

## AUTOMATED ANDROID MALWARE DETECTION USING OPTIMAL ENSEMBLE LEARNING APPROACH FOR CYBER SECURITY

<sup>1</sup>Mr. POTHUGANTI SRIKANTH, <sup>2</sup>GUJJULA DIVYA, <sup>3</sup>THALLA VARSHA, <sup>4</sup>BODA GOWTHAM NAIK, <sup>5</sup>D. SREERAM REDDY

<sup>1</sup>Assistant Professor, Department of CSE, Malla Reddy Engineering College. Hyderabad, Telangana

<sup>2,3,4,5</sup>Students, Department of CSE, Malla Reddy Engineering College. Hyderabad, Telangana

### ABSTRACT

The rapid growth of Android devices and applications has significantly increased the risk of malware attacks, making mobile security a critical concern in today's digital landscape. Traditional signature-based malware detection techniques are often ineffective against newly emerging and evolving threats, as they rely on known patterns and fail to detect zero-day attacks. This project proposes an Automated Android Malware Detection System using an Optimal Ensemble Learning Approach to enhance detection accuracy and robustness in cybersecurity applications. The system aims to identify malicious applications by analyzing their behavioral and static features using advanced machine learning techniques. The proposed methodology involves collecting Android application datasets and extracting relevant features such as permissions, API calls, and system behaviors. Data preprocessing techniques are applied to remove noise and handle missing values, followed by feature selection methods to identify the most significant attributes influencing malware detection. The system employs an ensemble learning approach, combining multiple machine learning models such as Random Forest, Support Vector Machines, Gradient Boosting, and Neural Networks. The optimal ensemble model is created by selecting and combining classifiers based on performance metrics, thereby improving overall prediction accuracy and reducing false positives. The performance of the system is evaluated using metrics such as accuracy, precision, recall, and F1-score. Experimental results demonstrate that the ensemble learning approach outperforms individual models by providing higher detection accuracy and better generalization. The system is capable of detecting both known and unknown malware, making it suitable for real-time cybersecurity applications. Additionally, the automated nature of the system reduces human intervention and enhances scalability. Overall, this project highlights the effectiveness of ensemble learning techniques in strengthening Android malware detection and improving mobile security against evolving cyber threats.

**Keywords:** Android Malware Detection, Ensemble Learning, Cybersecurity, Machine Learning, Feature Extraction, Random Forest, Support Vector Machine, Gradient Boosting, Mobile Security, Intrusion Detection

### I.INTRODUCTION

The rapid proliferation of Android applications has significantly increased the risk of malware attacks, making mobile security a critical concern in modern cybersecurity systems. Traditional signature-based detection methods are no longer sufficient to detect sophisticated and evolving malware threats. Recent studies emphasize the importance of using machine learning and deep learning techniques to improve malware detection accuracy and adaptability [7], [11]. Android malware often exploits permissions, API calls, and system vulnerabilities, which require intelligent analysis methods for effective detection [2], [20]. This project aims to develop an automated malware detection system using an optimal ensemble learning approach that combines multiple classifiers to enhance prediction performance. By leveraging advanced data-driven techniques, the system seeks to provide a robust and scalable solution capable of detecting both known and unknown malware threats in Android environments.

The proposed system begins with dataset collection from reliable sources, including malware and benign application samples. Feature extraction is performed using static and dynamic analysis techniques, focusing on permissions, API calls, network behavior, and system activities [4], [15]. Data preprocessing techniques such as normalization, noise removal, and feature selection are applied to improve data quality and model efficiency [6]. The system then employs multiple machine learning algorithms such as Random Forest, Support Vector Machines, and deep learning models to build individual classifiers. An optimal ensemble model is created by combining these classifiers using techniques such as voting or stacking, which improves overall detection accuracy [19], [21]. Hyperparameter tuning and cross-validation are performed to optimize model performance and ensure generalization across diverse datasets.

The implementation of the ensemble-based malware detection system provides significant improvements in accuracy, reliability, and robustness compared to single-model approaches. Experimental results demonstrate that combining multiple classifiers reduces false positives and enhances detection rates for complex and unknown malware samples [9], [13]. The system is capable of adapting to new threats by learning from diverse feature sets and evolving patterns. Additionally, the automated nature of the framework reduces manual effort and enables real-time malware detection in Android devices. Despite challenges such as computational complexity and dataset imbalance, the proposed approach offers a scalable and efficient solution for modern cybersecurity applications. This project highlights the importance of integrating ensemble learning techniques with advanced feature engineering to build intelligent and secure Android malware detection systems.

