

Evaluating the Vulnerability of ML-Based Ethereum Phishing Detectors to Single-Feature Adversarial Perturbations

AHOD ALGHURIED, University of Central Florida, USA
 ALI ALKINOON, University of Central Florida, USA
 ABDULAZIZ ALGHAMDI, University of Central Florida, USA
 SOOHYEON CHOI, University of Central Florida, USA
 MANAR MOHAISEN, Northeastern Illinois University, USA
 DAVID MOHAISEN*, University of Central Florida, USA

This paper explores the vulnerability of machine learning models to simple single-feature adversarial attacks in the context of Ethereum fraudulent transaction detection. Through comprehensive experimentation, we investigate the impact of various adversarial attack strategies on model performance metrics. Our findings, highlighting how prone those techniques are to simple attacks, are alarming, and the inconsistency in the attacks' effect on different algorithms promises ways for attack mitigation. We examine the effectiveness of different mitigation strategies, including adversarial training and enhanced feature selection, in enhancing model robustness and show their effectiveness.

Additional Key Words and Phrases: Ethereum, Adversarial Attacks, Machine Learning, Phishing Detection

ACM Reference Format:

Ahod Alghuried, Ali Alkinoon, Abdulaziz Alghamdi, Soohyeon Choi, Manar Mohaisen, and David Mohaisen. 2025. Evaluating the Vulnerability of ML-Based Ethereum Phishing Detectors to Single-Feature Adversarial Perturbations. 1, 1 (April 2025), 24 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Machine learning algorithms are being heavily employed in many applications, and issues concerning their robustness and security are becoming increasingly important and concerning. In particular, adversarial attacks on machine learning algorithms pose a serious risk where small perturbations are added to input data leading to misclassifications that undermine the reliability of these algorithms [40, 41]. Understanding the vulnerabilities of these models to adversarial attacks is an essential first step toward the development of robust and dependable machine learning-powered systems, especially in security applications.

One of the applications of machine learning algorithms has been detecting fraud in Ethereum, including phishing and fake initial coin offerings (ICOs). In Ethereum, phishing involves an adversary posing as an honest user to deceive benign users into revealing sensitive information or transferring

*Corresponding author: David Mohaisen

Authors' addresses: Ahod Alghuried, ah104940@ucf.edu, University of Central Florida, Florida, USA; Ali Alkinoon, alialkinoon@ucf.edu, University of Central Florida, Florida, USA; Abdulaziz Alghamdi, abdulaziz.alghamdi@ucf.edu, University of Central Florida, Florida, USA; Soohyeon Choi, soohyeon.choi@ucf.edu, University of Central Florida, Florida, USA; Manar Mohaisen, m-mohaisen@neiu.edu, Northeastern Illinois University, Chicago, USA; David Mohaisen, mohaisen@ucf.edu, University of Central Florida, Florida, USA.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2025 Association for Computing Machinery.

XXXX-XXXX/2025/4-ART \$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

digital currency through deceptive means such as fake websites or smart contracts. These phishing attempts exploit the decentralized and anonymous nature of the blockchain, making it challenging to trace and recover lost assets once the attack is executed [53]. In this context, machine learning classifiers that distinguish between benign and phishing attempts are shown to be effective [6]. However, despite such effectiveness, various studies have shown issues with the robustness of these algorithms, including adversarial attacks, where research efforts have introduced sophisticated manipulation (adversarial) attacks that lead to misclassification of phishing attempts, thus fooling benign users [1–5, 10, 23] (more in section 2).

This study explores how adversarial attacks affect machine learning models used for detecting fraudulent transactions, focusing on Random Forest (RF), Decision Tree (DT), and K-Nearest Neighbors (KNN) classifiers. By applying the Fast Gradient Sign Method (FGSM) [7, 47], a widely known adversarial attack technique, we examine how these models degrade in performance when exposed to subtle yet deceptive adversarial inputs. Beyond measuring accuracy loss, we also analyze the broader implications of these attacks, particularly their potential to disrupt automated systems that rely on these models. A single misclassification could lead to incorrect decisions, amplifying the overall risk. Our goal is to understand the vulnerabilities of these models to adversarial manipulation and explore strategies to strengthen their defenses, ensuring greater robustness against such attacks.

