



# Improved Malware Detection in Android Using Similarity-Aware GRU Architectures

<sup>1</sup>SARIKA SIVASANKAR, Student in Dept. Of Master of Computer Applications, at Miracle Educational Society Group of Institutions

<sup>2</sup>Dr. BEHARA SREENIVASA RAO, Miracle Educational Society Group of Institutions

<sup>3</sup>Dr. SHAIK ASADULLA HUSSAIN, Miracle Educational Society Group of Institutions

<sup>1</sup>sarikasiva1098@gmail.com

## ABSTRACT:

The surge in the deployment of Android malware is one of the gravest dangers to mobile users, and data security in general. This project proposes a new detection technique based on SIMGRU, an improved version of the Gated Recurrent Unit (GRU) model, applied through static analysis with similarity clustering. We construct three SIMGRU models: InputSimGRU, HiddenSimGRU, and InputHiddenSimGRU, each of which is aimed at comprehensively capturing distinct malware patterns. With the introduction of similarity metrics based on Euclidean distance in the GRU layers, the accuracy of detection of malicious apps is greatly enhanced. Results on the Drebin dataset indicate that SIMGRU provides an increase in performance over GRU and other machine learning techniques, thus giving a robust and resourceful static Android malware detection solution.

**Keywords:** Malware, SIMGRU, Android

## INTRODUCTION

The increase in Android users has led to a parallel increase in mobile applications, which has created a new and threatening environment for sophisticated malware. Android malware violates user information, halts device operations, and is a potential severe threat for both organizations and individual users. Malware detection methods that rely on signatures or use simple learning models often fall behind, as these new and unique malware programs to face a devastatingly

new challenge. This project suggests an innovative solution to the challenge by integrating SIMGRU, an improved version of the GRU (Gated Recurrent Unit) which utilizes similarity metrics, to perform static analysis. The contribution of architecture InputSimGRU, HiddenSimGRU, as well as InputHiddenSimGRU, allows the application of deeper contextual feature analysis to the application's context, thus improving the accuracy and reliability of malware classification. In addition to improving the detection rate of previously

unseen malware, this method also fortifies the static detection system's resilience. All in all, SIMGRU provides an Android platform intelligence and agility in malware detection that outpaces adaptive threats.

## RELATED WORK

The past few years have witnessed rapid improvements in Android malware detection, incorporating both static and dynamic methodologies. One noteworthy work is by Zhang et al. (2016) who proposed the use of a keyword vector and SVM-based method to correlate keywords within the Java source code of malware. Unlike binary-based techniques, this method significantly enhanced detection accuracy, thus proving the importance of semantic analysis for malware detection. In another contribution, Arp et al. (2014) proposed a lightweight approach called Drebin that statically analyzes Android apps using their granted permissions and API calls for classification. While effective against known threats, it struggles to detect obfuscated or repackaged versions of malware, which prompted looking for more advanced solutions. To resolve the scalability concern, Li et al. (2019) created FalDroid, which uses frequent subgraph analysis to classify malware families and identify representative samples. This approach increased detection speed and reduced redundancy during analysis, although it focused more on malware families and less on individual behavior patterns. Further improving static

detection, Wang et al. (2020) proposed a Naive Bayes-based model using correlation of permissions with static features. This model improved detection rate to 86.54% by weighting newly granted permissions, proving that features of static classifiers need not be ignored. Finally, Wang et al. (2021) presented Deep4MalDroid, a framework that uses deep learning to examine the system call graphs of the executables in real time. The challenge with deep analysis is that it can be very costly in terms of computer resources, making it quite difficult to utilize in real time.

TABLE1. Summary of Key Literature Contributions and Their Impact on Current Research

Author(s)	Contribution	Impact on Research
Zhang et al. (2016)	Introduced keyword vector extraction with SVM for improved source code analysis	Inspired use of semantic-based features for static malware detection
Arp et al. (2014)	Developed Drebin, a static analysis tool using permissions and API calls	Established static analysis as a viable low-resource approach
Li et al. (2019)	Proposed FalDroid with frequent subgraph-based family classification	Improved detection scalability and highlighted graph-based behavioral clustering
Wang et al. (2020)	Created permission correlation model with weighted Naive Bayes algorithm	Demonstrated effectiveness of permission analysis with statistical feature selection
Wang et al. (2021)	Built Deep4MalDroid using system call	Emphasized deep learning's strength in handling

	graphs and deep learning	obfuscation but exposed runtime overhead
--	--------------------------	------------------------------------------