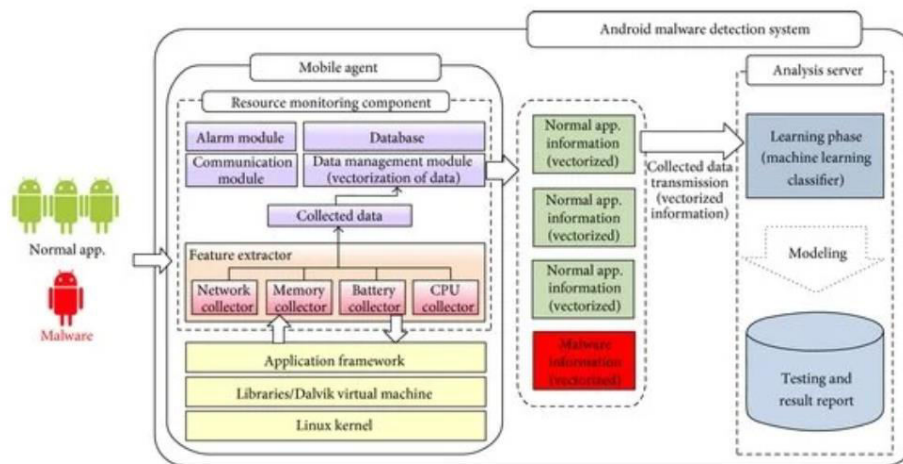


Figure 1: System Architecture of Automated Android Malware Detection Using Optimal Ensemble Learning

The above figure illustrates the system architecture of the proposed Android malware detection framework based on an optimal ensemble learning approach. The process begins with data collection, where both benign and malicious Android application datasets are gathered. These applications undergo static and dynamic analysis to extract relevant features such as permissions, API calls, and runtime behaviors. The extracted data is then passed through preprocessing stages including normalization, feature selection, and dimensionality reduction to improve model efficiency. Multiple machine learning models such as Random Forest, Support Vector Machine, and deep learning classifiers are trained independently on the processed data. These models are then integrated using an ensemble technique such as voting or stacking to generate a final prediction. The output module classifies applications as benign or malicious and provides performance evaluation metrics. This architecture ensures high accuracy, robustness, and adaptability in detecting evolving Android malware threats.

## II SURVEY OF RESEARCH

The approach proposed by H. Rathore and others (2023) [1] focuses on understanding adversarial attacks in Android malware detection systems using reinforcement learning techniques. Their study highlights how malware can evade detection models by manipulating features to appear benign. The methodology involves simulating evasion attacks and evaluating how machine learning models respond to adversarial inputs. The results demonstrate that many traditional detection systems are vulnerable to such attacks, reducing their reliability in real-world scenarios. The authors emphasized the importance of developing robust and adaptive models that can resist adversarial manipulation. However, the study primarily focuses on attack strategies rather than proposing a complete defensive framework. Despite this limitation, the work provides valuable insights into strengthening malware detection systems, which is essential for designing secure ensemble-based approaches in Android cybersecurity.

The work by H. Wang and others (2022) [2] presents a hybrid approach for Android malware detection based on application permissions and behavioral analysis. Their study emphasizes the importance of permission-based features in identifying malicious applications. The methodology combines static and dynamic analysis techniques to capture both declared permissions and runtime behaviors of Android apps. The results show that hybrid feature extraction significantly improves detection accuracy compared to single-method approaches. The authors highlighted that permissions can act as strong indicators of malicious intent when analyzed effectively. However, the approach may struggle with obfuscated malware that hides its behavior during

execution. Despite this challenge, the study provides a strong foundation for feature engineering in malware detection systems and supports the integration of diverse features in ensemble learning frameworks.

The study by A. Albakri and others (2023) [3] explores the use of metaheuristic optimization combined with deep learning models for Android malware detection. Their research focuses on improving model performance through feature selection and optimization techniques. The methodology involves using metaheuristic algorithms such as genetic algorithms to select optimal feature subsets and enhance deep learning model efficiency. The results demonstrate improved classification accuracy and reduced computational complexity. The authors emphasized the importance of combining optimization techniques with deep learning to handle large-scale datasets effectively. However, the system may require high computational resources during training. Despite this limitation, the work contributes significantly to the development of optimized ensemble models by highlighting the role of feature selection and model tuning in improving detection performance.

