

# Comparative Analysis of Phishing Detection in Ethereum Cryptocurrency Transactions Using SMOTE-Based Random Forest and LightGBM Algorithms

Ihsan Maulana<sup>1</sup>, Zulkifli<sup>2</sup>

<sup>1,2</sup>Informatics Engineering Study Program, Pamulang University, Kota Tangerang Selatan, Indonesia.  
Email : ihsanmaulana.app@gmail.com

The rapid growth of cryptocurrency adoption has increased the risk of phishing attacks targeting blockchain transactions, particularly within the Ethereum ecosystem. This study aims to conduct a comparative analysis of machine learning approaches for detecting phishing transactions by addressing the challenge of highly imbalanced data. A supervised classification framework is developed using Random Forest and Light Gradient Boosting Machine algorithms, combined with the Synthetic Minority Over sampling Technique to enhance minority class representation. The research process involves data preprocessing, feature normalization, class balancing, model training, and performance evaluation using appropriate classification metrics. The findings indicate that ensemble based methods are effective in capturing complex transaction patterns and distinguishing legitimate from phishing activities. Random Forest demonstrates more robust and consistent performance compared to Light Gradient Boosting Machine, particularly in minimizing misclassification of fraudulent transactions. These results highlight the importance of data balancing techniques and algorithm selection in improving blockchain security. The study contributes to the development of reliable phishing detection models that can support safer cryptocurrency transaction environments.

**Keywords:** Phishing Detection, Cryptocurrency Transactions, Ethereum, Machine Learning, Data Imbalance, Ensemble Algorithms

This is an open access article under the [CC BY-NC](#) license



## Corresponding Author:

Ihsan Maulana

Informatics Engineering Study Program, Pamulang University

Jl. Raya Puspitek, Buaran, Kec. Pamulang, Kota Tangerang Selatan, Banten 15310

ihsanmaulana.app@gmail.com

## 1. Introduction

Cryptocurrency has evolved from an experimental concept into a significant global financial instrument [1]. Ethereum, as the largest smart contract platform, serves as the backbone for thousands of decentralized applications and digital tokens. However, the pseudo anonymous nature, decentralization, and irreversibility of blockchain transactions act as a double edged sword; while providing privacy and autonomy, they also create fertile ground for cybercriminal activities, particularly phishing based fraud [2]

Recent cybersecurity reports indicate a sharp increase in attacks targeting crypto wallets, in which attackers deceive victims into signing malicious transactions or surrendering their private keys. The resulting losses not only affect individual finances but also undermine public trust in the Decentralized Finance ecosystem. Therefore, the need for early detection systems capable of identifying malicious transaction patterns before confirmation has become increasingly critical [3]

The main technical challenge in developing such detection systems lies in data imbalance. Among millions of transactions occurring on the Ethereum network, the proportion of phishing transactions is extremely small, often less than one percent compared to legitimate transactions [4]. Conventional machine learning algorithms tend to be biased toward the majority class, leading to a high rate of false negatives or failures to detect actual attacks [5]

This study proposes a machine learning based approach by employing the Synthetic Minority Over sampling Technique to balance the transaction data distribution [6]. It compares the effectiveness of two popular ensemble algorithms: Random Forest, known for its robustness in reducing variance through bagging, and LightGBM, recognized for its computational efficiency through the Gradient based One Side Sampling technique [7]. The primary focus of this research is to evaluate the extent to which the use of SMOTE can enhance model sensitivity to the minority class while maintaining accurate classification of normal transactions.

## 2. Literature Review and Problem Statement

### Literature Review

The rapid growth of blockchain-based financial ecosystems, particularly Ethereum, has significantly transformed digital transactions by offering decentralization, transparency, and trustless execution through smart contracts [8]. However, these advantages have simultaneously introduced new vectors for cybercrime, especially phishing attacks that exploit user behavior, wallet vulnerabilities, and transaction-level manipulation. Unlike traditional phishing, blockchain phishing often results in irreversible financial losses due to the immutability of transactions, making early and accurate detection critically important [9].

Prior studies have increasingly focused on leveraging machine learning techniques to detect anomalous and fraudulent behavior in blockchain transactions. Provide a comprehensive survey on blockchain anomaly detection, highlighting the effectiveness of supervised and ensemble learning methods in identifying abnormal transaction patterns [10]. Their work emphasizes that transaction-based features, such as value transfer, address interaction frequency, and temporal behavior, are particularly informative for detecting fraud-related activities. Nevertheless, the study also underlines persistent challenges related to class imbalance and scalability in real-world blockchain datasets.

