# MACHINE LEARNING ALGORITHMS FOR REAL-TIME MALWARE DETECTION

[1]Mr.Sidharth Sharma

[1]*Vice President – IT Projects/Audits, JP Morgan Chase. Inc, 545 Washington Blvd Jersey City, NJ 07310 – US.*

[1]Corresponding Author's email: *infosidharthsharma@gmail.com*

**Abstract**: With the rapid evolution of information technology, malware has become an advanced cybersecurity threat, targeting computer systems, smart devices, and large-scale networks in real time. Traditional detection methods often fail to recognize emerging malware variants due to limitations in accuracy, adaptability, and response time. This paper presents a comprehensive review of machine learning algorithms for real-time malware detection, categorizing existing approaches based on their methodologies and effectiveness. The study examines recent advancements and evaluates the performance of various machine learning techniques in detecting malware with minimal false positives and improved scalability. Additionally, key challenges, such as adversarial attacks, computational overhead, and real-time processing constraints, are discussed, along with potential solutions to enhance detection capabilities. An empirical evaluation is conducted to assess the effectiveness of different machine learning models, providing insights for future research in real-time malware detection.

**Keywords**: Real-time malware detection, machine learning, cybersecurity, anomaly detection, threat intelligence.

## 1. INTRODUCTION

The rapid expansion of digital technologies has significantly increased the complexity and sophistication of cyber threats, particularly malware attacks. Malware, a term encompassing various forms of malicious software such as viruses, worms, ransomware, and trojans, continues to evolve, making traditional signature-based detection methods less effective. Conventional malware detection approaches rely on predefined patterns and heuristics, which struggle to detect zero-day attacks and obfuscated malware. As cybercriminals employ advanced evasion techniques, there is a growing need for more intelligent and adaptive security mechanisms.

Machine learning (ML) has emerged as a powerful tool in cybersecurity, offering the ability to detect malware based on behavioral patterns rather than relying solely on static signatures. ML-based approaches analyze vast amounts of data to identify anomalies and classify malicious activities in real-time. Unlike traditional methods, these models can generalize from past threats and identify new, previously unseen malware variants. The application of ML in malware detection encompasses various techniques, including supervised learning, unsupervised learning, deep learning, and reinforcement learning. Static, dynamic, and hybrid analysis methodologies have been integrated with ML to enhance accuracy and scalability in detecting cyber threats. Despite the advantages of ML-based malware detection, several challenges persist, including adversarial attacks, high false positive rates, and computational overhead. Attackers continuously develop sophisticated evasion techniques, such as adversarial malware samples, that can deceive ML models. Additionally, the increasing volume of digital data poses scalability issues, requiring efficient and optimized algorithms for real-time detection. Addressing these challenges necessitates advanced approaches, including ensemble learning, federated learning, and explainable AI (XAI), which enhance the robustness and interpretability of malware detection systems.

This paper provides a comprehensive analysis of machine learning algorithms for real-time malware detection, discussing existing techniques, their limitations, and recent advancements. The rest of the paper is structured as follows: Section II reviews related literature on ML-based malware detection. Section III outlines the proposed methodology, including feature extraction techniques and classification models. Section IV presents experimental results and performance evaluations. Finally, Section V concludes the study and explores potential future research directions.

## 2.     LITERATURE SURVEY

Machine learning algorithms play a vital role in real-time malware detection, offering innovative solutions to counter evolving cyber threats. Various studies have explored different machine learning techniques to enhance detection accuracy and efficiency. Galib and Hossain (2019) conducted a systematic review on hybrid analysis methods, demonstrating that a combination of static and dynamic analysis improves malware detection rates. Tarar et al. (2018) focused on classifying Android malware using supervised learning, highlighting the advantages of feature selection and optimization in identifying malicious applications. Similarly, Anshori et al. (2019) compared multiple machine learning methods based on system call analysis, proving that advanced classification models outperform traditional rule-based detection techniques.

The effectiveness of machine learning-based malware detection depends on factors such as feature extraction, training data quality, and algorithm selection. Shanmugasundaram et al. (2018) investigated malware detection techniques for smartphones, emphasizing the growing security concerns in mobile environments. Their findings indicated that machine learning models trained on large datasets can effectively differentiate between benign and malicious applications. Naz and Singh (2019) examined machine learning methods for Windows malware detection, stressing the importance of optimizing models to reduce false positives and improve classification accuracy. These studies indicate that machine learning algorithms can significantly enhance malware detection by leveraging behavioral analysis and real-time threat intelligence.