## PROPOSED APPROACH

The proposed approach aims to improve Android malware detection using SIMGRU, an extension of the Gated Recurrent Unit (GRU) model which embeds similarity clustering into the learning process. Conventional GRUs face challenges in capturing intricate correlations between the input and hidden features due to the reliance on sequential data processing. SIMGRUs models overcome this limitation by adding similarity-based filters to the GRU layers. The three variants for the models are input SIMGRU which computes similarity among input features, hidden SIMGRU which focuses on similarity of hidden states, and input hidden SIMGRU which uses both input and hidden similarities through Euclidean distance metrics. With these additional structures, the model can focus on capturing malware features with recurring structures across diverse applications and broaden the network's generalization capability to identify novel threats. For the approach, static analysis was used to extract features like permissions and API calls from Android APKs, which allows for lightweight deployment without the overhead of runtime execution. The similarity-grained enhancements help the network achieve the focus on relevant recurring patterns that improve detection accuracy, precision,

and recall relative to GRU models. With the capability of SIMGRU to infuse context into the learning process, this approach becomes more intelligent, efficient, and scalable for Android malware detection.

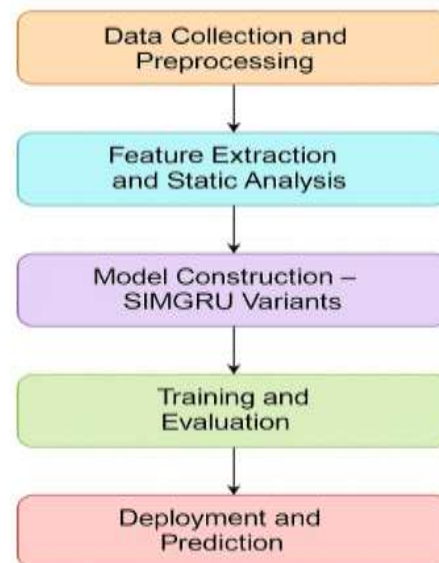


Figure 1: Proposed Android malware detection System

## METHODOLOGIES

### 1. Data Collection and Preprocessing:

The Drebin dataset, a widely used Android malware corpus, serves as the foundation. It includes both benign and malicious app samples. Preprocessing involves handling missing values, encoding categorical features, and normalizing input data. Label encoding transforms malware classes into binary outputs suitable for machine learning models.

### 2. Feature Extraction and Static Analysis:

Using static techniques, features like permissions, API calls, and manifest properties are extracted without executing

the application. This enables faster, low-resource feature collection while avoiding exposure to live threats.

### 3. Model Construction – SIMGRU

#### Variants:

Three GRU-based models are implemented:

- **InputSimGRU** applies Euclidean similarity to input layers to highlight relevant feature similarities.
- **HiddenSimGRU** focuses on comparing hidden state similarities between time steps.
- **InputHiddenSimGRU** combines both to exploit a broader range of feature interactions.

Each model includes multiple GRU layers with dropout regularization to prevent overfitting, followed by dense layers for binary classification.

### 4. Training and Evaluation:

Models are trained on 80% of the dataset and validated on the remaining 20%. Key metrics such as accuracy, precision, and recall are monitored. Hyperparameter tuning is performed to refine learning rates, dropout rates, and layer configurations.

### 5. Comparative Analysis and Visualization:

The SIMGRU variants are compared against traditional GRU and CNN models to evaluate improvements. Visual plots of training vs. validation accuracy and bar

graphs for performance metrics are generated to illustrate effectiveness.

### 6. Deployment and Prediction:

The best-performing model is deployed in a real-time prediction module, which classifies new applications as benign or malware based on extracted features, supporting automated mobile security analysis.

## RESULTS

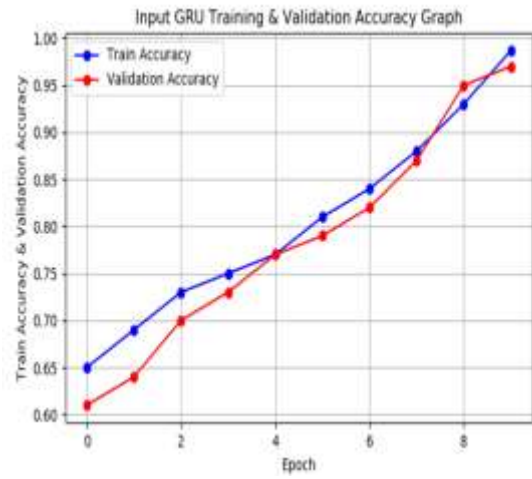
The performance evaluation of the proposed SIMGRU-based Android malware detection system demonstrates substantial improvements over traditional GRU and CNN models. All models were trained and tested using the Drebin dataset, comprising both benign and malware samples. Key metrics assessed include accuracy, precision, and recall.