Specific to Ethereum phishing, investigate phishing scams by analyzing transaction behavior patterns on the Ethereum network [11]. Their findings demonstrate that machine learning classifiers can effectively distinguish phishing-related transactions from legitimate ones when behavioral features are properly engineered. However, their approach does not explicitly address severe data imbalance, which is a common characteristic of phishing datasets where fraudulent transactions represent only a very small fraction of total observations. This limitation potentially affects model generalizability and real-world deployment.

To address class imbalance, several studies have incorporated resampling techniques such as the Synthetic Minority Oversampling Technique (SMOTE) [12].) SMOTE in combination with LightGBM for real-time detection of Ponzi schemes on Ethereum. Their results show significant performance improvements, particularly in recall and F1-score, indicating that SMOTE can enhance a model's ability to identify rare fraudulent events. Despite these promising results, their study focuses on Ponzi schemes rather than phishing attacks and evaluates only a single ensemble algorithm, limiting comparative insights.

Ensemble learning methods, such as Random Forest and gradient boosting frameworks, have been widely recognized for their robustness and high predictive accuracy in fraud detection contexts. Demonstrates that ensemble-based models outperform single classifiers in decentralized finance (DeFi) fraud detection due to their ability to capture complex nonlinear relationships and reduce overfitting [13]. Nonetheless, the study primarily addresses DeFi protocols at a macro level and does not specifically evaluate transaction-level phishing detection on Ethereum.

In the context of traditional financial fraud, conduct a comparative analysis of Random Forest and XGBoost using SMOTE to handle imbalanced credit card fraud datasets [14]. Their findings reveal that ensemble models combined with SMOTE significantly improve classification performance. Although conducted

outside the blockchain domain, this study provides empirical support for the effectiveness of combining ensemble learning with resampling techniques in highly imbalanced fraud detection problems.

Despite the growing body of literature, several gaps remain evident. First, comparative studies that directly evaluate multiple ensemble algorithms for Ethereum phishing detection under identical experimental conditions are still limited. Second, while SMOTE has been widely adopted, its specific impact on different ensemble models, particularly Random Forest versus LightGBM, has not been sufficiently examined in the context of blockchain phishing. Third, many existing studies prioritize overall accuracy or ROC-AUC without adequately discussing trade-offs between precision and recall, which are crucial in fraud detection scenarios where false negatives can lead to substantial financial losses [15].

### **Problem Statement**

Although machine learning-based approaches have shown strong potential in detecting fraudulent activities on blockchain networks, the detection of phishing transactions on Ethereum remains a challenging problem. One of the primary challenges lies in the highly imbalanced nature of blockchain transaction datasets, where phishing transactions constitute a very small minority compared to legitimate transactions [9]. This imbalance often biases classification models toward the majority class, resulting in poor detection of actual phishing activities.

Furthermore, existing studies on Ethereum phishing detection tend to focus on single-model approaches or lack systematic comparisons between different ensemble algorithms under the same preprocessing and evaluation framework. As a result, it remains unclear which ensemble method offers superior performance, robustness, and practical suitability when combined with imbalance-handling techniques such as SMOTE.

Another unresolved issue concerns the trade-off between detection performance and computational efficiency [16]. While advanced ensemble models may achieve high accuracy, they may also require longer training times or higher computational resources, which can limit their applicability in real-time or large-scale blockchain monitoring systems.

Based on these considerations, this study addresses the following research problems:

- a. How does the application of SMOTE affect the performance of ensemble learning models in detecting phishing transactions on the Ethereum blockchain?
- b. Is there a significant performance difference between Random Forest and LightGBM when applied to Ethereum phishing detection under an imbalanced data setting?
- c. Which model provides a better balance between detection effectiveness (precision, recall, F1-score, and ROC-AUC) and computational efficiency for practical deployment?

By addressing these problems, this research aims to contribute empirical evidence to the blockchain cybersecurity literature and provide practical insights for selecting appropriate machine learning models for phishing detection in Ethereum-based transaction systems.

### **3. Method**

This study adopts a quantitative research design using a supervised machine learning approach to detect phishing activities in Ethereum cryptocurrency transactions. The methodological framework is structured to ensure systematic data processing, model development, and performance evaluation. The overall workflow consists of data collection, preprocessing, data balancing, model training, and comparative evaluation. This design allows for an objective assessment of the effectiveness of different ensemble algorithms in handling imbalanced blockchain transaction data.