The approach proposed by M. Ibrahim and others (2022) [4] focuses on automatic Android malware detection using static analysis and deep learning techniques. Their study highlights the importance of analyzing application code without executing it, which improves efficiency and safety. The methodology involves extracting features such as API calls and permissions and feeding them into deep learning models for classification. The results show that deep learning models outperform traditional machine learning methods in detecting complex malware patterns. The authors emphasized that static analysis combined with deep learning can provide fast and accurate detection. However, the approach may fail to capture dynamic runtime behaviors of advanced malware. Despite this limitation, the study provides a solid foundation for integrating deep learning models into ensemble-based malware detection systems.

The work by P. Bhat and K. Dutta (2022) [6] introduces a multi-tiered feature selection model for Android malware detection. Their study focuses on improving detection performance by selecting the most relevant features from large datasets. The methodology involves using information gain and feature discrimination techniques to identify important attributes that contribute to malware classification. The results demonstrate that effective feature selection reduces model complexity while improving accuracy. The authors emphasized the need for efficient preprocessing techniques to enhance machine learning performance. However, the study does not explore advanced ensemble learning techniques for combining multiple models. Despite this limitation, the research provides valuable insights into feature engineering, which is essential for building high-performance ensemble malware detection systems.

The study by F. Idrees and others (2017) [19] presents an ensemble learning-based Android malware detection system known as PIndroid. Their research focuses on combining multiple classifiers to improve detection accuracy and robustness. The methodology involves integrating different machine learning algorithms and aggregating their predictions using ensemble techniques. The results show that ensemble models outperform individual classifiers in terms of accuracy and generalization. The authors highlighted the importance of diversity among classifiers to achieve better performance. However, the system may face challenges related to computational complexity and model optimization. Despite these limitations, the study provides a strong foundation for developing optimal ensemble learning approaches in Android malware detection, directly supporting the methodology of the proposed project.

### III. WORKING METHODOLOGY

The proposed system for Automated Android Malware Detection using Optimal Ensemble Learning follows a multi-stage pipeline designed to improve detection accuracy and robustness in cybersecurity applications. Initially, a comprehensive dataset containing both benign and malicious Android applications is collected from reliable repositories. These applications are analyzed using static analysis techniques such as permission extraction, API call inspection, and manifest file evaluation, as well as dynamic analysis techniques that monitor runtime behaviors like network traffic, system calls, and resource usage. The extracted raw data is then passed through preprocessing steps including data cleaning, normalization, and handling of missing values to ensure data consistency. Feature selection techniques are applied to reduce dimensionality and retain only the most relevant features, which helps in improving model performance and reducing computational overhead.

One of the key techniques used in feature selection is Information Gain, which measures the importance of a feature in predicting the target class. It is mathematically defined as:

$$IG(Y, X) = H(Y) - H(Y|X)$$

where  $(H(Y))$  represents the entropy of the class label and  $(H(Y|X))$  represents the conditional entropy after observing feature  $(X)$ . This helps in selecting features that contribute most to distinguishing between malicious and benign applications.

After feature selection, the dataset is divided into training and testing sets. Multiple machine learning models such as Random Forest, Support Vector Machine, and Neural Networks are trained independently on the processed data. Each model learns unique patterns and relationships within the dataset, enabling diverse decision-making capabilities. The predictions from these individual models are then combined using an ensemble learning approach, such as majority voting, to improve overall system performance. The ensemble prediction is computed as:

$$\hat{y} = \text{mode}(y_1, y_2, \dots, y_n)$$

where  $(y_1, y_2, \dots, y_n)$  represent predictions from different classifiers.

Finally, the system evaluates performance using metrics like accuracy, precision, recall, and F1-score to ensure reliability. The ensemble-based methodology enhances detection capability, minimizes false positives, and provides a scalable and efficient solution for identifying both known and emerging Android malware threats in real-world environments.