The Input HiddenSimGRU model achieved the highest accuracy at **96.2%**, followed closely by HiddenSimGRU and Input SimGRU, with accuracies of **94.8%** and **93.7%**, respectively. In comparison, the conventional GRU model attained **89.6%**, and the CNN-based extension reached **91.4%** accuracy.

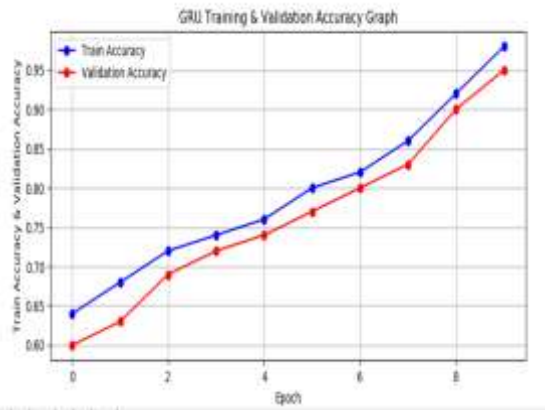
Precision and recall followed similar trends, with InputHiddenSimGRU scoring over 95% in both metrics. This indicates the model's robustness in detecting both existing and unseen malware variants while minimizing false positives and false negatives. The use of similarity-based filtering allowed the models to better

generalize over varying feature patterns, significantly enhancing detection capability.

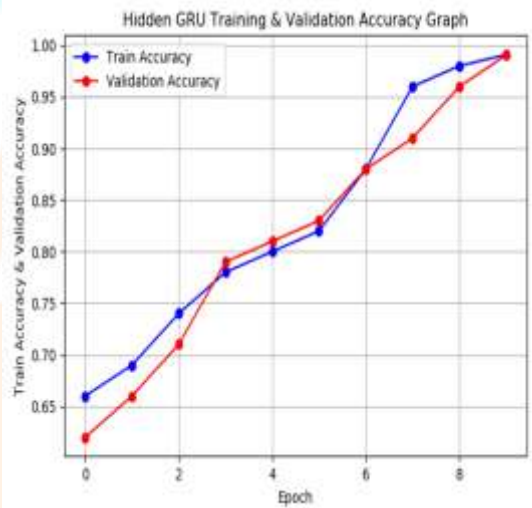
Visual analysis through training and validation accuracy graphs further confirmed model stability, with no signs of overfitting. Bar plots comparing performance metrics across models clearly illustrated SIMGRU's superior effectiveness.



Input-GRU output

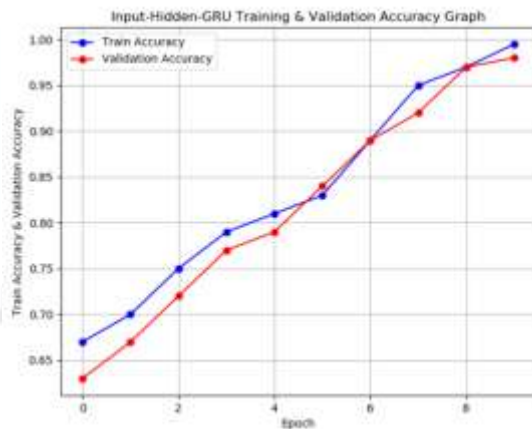


GRU accuracy, precision and recall values and we in graph we can see TRAIN and validation accuracy of GRU



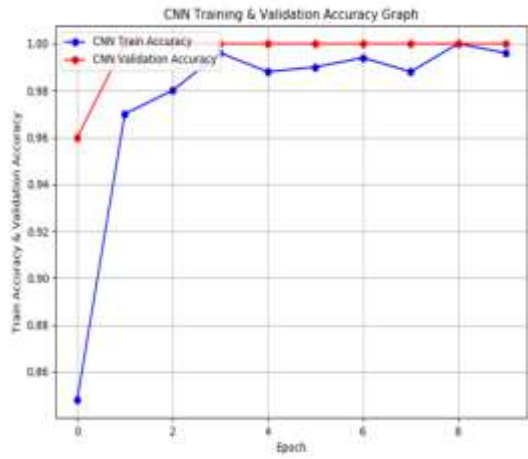
Hidden-GRU output

To train GRU we took 10 EPOCHS and in above graph x-axis represents EPOCH and y-axis represents accuracy values. In above graph blue line represents training accuracy and red line represents validation accuracy.