Recent studies have pointed out that machine learning models can be vulnerable to adversarial attacks, a problem that is not just limited to certain algorithms or use cases but is a widespread issue across the field [36]. Our study builds on this understanding by providing an in-depth analysis of how these attacks affect models designed to detect fraud and offering practical solutions to lessen their impact. Our result indicate that not all algorithms handle such manipulations equally highlighting the critical need to carefully select the models and to design robust defense strategies that are tailored to the specific needs of the application. This research also brings to light the ongoing necessity to continually test and update machine learning models, ensuring they can stand up to the ever evolving threats in dynamic environments.

2 LITERATURE REVIEW

We review works that focus on detecting and mitigating threats within cryptocurrency networks, particularly focusing on Ethereum and Bitcoin. As shown in Table 1, the review covers key areas, including the detection of malicious transactions, the vulnerabilities of machine learning models to adversarial attacks, and the role of Generative Adversarial Networks (GANs) in enhancing detection capabilities. Additionally, it highlights recent innovations in fraud and phishing detection.

2.1 Detection and Analysis of Malicious Transactions

Malicious Transactions in Ethereum. Agarwal *et al.* [8] analyzed malicious Ethereum transactions, revealing how subtle data manipulations can mislead advanced machine learning models, posing security challenges. Rabieinejad *et al.* [37] leveraged Ethereum data to improve cyber threat detection, stressing the need for representative datasets. Sanjalawe *et al.* [42] utilized the Benchmark Labeled Transactions Ethereum dataset to detect illicit activities, highlighting the role of feature engineering and semi-supervised learning in improving accuracy.

Entity Classification in Cryptocurrency Networks. Zola *et al.* [57] used WalletExplorer data to classify Bitcoin addresses, reducing anonymity and helping to identify illicit entities while highlighting the need for privacy-preserving techniques. Yang *et al.* [52] used transaction data to detect Bitcoin spam attacks, demonstrating the effectiveness of GRU-based models.

Table 1. Summary of some of the prior work.

Paper Title	Year	Adversarial Techniques	Applications in Cryptocurrency
Li <i>et al.</i> [27]	2018	GANs	Anomaly detection, Secure Water Treatment
Qingyu <i>et al.</i> [24]	2019	IFCM, AIS, R3	Fraud detection, TaoBao
Ba <i>et al.</i> [12]	2019	GANs	Credit card fraud, 31-feature dataset
Ngo <i>et al.</i> [34]	2019	GANs	Anomaly detection, MNIST, CIFAR10
Zola <i>et al.</i> [56]	2020	GANs, data augmentation	Bitcoin, WalletExplorer (categorized addresses)
Yang <i>et al.</i> [52]	2020	WGAN-div, GRU-based detection	Bitcoin, custom spam transaction dataset
Shu <i>et al.</i> [43]	2020	GANs	Intrusion detection, network traffic
Mozo <i>et al.</i> [32]	2021	WGANs, synthetic traffic	Monero, custom cryptomining dataset
Fursov <i>et al.</i> [21]	2021	Black-box attacks	Transaction records
Agarwal <i>et al.</i> [8]	2022	CTGAN, K-Means Clustering	Ethereum, Etherscan dataset (2,946 malicious accounts)
Fidalgo <i>et al.</i> [19]	2022	GANs, data augmentation	Bitcoin, Elliptic dataset (200K+ transactions)
Zola <i>et al.</i> [57]	2022	Various GANs, adversarial learning	Bitcoin, WalletExplorer (16M+ addresses)
Rabieinejad <i>et al.</i> [37]	2023	CTGAN, WGAN	Ethereum, 57K normal, 14K abnormal transactions
Sanjalawe <i>et al.</i> [42]	2023	GANs, feature extraction	Ethereum, labeled transactions (normal, abnormal)

Detection of Cryptomining Attacks. Mozo *et al.* [32] applied transaction datasets to detect cryptomining attacks within the Monero network. Their research highlighted that leveraging synthetic network traffic data generated through advanced GAN architectures enables high-precision detection of cryptomining activities, even in privacy-focused networks like Monero.