The dataset used in this study is obtained from a publicly available Ethereum transaction repository, which contains labeled records representing legitimate and fraudulent activities. Each transaction is described by multiple numerical features that capture transaction behavior within the Ethereum network. Prior to analysis, irrelevant attributes are removed to improve computational efficiency and reduce noise, while ensuring that all retained features are theoretically and empirically relevant to phishing detection.

Data preprocessing is conducted to enhance data quality and model performance. This stage includes handling missing values, normalizing feature scales using a min-max transformation, and ensuring consistent data formatting. Normalization is applied to prevent features with larger numerical ranges from dominating the learning process, particularly in tree-based ensemble models that may still be influenced by feature distribution.

Given the highly imbalanced nature of phishing transaction data, a data balancing strategy is essential. This study employs the Synthetic Minority Over-sampling Technique to address class imbalance by generating synthetic samples of the minority class based on feature space similarities. The oversampling process is applied exclusively to the training data to avoid information leakage and to preserve the integrity of model evaluation.

Following preprocessing and data balancing, the dataset is divided into training and testing subsets using a stratified splitting strategy. Stratification ensures that the class distribution in the testing data reflects the original dataset, enabling a fair assessment of model generalization. The training subset is used to build the predictive models, while the testing subset is reserved solely for performance evaluation.

Two ensemble learning algorithms are implemented and compared in this study, namely Random Forest and Light Gradient Boosting Machine. Random Forest is selected due to its robustness in handling noisy data and its ability to reduce overfitting through bootstrap aggregation. Light Gradient Boosting Machine is chosen for its computational efficiency and capability to model complex nonlinear relationships using gradient boosting techniques.

Model performance is evaluated using classification metrics that are suitable for imbalanced data scenarios, including precision, recall, F score, and the area under the receiver operating characteristic curve. These metrics provide a comprehensive understanding of each model's ability to correctly identify phishing transactions while maintaining reliable performance on legitimate transactions. The comparative analysis focuses on identifying the model that offers the most balanced and reliable detection capability for practical implementation in blockchain security systems.

## 4. Results And Discussion

### Data Distribution Analysis

Before the application of SMOTE, the dataset exhibited extreme class imbalance. After SMOTE was applied to the training data, the class distribution became balanced with a ratio of 1:1, providing a fair opportunity for the algorithms to learn the characteristics of both classes (Normal and Fraud).

### Random Forest Classification Results

The evaluation of the Random Forest model on the testing data, which was not affected by SMOTE, demonstrates highly impressive results.

- a. Precision: 0.98
- b. Recall: 0.96
- c. F1 Score: 0.97
- d. ROC AUC: 0.992

The high ROC AUC value of 0.992 confirms that the Random Forest model has excellent separability in distinguishing between legitimate transactions and phishing attempts. The minimal occurrence of false positives makes this model particularly safe for deployment, as it reduces the risk of unnecessary transaction blocking that could negatively affect user experience.

### LightGBM Classification Results

The LightGBM model also exhibits strong performance, although slightly lower than Random Forest in terms of precision.

- a. Precision: 0.95
- b. Recall: 0.94
- c. F1 Score: 0.945
- d. ROC AUC: 0.988

The primary advantage of LightGBM lies in its training time, which is approximately 40% faster than that of Random Forest. However, the model shows a slightly higher tendency to produce false positives.

### Comparative Analysis

Table 1 presents a head to head comparison between the two algorithms.

**Table 1.** Comparison of Model Performance for the Fraud Class

Algorithm	Precision	Recall	F Score	ROC AUC	Stability
Random Forest	0.98	0.96	0.97	0.992	High
LightGBM	0.95	0.94	0.945	0.988	Moderate

Based on the comparative analysis, Random Forest is more strongly recommended for financial fraud detection scenarios in which the cost of prediction errors is extremely high, particularly in cases involving undetected fraud or the incorrect blocking of legitimate users. The voting mechanism employed by Random Forest demonstrates greater robustness to noise compared to the more aggressive boosting strategy used by LightGBM within this dataset.