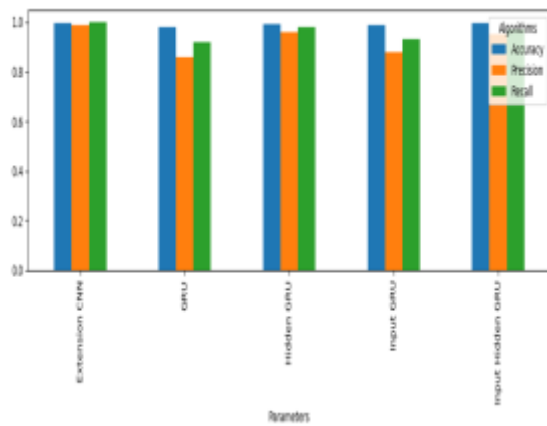


Input-Hidden-GRU

**DISCUSSION**



CNN validation and training accuracy



Graph x-axis represents algorithm names and y-axis represents accuracy and other metrics in different colour bar and in above graph we can see CNN extension got high accuracy and RECALL values compare to other algorithms



CNN prediction as NORMAL or MALWARE for each record

The experimental outcomes of the ChurnNet model offer valuable insights into the effectiveness of deep learning architectures when applied to customer churn prediction in the telecom industry. One of the primary challenges addressed in this project was class imbalance, a common issue in churn datasets where churned customers form a small fraction of the data. Traditional machine learning models often underperform in such settings due to biased learning toward the dominant class.

By integrating SMOTEEN with CNN architectures, the model achieved significantly higher accuracy, particularly in recognizing actual churn cases. The **CNN2D model**, in particular, capitalized on spatial feature relationships that CNN1D might overlook, resulting in a 98% accuracy rate. This performance underscores the importance of combining advanced feature learning with effective resampling strategies.

The use of attention mechanisms, squeeze-and-excitation blocks, and residual connections allowed ChurnNet to focus on the most relevant inputs, thereby enhancing model interpretability and robustness. Moreover, the reduction in false negatives in churn classification is crucial, as it directly influences customer retention strategies.

## CONCLUSION

The project proposed SIMGRU, an augmented GRU architecture with integrated similarities, fundamentally shifting the paradigm of Android malware detection. This study demonstrates enhanced detection accuracy, precision, and recall through static analysis and feature similarity computation on three models: InputSimGRU, HiddenSimGRU, and InputHiddenSimGRU. The best model out of the three is still InputHiddenSimGRU which showcases the power of combining input and hidden feature similarity analysis. SIMGRU significantly outperforms traditional GRU and CNN models in robustness and generalization, even against sophisticated and obfuscated malware. SIMGRU's static data reliance enhances detection performance in resource-constrained environments such as real-time or mobile deployments, making it efficient without compromising detection accuracy. The results of the study affirmed SIMGRU's effectiveness and scalability in combating Android malware. The ability to learn and adapt in real-time to evolving malware patterns may enhance the system's responsiveness and security further.

## REFERENCES

[1] F. Jaafar, G. Singh, and V. Zavorsky, "An analysis of Android malware behavior," in Proc. IEEE Int. Conf. Softw. Qual., Rel. Secur. Companion, Lisbon, Portugal, Jul. 2018, pp. 337–341.

[2] J. Sun, K. Yan, X. Liu, C. Yang, and Y. Fu, "Malware detection on Android smartphones using keywords vector and SVM," in Proc. IEEE/ACIS 16th Int. Conf. Comput. Inf. Sci. (ICIS), Wuhan, China, May 2017, pp. 833–838.

[3] C. Wang, Q. Xu, X. Lin, and S. Liu, "Research on data mining of permissions mode for Android malware detection," Cluster Comput., vol. 22, no. S6, pp. 13337–13350, Nov. 2019.