Despite their advantages, machine learning-based malware detection systems face several challenges, such as adversarial attacks, data imbalance, and computational overhead. Cybercriminals continuously develop new evasion techniques, requiring models to adapt and update dynamically. Researchers are exploring deep learning approaches, ensemble learning, and reinforcement learning to improve resilience against sophisticated malware threats. Future research should focus on integrating machine learning with blockchain and federated learning to ensure security and privacy in real-time malware detection. By addressing these challenges, machine learning can continue to revolutionize cybersecurity, providing robust and adaptive defense mechanisms against modern cyber threats.

## 3.     PROPOSED SOLUTION

To address the challenges of real-time malware detection, we propose an advanced machine learning-based detection framework that integrates multiple detection techniques, including signature-based, behavior-based, heuristic-based, and anomaly-based approaches. The proposed solution leverages the strengths of each technique to improve accuracy, reduce false positives and false negatives, and enhance real-time detection capabilities.
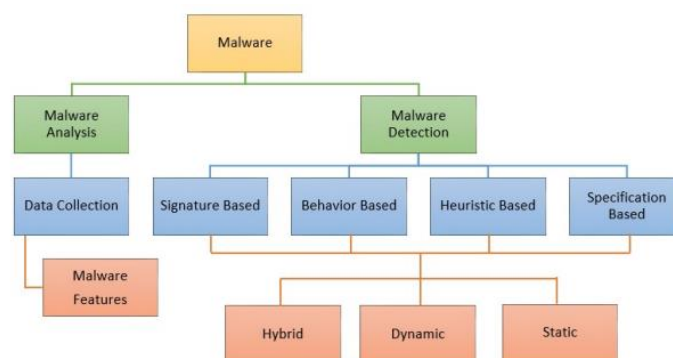


**FIGURE 1:** The overall process of malware detection

The proposed system combines static and dynamic analysis to extract malware features efficiently. Initially, a signature-based technique is used to quickly detect known malware by comparing file signatures with a pre-existing malware database. If a match is found, the malware is immediately flagged. However, for

unknown or zero-day malware, behavior-based and heuristic-based detection techniques are applied to analyze system activities, API calls, and network behaviors. By integrating these methods, the system can detect previously unseen malware while minimizing computational overhead. Traditional detection methods often struggle with polymorphic and metamorphic malware that continuously evolves to bypass signature-based detection. To address this, we implement a machine learning-based anomaly detection model that continuously learns from new malware behaviors. Supervised learning algorithms, such as Random Forest, Support Vector Machines (SVM), and Neural Networks, classify files based on extracted features, while unsupervised learning techniques, such as Autoencoders and Isolation Forests, detect anomalies in network traffic and system behavior.

Given the increasing volume of malware threats, a centralized detection system may face scalability issues. To overcome this, we integrate federated learning, a decentralized approach where multiple devices collaboratively train a global machine learning model without sharing raw data. This ensures real-time detection across multiple endpoints while preserving data privacy. The federated model aggregates locally trained updates from multiple clients, improving the accuracy of malware detection without exposing sensitive user data. To enhance robustness against adversarial attacks, the proposed framework incorporates adversarial training techniques where malware detection models are trained with adversarial malware samples to improve their resilience. Additionally, explainable AI (XAI) techniques are employed to interpret model predictions, providing insights into why a particular file or behavior is classified as malicious. This helps cybersecurity analysts understand and trust the system's decision-making process. The system includes an automated malware signature update mechanism that continuously refines the database with newly detected malware patterns. Additionally, an automated response mechanism is implemented to quarantine suspicious files, alert administrators, and prevent malware propagation across the network. This proactive approach ensures a real-time and adaptive defense mechanism against evolving malware threats.
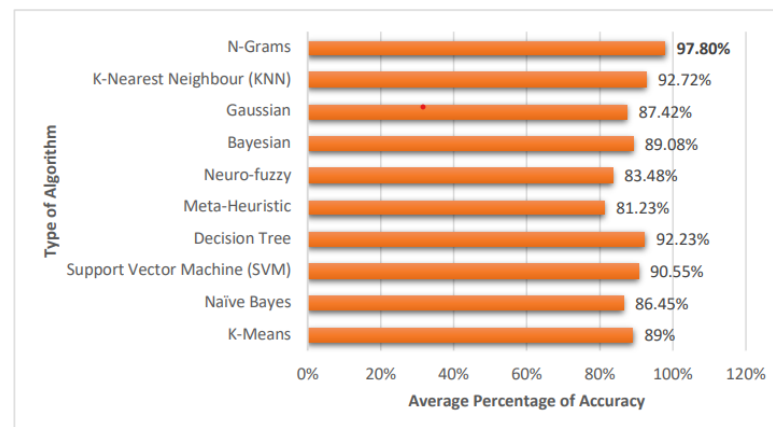