## 5. Conclusion

This study provides a comprehensive comparative analysis of machine learning based approaches for detecting phishing activities in Ethereum cryptocurrency transactions under highly imbalanced data conditions. By integrating a data balancing strategy with ensemble learning algorithms, the research demonstrates that appropriate preprocessing and model selection play a crucial role in improving fraud detection performance within blockchain environments. The application of synthetic oversampling enables the models to better capture the behavioral patterns of fraudulent transactions, which are often underrepresented in real world datasets. The findings indicate that ensemble methods are particularly effective for this classification task, as they are capable of modeling complex and nonlinear transaction characteristics. Among the evaluated approaches, the Random Forest algorithm exhibits superior robustness and consistency in distinguishing between legitimate and fraudulent transactions. Its aggregation based learning mechanism reduces sensitivity to noise and minimizes misclassification risks, making it well suited for security critical financial applications. Light Gradient Boosting Machine also delivers strong performance and offers advantages in computational efficiency, although it shows a slightly higher tendency to generate false alerts. Overall, this study highlights the importance of balancing accuracy, sensitivity, and reliability when designing phishing detection systems for cryptocurrency platforms. The proposed framework contributes valuable insights for practitioners and researchers seeking to enhance blockchain security and can serve as a foundation for future work involving real time detection systems and the integration of additional behavioral and network based features.

## 6. References

- [1] S. Saksonova and I. Kuzmina-Merlino, "Cryptocurrency as an investment instrument in a modern financial market," *Вестник Санкт-Петербургского университета. Экономика*, vol. 35, no. 2, pp. 269–282, 2019.
- [2] J. Yang *et al.*, "Who stole my nft? investigating web3 nft phishing scams on ethereum," *IEEE Trans. Inf. Forensics Secur.*, 2024.
- [3] Y. Qi, J. Wu, H. Xu, and M. Guizani, "Blockchain data mining with graph learning: A survey," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 46, no. 2, pp. 729–748, 2023.
- [4] Y. Wan, F. Xiao, and D. Zhang, "Early-stage phishing detection on the Ethereum transaction network," *Soft Comput.*, vol. 27, no. 7, pp. 3707–3719, 2023.
- [5] S. Siddique, M. A. Haque, R. George, K. D. Gupta, D. Gupta, and M. J. H. Faruk, "Survey on machine learning biases and mitigation techniques," *Digital*, vol. 4, no. 1, pp. 1–68, 2023.
- [6] F. A. Ghaleb, F. Saeed, M. Al-Sarem, S. N. Qasem, and T. Al-Hadhrani, "Ensemble synthesized minority oversampling-based generative adversarial networks and random forest algorithm for credit card fraud detection," *IEEE Access*, vol. 11, pp. 89694–89710, 2023.
- [7] R. Bakir, C. Orak, and A. Yüksel, "Optimizing hydrogen evolution prediction: A unified approach using random forests, lightGBM, and Bagging Regressor ensemble model," *Int. J. Hydrogen Energy*, vol. 67, pp. 101–110, 2024.
- [8] A. Parhi, "Blockchain in finance: Applications, platforms, and global trends in a decentralizing ecosystem," *Platforms, Glob. Trends a Decentralizing Ecosyst. (June 06, 2025)*, 2025.
- [9] R. O. Ogundokun, M. O. Arowolo, R. Damaševičius, and S. Misra, "Phishing detection in blockchain transaction networks using ensemble learning," in *Telecom*, MDPI, 2023, pp. 279–297.
- [10] S. Hisham, M. Makhtar, and A. A. Aziz, "Combining multiple classifiers using ensemble method for anomaly detection in blockchain networks: A comprehensive review," *Int. J. Adv. Comput. Sci. Appl.*, vol. 13, no. 8, 2022.
- [11] M. Ghosh, D. Ghosh, R. Halder, and J. Chandra, "Investigating the impact of structural and temporal behaviors in Ethereum phishing users detection," *Blockchain Res. Appl.*, vol. 4, no. 4, p. 100153, 2023.
- [12] M. Carvalho, A. J. Pinho, and S. Brás, "Resampling approaches to handle class imbalance: a review from a data perspective," *J. Big Data*, vol. 12, no. 1, p. 71, 2025.
- [13] A. A. Alhashmi, A. M. Alashjaee, A. A. Darem, A. F. Alanazi, and R. Effghi, "An ensemble-based fraud detection model for financial transaction cyber threat classification and countermeasures," *Eng. Technol. Appl. Sci. Res.*, vol. 13, no. 6, pp. 12433–12439, 2023.
- [14] M. Riham, "A Comparative Study of Random Forest and XGBoost for Detecting Credit Card Fraud Transactions using Big Data," *BCAS Campus*, 2025.
- [15] C. Müller-Bloch and J. Kranz, "A framework for rigorously identifying research gaps in qualitative literature reviews," 2015.
- [16] J. Huang *et al.*, "Speed/accuracy trade-offs for modern convolutional object detectors," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 7310–7311.