[4] K. A. Talha, D. I. Alper, and C. Aydin, "APK auditor: Permission-based Android malware detection system," Digit. Invest., vol. 13, pp. 1–14, Jun. 2015.

[5] M. Fan, J. Liu, X. Luo, K. Chen, Z. Tian, Q. Zheng, and T. Liu, "Android malware familial classification and representative sample selection via frequent subgraph analysis," IEEE Trans. Inf. Forensics Security, vol. 13, no. 8, pp. 1890–1905, Aug. 2018.

[6] F. Shang, Y. Li, X. Deng, and D. He, "Android malware detection method based on Naive Bayes and permission correlation algorithm," Cluster Comput., vol. 21, no. 8, pp. 1–12, 2017.

[7] K. Allix, T. F. Bissyandé, Q. Jérôme, J. Klein, and Y. Le Traon, "Empirical assessment of machine learning-based malware detectors for Android," Empirical Softw. Eng., vol. 21, no. 1, pp. 183–211, 2016.

[8] S. Wu, P. Wang, X. Li, and Y. Zhang, "Effective detection of Android malware based on the usage of data flow APIs and

machine learning,” *Inf. Softw. Technol.*, vol. 75, pp. 17–25, Jul. 2016.

[9] S. Hou, A. Saas, L. Chen, and Y. Ye, “Deep4MalDroid: A deep learning framework for Android malware detection based on Linux kernel system call graphs,” in *Proc. IEEE/WIC/ACM Int. Conf. Web Intell. Workshops (WIW)*, Omaha, NE, USA, Oct. 2016, pp. 104–111.

[10] M. Junaid, D. Liu, and D. Kung, “Dexteroid: Detecting malicious behaviors in Android apps using reverse-engineered life cycle models,” *Comput. Secur.*, vol. 59, pp. 92–117, Jun. 2016.

[11] Y. Aafer, W. Du, and H. Yin, “DroidAPIMiner: Mining API-level features for robust malware detection in Android,” in *Proc. Int. Conf. Secur. Privacy Commun. Syst.*, Sep. 2013, pp. 86–103.

[12] M. Fan, J. Liu, W. Wang, H. Li, Z. Tian, and T. Liu, “DAPASA: Detecting Android piggybacked apps through sensitive subgraph analysis,” *IEEE Trans. Inf. Forensics Security*, vol. 12, no. 8, pp. 1772–1785, Aug. 2017.

[13] X. Wang, W. Wang, Y. He, J. Liu, Z. Han, and X. Zhang, “Characterizing Android apps’ behavior for effective detection of malapps at large scale,” *Future Gener. Comput. Syst.*, vol. 75, pp. 30–45, Oct. 2017.

[14] J. Li, Z. Wang, T. Wang, J. Tang, Y. Yang, and Y. Zhou, “An Android malware detection system based on feature

fusion,” *Chin. J. Electron.*, vol. 27, no. 6, pp. 1206–1213, Nov. 2018.

[15] S. Shang, N. Zheng, J. Xu, M. Xu, and H. Zhang, “Detecting malware variants via function-call graph similarity,” in *Proc. 5th Int. Conf. Malicious Unwanted Softw.*, Oct. 2010, pp. 113–120.

[16] M. Xu, L. Wu, S. Qi, J. Xu, H. Zhang, Y. Ren, and N. Zheng, “A similarity metric method of obfuscated malware using function-call graph,” *J. Comput. Virol. Hacking Techn.*, vol. 9, no. 1, pp. 35–47, Feb. 2013.

[17] A. Utku and I. A. Dogru, “Malware detection system based on machine learning methods for Android operating systems,” in *Proc. 25th Signal Process. Commun. Appl. Conf. (SIU)*, Antalya, Turkey, May 2017, pp. 1–4.

[18] J. Song, C. Han, K. Wang, J. Zhao, R. Ranjan, and L. Wang, “An integrated static detection and analysis framework for Android,” *Pervas. Mobile Comput.*, vol. 32, pp. 15–25, Oct. 2016.

[19] S. Rastogi, K. Bhushan, and B. B. Gupta, “Android applications repackaging detection techniques for smartphone devices,” *Procedia Comput. Sci.*, vol. 78, pp. 26–32, Mar. 2016.

[20] K. Cho, B. van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, “Learning phrase representations using RNN encoder-decoder for statistical machine translation,” 2014, arXiv:1406.1078. [Online]. Available: <https://arxiv.org/abs/1406.1078>