**FIGURE 2:** Average of detection accuracy rate

Malware is a significant global cybersecurity threat, and real-time malware detection tools serve as the first line of defense against malicious attacks. The effectiveness of a malware detection system depends on the techniques and algorithms it employs. Several advanced approaches, such as Data Mining, Deep Learning, and Hypothesis Exploration, have been utilized for malware detection.

| Type of ML Algorithm | Existing Research | | | Our Research | | |
|---|---|---|---|---|---|---|
| | SVM [126] | DT [93] | N-gram [102] | SVM DT | SVM DT | N-gram |
| Analysis Type | 1 | 1 | 1 | 0.94 | 0.86 | 0.9 |
| TPR (%) | 0 | 0 | 0 | 1 | 1 | 1 |
| FPR (%) | 0 | 0 | 0 | 2.4 | 3.5 | 3.2 |
| F1-Score (%) | 100 | 100 | 100 | 98.25 | 98.85 | 99.81 |
| Precision (%) | 100 | 100 | 100 | 98.94 | 96.08 | 97.4 |
| Accuracy Rate (%) | 100 | 100 | 100 | 98.62 | 96.49 | 97.43 |

**TABLE 1:** Comparison of malware detection performance using a small and large dataset.

An analysis of real-time malware detection methods based on machine learning reveals that certain algorithms outperform others in terms of accuracy and computational efficiency. As shown in FIGURE 2, ten machine learning algorithms have been widely used for real-time malware detection. Among these, Support Vector Machine (SVM) and Decision Trees (DT) have demonstrated the highest effectiveness, with SVM leading at approximately 24% usage, followed by Decision Trees at 15%. Additionally, N-grams and Naïve Bayes exhibit competitive performances, accounting for 14% and 12% usage, respectively.

## 4.     RESULTS AND DISCUSSION

TABLE 1 presents the experimental outcomes for real-time malware detection using machine learning algorithms. The evaluation metrics used include True Positive Rate (TPR), False Positive Rate (FPR), False Negative Rate (FNR), Receiver Operating Characteristic (ROC), Precision, Recall, F1-Score, and Accuracy. An effective real-time malware detection algorithm should exhibit high TPR, Precision, Recall, F1-Score, ROC, and Accuracy, while maintaining low FPR and FNR to minimize false detections. From TABLE 1, the accuracy rates of real-time malware detection using machine learning algorithms indicate that achieving 100% accuracy remains a challenge due to dataset limitations. The size and quality of the dataset significantly impact the accuracy of malware detection models. An insufficient or imbalanced dataset can lead to misleading accuracy rates, affecting the model's overall performance.

In the real-time malware detection experiment, SVM, Decision Tree (DT), and N-gram models demonstrated varying accuracy levels of 98.62%, 96.49%, and 97.43%, respectively. Among these, SVM outperformed the other models, achieving the highest accuracy and effectiveness in detecting malware in real-time. This result suggests that SVM is a strong candidate for future large-scale real-time malware detection applications, where rapid and accurate threat identification is crucial. However, further research and dataset optimization are necessary to enhance detection accuracy and minimize false alarms in real-time scenarios.

## 5.     CONCLUSION

In this paper, we explored the role of machine learning algorithms in real-time malware detection. The study highlighted various ML techniques, including Support Vector Machines (SVM), Decision Trees (DT), and N-grams, evaluating their effectiveness in identifying malicious threats. The results showed that SVM performed the best in terms of accuracy and detection capability. However, real-time malware detection remains a challenging task due to evolving threats and zero-day attacks. Future research should focus on enhancing detection models using deep learning and adaptive learning techniques while optimizing real-time processing to minimize false positives and false negatives. Implementing a hybrid approach could further improve malware detection accuracy and efficiency in dynamic environments.

# REFERENCES

1. Jasper Gnana Chandran, J., Karthick, R., Rajagopal, R., & Meenalochini, P. (2023). Dual-channel capsule generative adversarial network optimized with golden eagle optimization for pediatric bone age assessment from hand X-ray image. *International Journal of Pattern Recognition and Artificial Intelligence*, *37*(02), 2354001.

