2× Intel Xeon Platinum 8160 ,384GB RAM,6TB SSD	97.65%
k-means + feature selection	2,00,000	448	Google Colaboratory	96.77%
AutoML (Static and Online)	2,00,000	2381	2 vCPUs,448 GB RAM,16 GB of VRAM, 8 Tesla V100 GPUs.	95.80%
Ensemble Learning (Bagging)	1,60,000	2381	Nvidia DGX Station A100 ,AMD EPYC 7742 64-core processor, 4X Tesla A100 ,40GB of GPU memory, and 512 GB of DDR4 RAM.	96.56%
SVM (Linear SVC)	1,32,000	2381	-	92.6%
<b>Our Approach</b>	<b>1,99,956</b>	<b>384</b>	<b>Intel Xeon 2.20 GHz CPU,30GB RAM,19.5 GB disk space</b>	<b>97.25%</b>

#### ACKNOWLEDGMENT

The authors sincerely thank the Data Science and Artificial Intelligence (DSAI) Club of COEP Technological University for their invaluable support and for providing the opportunity to undertake this research. The club's consistent guidance and its commitment to fostering innovation in data science and artificial intelligence played a pivotal role in the successful completion of this work.

#### REFERENCES

- [1] R. Thomas, 'Lief-library to instrument executable formats', Retrieved February, vol. 22, 2017.
- [2] H. S. Anderson and P. Roth, 'EMBER: An Open Dataset for Training Static PE Malware Machine Learning Models', arXiv [cs.CR]. 2018.
- [3] E. Raff, 'Malware detection by eating a whole exe', in Workshops at the thirty-second AAAI conference on artificial intelligence, 2018.
- [4] C. Wu, J. Shi, Y. Yang, and W. Li, 'Enhancing machine learning based malware detection model by reinforcement learning', in Proceedings of the 8th International Conference on Communication and Network Security, 2018, pp. 74–78.
- [5] Y. Oyama, T. Miyashita, and H. Kokubo, 'Identifying useful features for malware detection in the ember dataset', in seventh international symposium on computing and networking workshops, IEEE, 2019.
- [6] R. Vinayakumar, M. Alazab, K. P. Soman, P. Poornachandran, and S. Venkatraman, 'Robust Intelligent Malware Detection Using Deep Learning', IEEE Access, vol. 7, pp. 46717–46738, 2019.
- [7] S. Pramanik and H. Teja, 'Ember-analysis of malware dataset using convolutional neural networks', in 2019 Third International Conference on Inventive Systems and Control (ICISC), 2019, pp. 286–291.
- [8] C. Galen and R. Steele, 'Evaluating performance maintenance and deterioration over time of machine learning-based malware detection models on the ember pe dataset', in 2020 Seventh International Conference on Social Networks Analysis, Management and Security (SNAMS), 2020, pp. 1–7.
- [9] N. Loi, C. Borile, and D. Ucci, 'Towards an automated pipeline for detecting and classifying malware through machine learning', arXiv preprint arXiv:2106.05625, 2021.
- [10] P. P. Kundu, L. Anatharaman, and T. Truong-Huu, 'An empirical evaluation of automated machine learning techniques for malware detection', in Proceedings of the 2021 ACM Workshop on Security and Privacy Analytics, 2021, pp. 75–81.
- [11] N. Erickson et al., 'AutoGluon-Tabular: Robust and Accurate AutoML for Structured Data', arXiv preprint arXiv:2003.06505, 2020.
- [12] Microsoft, Neural Network Intelligence. 1 2021.
- [13] K. Thosar, P. Tiwari, R. Jyothula, and D. Ambawade, 'Effective malware detection using gradient boosting and convolutional neural network', in 2021 IEEE Bombay Section Signature Conference (IBSSC), 2021, pp. 1–4.
- [14] S. S. Lad and A. C. Adamuthe, 'Improved deep learning model for static PE files malware detection and classification', International Journal of Computer Network and Information Security, vol. 11, no. 2, p. 14, 2022.
- [15] H. V. Vo, P. H. Nguyen, H. T. Nguyen, D. B. Vu, and H. N. Nguyen, 'Enhancing AI-Powered Malware Detection by Parallel Ensemble Learning', in 2023 RIVF International Conference on Computing and Communication Technologies (RIVF), 2023, pp. 503–508.
- [16] O. Shinde, A. Khobragade, and P. Agrawal, 'Static malware detection of Ember windows-PE API call using machine learning', in AIP Conference Proceedings, 2023, vol. 2724.
- [17] M. Dener and S. Gulburun, 'Clustering-Aided Supervised Malware Detection with Specialized Classifiers and Early Consensus', CMC-COMPUTERS MATERIALS & CONTINUA, vol. 75, no. 1, 2023.
- [18] C. Connors and D. Sarkar, 'Machine learning for detecting malware in pe files', in 2023 International Conference on Machine Learning and Applications (ICMLA), 2023, pp. 2194–2199.
- [19] A. Manikandaraja, P. Aaby, and N. Pitropakis, 'Rapidrift: Elementary Techniques to Improve Machine Learning-Based Malware Detection', Computers, vol. 12, 2023.
- [20] B. C. Rayankula, 'An Evaluation and Performance study on BODMAS dataset for Malware Analysis', Dublin, National College of Ireland, 2023.
- [21] L. Yang, A. Ciptadi, I. Laziuk, A. Ahmadzadeh, and G. Wang, 'BODMAS: An Open Dataset for Learning based Temporal Analysis of PE Malware', in 2021 IEEE Security and Privacy Workshops (SPW), 2021, pp. 78–84.
- [22] A. Brown, M. Gupta, and M. Abdelsalam, 'Automated machine learning for deep learning based malware detection', Computers & Security, vol. 137, p. 103582, 2024.4
- [23] R. Harang and E. M. Rudd, 'SOREL-20M: A large scale benchmark dataset for malicious PE detection', arXiv preprint arXiv:2012.07634, 2020.
- [24] N. S. Shashank, M. K. Sd, and Others, 'Enhancing Malware Detection: A Comparative Analysis of Ensemble Learning Approaches', in 2024 First International Conference on Technological Innovations and Advance Computing (TIACOMP), 2024, pp. 35–40.
- [25] Z. Maryam, S. J. I. Ismail, and Others, 'Intelligence Malware Detection System with LinearSVC', in 2024 10th International Conference on Wireless and Telematics (ICWT), 2024, pp. 1–6.
- [26] S. Bhardwaj, A. S. Li, M. Dave, and E. Bertino, 'Overcoming the lack of labeled data: Training malware detection models using adversarial domain adaptation', Computers & Security, vol. 140, p. 103769, 2024.
- [27] A. Buriro, A. Rafi, M. A. Yaqub, and F. Luccio, 'Malware Detection using Anomaly Detection Algorithms', in 2024 Fifteenth International Conference on Ubiquitous and Future Networks (ICUFN), 2024, pp. 330–335.
- [28] M. E. Farfoura, I. Mashal, A. Alkhatib, and R. M. Batyha, 'A lightweight machine learning methods for malware classification', Cluster Computing, vol. 28, no. 1, pp. 1–14, 2025.
- [29] L. Jia, 'ERMDS: A obfuscation dataset for evaluating robustness of learning-based malware detection system', BenchCouncil Transactions on Benchmarks, Standards and Evaluations, vol. 3, 2023.
- [30] C. Wang, 'Flaml: A fast and lightweight automl library', Proceedings of Machine Learning and Systems, vol. 3, pp. 434–447, 2021.