2.2 Vulnerabilities of ML Models to Adversarial Attacks

As cryptocurrencies like Ethereum gain traction, securing ML models for transaction classification is crucial for detecting fraud and predicting outcomes [44, 45]. However, their vulnerability to adversarial attacks remains a major concern [39]. Chen *et al.* [49], Li *et al.* [29], and Oliveira *et al.* [35] showed how manipulated inputs deceive ML models, threatening classification accuracy. Narodytska and Kasiviswanathan [33] highlighted CNN susceptibility, while Cartella *et al.* [15] optimized fraud detection on imbalanced data, achieving a perfect attack success rate. Bhagoji *et al.* [13] introduced Gradient Estimation black-box attacks, achieving near-perfect success on DNNs.

2.3 GANs Enhanced Detection and Defense Strategies

Data Augmentation and Synthetic Data Generation. Generative Adversarial Networks (GANs) generate realistic synthetic data, addressing labeled data scarcity and enhancing ML model training. They are instrumental in creating Adversarial Examples (AEs) and handling data perturbations in cryptocurrency transactions. Fidalgo *et al.* [19] leveraged GANs to mitigate class imbalance in Bitcoin, improving model performance. Agarwal *et al.* [8] used Conditional GANs (CTGAN) to generate adversarial Ethereum data, enhancing model robustness.

Enhancing Detection Capabilities with GANs. Rabieinejad *et al.* [37] employed CTGAN and Wasserstein GANs (WGAN) to augment Ethereum transaction datasets, improving cyber threat detection. Sanjalawe *et al.* [42] used Semi-Supervised GANs and feature extraction for Ethereum dataset perturbation. Zola *et al.* [56, 57] applied GAN-based augmentation to enrich the underrepresented Bitcoin transaction classes.

Adversarial Examples for Robust Detection. Yang *et al.* [52] utilized WGAN-div to generate adversarial examples for Bitcoin spam detection, ensuring stable training and high-quality synthetic data. Mozo *et al.* [32] employed WGANs to create synthetic network traffic for Monero cryptomining attack detection, enhancing accuracy with realistic data.

2.4 Fraud and Phishing Detection in Ethereum

Fraud Detection Techniques. Cartella *et al.* [15] explores how adversarial attacks can be adapted for fraud detection systems dealing with imbalanced data, showing that these attacks can successfully bypass AI models while remaining hard to detect. Ravindranath *et al.* [38] investigates ensemble machine learning models such as CATBoost and LGBM to detect fraud in Ethereum, emphasizing how oversampling techniques improve model accuracy and robustness. Kabla *et al.* [25] reviews the applicability of Intrusion Detection Systems (IDS) in detecting attacks on Ethereum-based Decentralized Applications (DApps), discussing vulnerabilities, existing detection methods, and future directions.

Phishing Detection in Ethereum Networks. Tan *et al.* [48] proposes using Graph Convolutional Networks to detect fraudulent transactions by leveraging network embeddings derived from Ethereum transaction records. Yang *et al.* [53] introduces a phishing detection method that enhances interpretability by extracting detailed features from transaction networks, improving detection precision and clarity. Chen *et al.* [17] describes a hybrid graph neural network model combined with data augmentation to detect phishing scams on Ethereum, showing superior results by integrating temporal and structural features. Luo *et al.* [30] presents a model using the bias2vec network embedding algorithm to effectively classify Ethereum accounts as phishing or benign.

Graph-Based Phishing Detection Methods. Lv *et al.* [31] proposes a graph-based method using *imgraph2vec* for phishing detection, showing improved feature extraction and classification performance. Yin *et al.* [54] introduces a community-enhanced graph neural network model that improves phishing detection by analyzing the community structure within Ethereum transaction networks. Wu *et al.* [50] proposes the *trans2vec* algorithm to detect phishing scams on Ethereum by embedding features from transaction records, demonstrating effectiveness in classification tasks.

Research Gap. Despite advancements in machine learning and Despite significant progress in machine learning and blockchain technologies, we still face challenges in defending against adversarial attacks and new threats within cryptocurrency networks. Previous studies have tackled issues like fraud detection and categorizing transactions. However the problem of adversarial examples (AEs) specifically designed to exploit weaknesses unique to cryptocurrencies has not been thoroughly examined. This is especially true for AEs that manipulate features to target transaction fraud, smart contract vulnerabilities, and different types of attacks, which all require more in-depth study.