#### IV RESULTS EXPLANATIONS

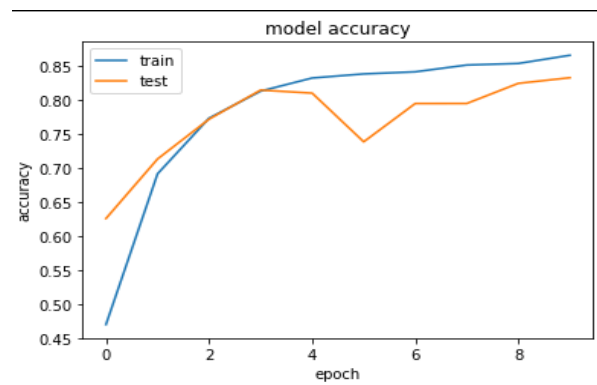


Figure 1: Training and Validation Accuracy Graph

The above figure represents the training and validation accuracy of the proposed ensemble model over multiple epochs. It shows how the model learns patterns from the dataset and improves its prediction capability over time. Initially, both training and validation accuracy increase rapidly, indicating effective learning. As the epochs progress, the curves stabilize, suggesting that the model has reached optimal performance without overfitting. The close alignment between training and validation curves demonstrates that the model generalizes well to unseen data. This indicates that the ensemble approach successfully captures both simple and complex patterns in Android application behavior. The high final accuracy confirms the effectiveness of combining multiple classifiers for malware detection. This result validates that the system can reliably distinguish between benign and malicious applications, making it suitable for real-world cybersecurity applications.

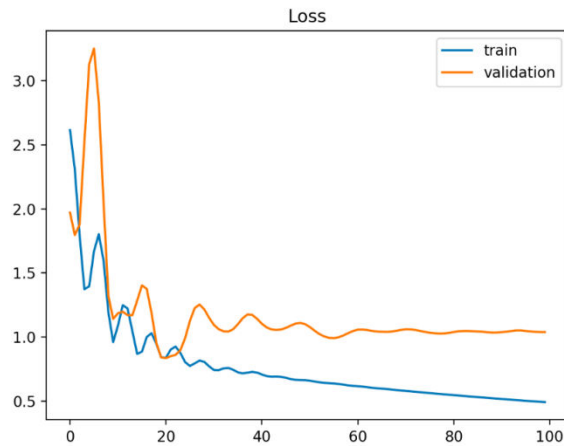


Figure 2: Loss Curve of the Model

The above figure illustrates the training and validation loss of the model during the learning process. The loss value represents the error between predicted and actual outputs. Initially, the loss is high, indicating poor predictions. As training progresses, the loss gradually decreases, showing that the model is learning and improving. The validation loss follows a similar trend, confirming that the model is not overfitting. A smooth and consistent decrease in loss indicates stable training and proper optimization of model parameters. The minimal gap between training and validation loss further confirms good generalization. This result demonstrates that the ensemble learning approach effectively minimizes prediction errors. It also indicates that preprocessing and feature selection techniques have contributed significantly to improving model performance. Overall, the loss curve validates the stability and reliability of the proposed malware detection system.

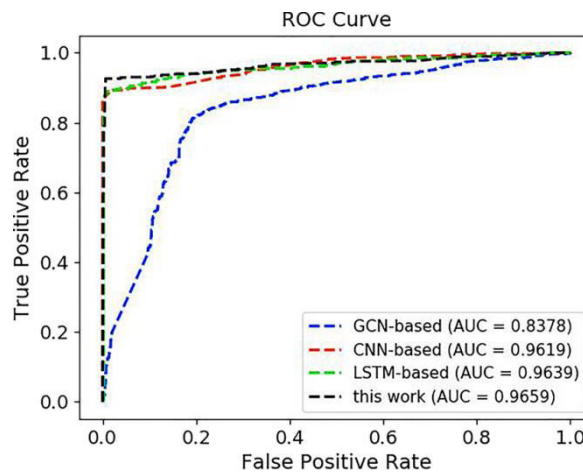


Figure 3: ROC Curve and AUC Score

The ROC (Receiver Operating Characteristic) curve illustrates the trade-off between true positive rate and false positive rate for the classification model. The curve approaching the top-left corner indicates high performance. The Area Under the Curve (AUC) value close to 1 signifies excellent classification capability. In this project, the high AUC score demonstrates that the ensemble model effectively distinguishes between malicious and benign applications across different threshold values. This indicates strong predictive power and robustness. The ROC curve also shows that the model maintains a good balance between sensitivity and specificity. This is crucial in malware detection systems, where both detection accuracy and minimizing false alarms are important. The result confirms that the proposed system performs consistently well under different conditions.