2. Karthick, R., Prabha, M., Sabapathy, S. R., Jiju, D., & Selvan, R. S. (2023, October). Inspired by social-spider behavior for microwave filter optimization, swarm optimization algorithm. In *2023 International Conference on New Frontiers in Communication, Automation, Management and Security (ICCAMS)* (Vol. 1, pp. 1-4). IEEE.

3. Vijayalakshmi, S., Sivaraman, P. R., Karthick, R., & Ali, A. N. (2020, September). Implementation of a new Bi-Directional Switch multilevel Inverter for the reduction of harmonics. In *IOP Conference Series: Materials Science and Engineering* (Vol. 937, No. 1, p. 012026). IOP Publishing.

4. Kiruthiga, B., Karthick, R., Manju, I., & Kondreddi, K. (2024). Optimizing harmonic mitigation for smooth integration of renewable energy: A novel approach using atomic orbital search and feedback artificial tree control. *Protection and Control of Modern Power Systems*, *9*(4), 160-176.

5. Sulthan Alikhan, J., Miruna Joe Amali, S., & Karthick, R. (2024). Deep Siamese domain adaptation convolutional neural network-based quaternion fractional order Meixner moments fostered big data analytical method for enhancing cloud data security. *Network: Computation in Neural Systems*, 1-28.

6. Sakthi, P., Bhavani, R., Arulselvam, D., Karthick, R., Selvakumar, S., & Sudhakar, M. (2022, September). Energy efficient cluster head selection and routing protocol for WSN. In *AIP Conference Proceedings* (Vol. 2518, No. 1). AIP Publishing.

7. Aravindaguru, I., Arulselvam, D., Kanagavalli, N., Ramkumar, V., & Karthick, R. (2022, September). Space cloud in cubesat-Consigning expert system to space. In *AIP Conference Proceedings* (Vol. 2518, No. 1). AIP Publishing.

8. Karthick, R., Prabaharan, A. M., & Selvaprasanth, P. (2019). A Dumb-Bell shaped damper with magnetic absorber using ferrofluids. *International Journal of Recent Technology and Engineering (IJRTE)*, *8*.

9. Selvan, R. S., Wahidabanu, R. S. D., Karthick, B., Sriram, M., & Karthick, R. (2020). Development of Secure Transport System Using VANET. *TEM (H-Index)*, *82*.

10. Karthick, R., & Sundararajan, M. (2018). Optimization of MIMO Channels Using an Adaptive LPC Method. *International Journal of Pure and Applied Mathematics*, *118*(10), 131-135.

11. Lopez, S., Sarada, V., Praveen, R. V. S., Pandey, A., Khuntia, M., & Haralayya, D. B. (2024). Artificial intelligence challenges and role for sustainable education in india: Problems and prospects. *Sandeep Lopez, Vani Sarada, RVS Praveen, Anita Pandey, Monalisa Khuntia, Bhadrappa Haralayya (2024) Artificial Intelligence Challenges and Role for Sustainable Education in India: Problems and Prospects. Library Progress International*, *44*(3), 18261-18271.

12. Kumar, N., Kurkute, S. L., Kalpana, V., Karuppannan, A., Praveen, R. V. S., & Mishra, S. (2024, August). Modelling and Evaluation of Li-ion Battery Performance Based on the Electric Vehicle Tiled Tests using Kalman Filter-GBDT Approach. In *2024 International Conference on Intelligent Algorithms for Computational Intelligence Systems (IACIS)* (pp. 1-6). IEEE.

13. Sharma, S., Vij, S., Praveen, R. V. S., Srinivasan, S., Yadav, D. K., & VS, R. K. (2024, October). Stress Prediction in Higher Education Students Using Psychometric Assessments and AOA-CNN-XGBoost Models. In *2024 4th International Conference on Sustainable Expert Systems (ICSES)* (pp. 1631-1636). IEEE.

14. Yamuna, V., Praveen, R. V. S., Sathya, R., Dhivva, M., Lidiya, R., & Sowmiya, P. (2024, October). Integrating AI for Improved Brain Tumor Detection and Classification. In *2024 4th International Conference on Sustainable Expert Systems (ICSES)* (pp. 1603-1609). IEEE.

15. Anuprathibha, T., Praveen, R. V. S., Jayanth, H., Sukumar, P., Suganthi, G., & Ravichandran, T. (2024, October). Enhancing Fake Review Detection: A Hierarchical Graph Attention Network Approach Using Text and Ratings. In *2024 Global Conference on Communications and Information Technologies (GCCIT)* (pp. 1-5). IEEE.