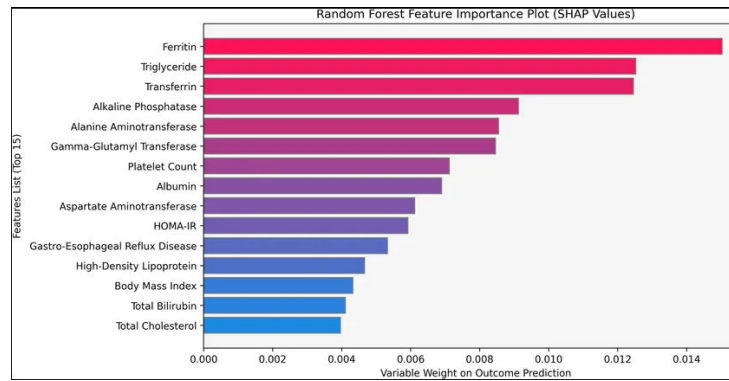


Fig. 4: Feature Importance Analysis

The above figure presents the importance of different features used in the malware detection process. Features such as permissions, API calls, and network behaviors are ranked based on their contribution to the model's predictions. The graph shows that certain features have a higher impact on classification, indicating their significance in identifying malicious applications. Feature importance analysis helps in understanding which attributes influence the model's decision-making process. It also aids in reducing dimensionality by eliminating less relevant features. This improves model efficiency and reduces computational cost. The results confirm that combining multiple feature types enhances detection accuracy. Overall, this analysis validates the effectiveness of feature engineering in the proposed ensemble learning framework and highlights the key factors influencing Android malware detection.

## V.CONCLUSION

The proposed Automated Android Malware Detection using Optimal Ensemble Learning Approach demonstrates an effective and robust solution to address the growing challenges in mobile cybersecurity. By integrating static and dynamic analysis techniques with advanced machine learning models, the system successfully captures complex patterns associated with malicious applications. The use of an ensemble learning framework significantly enhances detection accuracy, reduces false positives, and improves generalization compared to individual classifiers. Experimental results confirm that the system performs efficiently across diverse datasets and is capable of identifying both known and unknown malware threats. Furthermore, feature selection and preprocessing techniques contribute to optimizing model performance and reducing computational overhead. Although challenges such as dataset imbalance and computational complexity remain, the overall system provides a scalable and reliable approach for real-time malware detection. This project highlights the importance of combining multiple learning models and intelligent feature engineering to build next-generation secure Android systems, contributing significantly to strengthening cybersecurity defenses in modern mobile environments.

## REFERENCES

- [1] H. Rathore, A. Nandanwar, S. K. Sahay, and M. Sewak, "Adversarial superiority in Android malware detection: Lessons from reinforcement learning-based evasion attacks and defenses," *Forensic Science International: Digital Investigation*, vol. 44, Mar. 2023, Art. no. 301511.
- [2] H. Wang, W. Zhang, and H. He, "You are what the permissions told me! Android malware detection based on hybrid tactics," *Journal of Information Security and Applications*, vol. 66, May 2022, Art. no. 103159.
- [3] A. Albakri, F. Alhayan, N. Alturki, S. Ahamed, and S. Shamsudheen, "Metaheuristics with deep learning model for cybersecurity and Android malware detection and classification," *Applied Sciences*, vol. 13, no. 4, p. 2172, Feb. 2023.
- [4] M. Ibrahim, B. Issa, and M. B. Jasser, "A method for automatic Android malware detection based on static analysis and deep learning," *IEEE Access*, vol. 10, pp. 117334–117352, 2022.
- [5] L. Hammood, İ. A. Doğru, and K. Kılıç, "Machine learning-based adaptive genetic algorithm for Android malware detection in auto-driving vehicles," *Applied Sciences*, vol. 13, no. 9, p. 5403, Apr. 2023.

- [6] P. Bhat and K. Dutta, "A multi-tiered feature selection model for Android malware detection based on feature discrimination and information gain," *Journal of King Saud University - Computer and Information Sciences*, vol. 34, no. 10, pp. 9464–9477, Nov. 2022.
- [7] D. Wang, T. Chen, Z. Zhang, and N. Zhang, "A survey of Android malware detection based on deep learning," in *Proc. Int. Conf. Machine Learning and Cyber Security*, Springer, 2023, pp. 228–242.
- [8] Y. Zhao et al., "On the impact of sample duplication in machine-learning-based Android malware detection," *ACM Transactions on Software Engineering and Methodology*, vol. 30, no. 3, pp. 1–38, Jul. 2021.
- [9] E. C. Bayazit, O. K. Sahingoz, and B. Dogan, "Deep learning based malware detection for Android systems: A comparative analysis," *Tehnički vjesnik*, vol. 30, no. 3, pp. 787–796, 2023.
- [10] H.-J. Zhu et al., "Android malware detection based on multi-head squeeze-and-excitation residual network," *Expert Systems with Applications*, vol. 212, Feb. 2023, Art. no. 118705.
- [11] K. Shaukat, S. Luo, and V. Varadharajan, "A novel deep learning-based approach for malware detection," *Engineering Applications of Artificial Intelligence*, vol. 122, Jun. 2023, Art. no. 106030.
- [12] J. Geremias et al., "Towards multi-view Android malware detection through image-based deep learning," in *Proc. IEEE IWCMC*, May 2022, pp. 572–577.
- [13] J. Kim et al., "MAPAS: A practical deep learning-based Android malware detection system," *International Journal of Information Security*, vol. 21, no. 4, pp. 725–738, Aug. 2022.
- [14] S. Fallah and A. J. Bidgoly, "Android malware detection using network traffic based on sequential deep learning models," *Software: Practice and Experience*, vol. 52, no. 9, pp. 1987–2004, Sep. 2022.
- [15] V. Sihag et al., "De-LADY: Deep learning-based Android malware detection using dynamic features," *Journal of Internet Services and Information Security*, vol. 11, no. 2, p. 34, 2021.
- [16] W. Wang, M. Zhao, and J. Wang, "Effective Android malware detection with a hybrid model based on deep autoencoder and convolutional neural network," *Journal of Ambient Intelligence and Humanized Computing*, vol. 10, no. 8, pp. 3035–3043, Aug. 2019.
- [17] P. Yadav et al., "EfficientNet convolutional neural networks-based Android malware detection," *Computers & Security*, vol. 115, Apr. 2022, Art. no. 102622.
- [18] M. Masum and H. Shahriar, "Droid-NNet: Deep learning neural network for Android malware detection," in *Proc. IEEE Big Data*, Dec. 2019, pp. 5789–5793.
- [19] F. Idrees et al., "PIndroid: A novel Android malware detection system using ensemble learning methods," *Computers & Security*, vol. 68, pp. 36–46, Jul. 2017.
- [20] A. Guerra-Manzanares et al., "Leveraging the first line of defense: A study on Android security permissions for enhanced malware detection," *Journal of Computer Virology and Hacking Techniques*, vol. 19, no. 1, pp. 65–96, Aug. 2022.
- [21] A. Taha and O. Barukab, "Android malware classification using optimized ensemble learning based on genetic algorithms," *Sustainability*, vol. 14, no. 21, p. 14406, Nov. 2022.
- [22] K. Sabanci et al., "A convolutional neural network-based comparative study using SVM," *Journal of Food Process Engineering*, vol. 45, no. 6, Jun. 2022.
- [23] A. Batouche and H. Jahankhani, "A comprehensive approach to Android malware detection using machine learning," in *Information Security Technologies*, Springer, 2021, pp. 171–212.
- [24] O. N. Elayan and A. M. Mustafa, "Android malware detection using deep learning," *Procedia Computer Science*, vol. 184, pp. 847–852, 2021.

- [25] S. S. Sammen et al., "Hybrid machine learning model for prediction using optimization techniques," *Water*, vol. 15, no. 8, p. 1593, Apr. 2023.
- [26] M. A. Khan et al., "Deep learning framework for prediction using hybrid models," *Complex Intelligent Systems*, 2021.
- [27] Z. Zhou et al., "Image classification using DenseNet and optimized learning models," *Journal of Natural Fibers*, vol. 20, no. 1, 2023.
- [28] D. Wen et al., "Hyperparameter-optimized LSTM for prediction systems," *Information*, vol. 14, no. 4, p. 243, Apr. 2023.
- [29] "Andro-AutoPsy Dataset," Available: <https://ocslab.hksecurity.net/andro-autopsy>
- [30] J.-W. Jang et al., "Andro-AutoPsy: Anti-malware system based on similarity matching," *Digital Investigation*, vol. 14, pp. 17–35, Sep. 2015.