Our research fills this gap by testing the resilience of machine learning based phishing detection algorithms against *simple manipulations*. Focusing on subtle but realistic changes to features to highlight underlying weaknesses in the models and to develop practical, effective countermeasures. Our method involves assessing how vulnerable these systems are to adversarial and identifying which algorithms can best withstand such challenges through a detailed comparative analysis.

3 RESEARCH QUESTIONS

This research examines how well machine learning algorithms can detect phishing activities in Ethereum transactions, especially when facing adversarial threats. Since securing these transactions is essential, it is important to assess how different models perform under such attacks and explore ways to strengthen their reliability. This study aims to answer key questions that highlight major challenges and emphasize the need to incorporate insights from previous studies into our analysis.

RQ1. Are machine learning-based phishing detection algorithms robust against simple feature manipulations? ML models are vulnerable to adversarial attacks, where minor input modifications can alter classification outcomes [13, 15, 33, 44, 49]. Small changes in transaction features, such as amounts or timestamps, may enable adversaries to evade detection. Evaluating

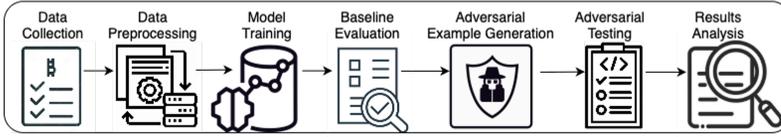


Fig. 1. Pipeline in Ethereum transactions and adversarial testing.

model robustness against such manipulations is essential to ensure reliable Ethereum phishing detection.

RQ2. How do different machine learning algorithms compare in robustness against adversarial manipulations in phishing detection? ML algorithms vary in resilience to adversarial attacks [8, 18, 19, 35, 37, 45]. Some models withstand adversarial manipulations better than others [20, 27, 46]. A comparative analysis is crucial to identifying the most robust algorithms for securing Ethereum transaction classification.

RQ3. How can adversarial manipulations be mitigated in ML-based Ethereum phishing detection? Defensive strategies are needed to counter adversarial attacks [14, 22, 32, 51, 52]. Mitigation techniques include enhanced feature selection, data augmentation, and adversarial training [29, 37, 56]. Strengthening these defenses improves detection accuracy and enhances Ethereum security.

4 METHODOLOGY

The components of our pipeline shown in Figure 1 are data collection, data preprocessing, modeling training, baseline evaluation, AEs testing, and results and analysis. In the following, we review these components in more detail.

4.1 Data Preparation

4.1.1 Dataset Description. We used two datasets in our analysis. The first dataset, described by Kabla et al. [26], is implemented for binary classification, phishing and legitimate transaction. It contains different features associated with each transaction, such as TxHash, a unique number that the block assigns to each transaction, and BlockHeight, which indicates the position where the transaction was recorded. The dataset also has TimeStamp which indicates the time when the transaction was approved and written into the blockchain. In addition, it contains information about the sender and receiver as the From and To address fields, storing the involved parties' Ethereum addresses. The Value attribute indicates how much Ether was transferred in each transaction, and ContractAddress specifies any smart contract related to the transaction. Moreover, the Input captures any extra data in the transaction. The dataset labels each transaction under the Class feature, marking them either phishing (1) or benign (0). In total, this dataset comprises 23,472 transactions, with 15,989 categorized as benign and 7,483 as phishing.

The second dataset is intended for multi-class classification, categorizing transactions into Phishing, Scam, Fake ICO, or Benign [9]. The features of this dataset include hash, a unique identifier for each transaction, and nonce, a counter ensuring that each transaction is processed only once. The transaction_index indicates the transaction's position within the block, while from_address and to_address denote the blockchain addresses of the sender and receiver, respectively. The value field specifies the amount of cryptocurrency transferred. The dataset also includes gas, representing the gas limit provided for the transaction, and gas_price, indicating the price per gas unit. The input field contains additional data attached to the transaction. Moreover, receipt_cumulative_gas_used provides the total gas used by all transactions up to and

including the current one within the block, and `receipt_gas_used` specifies the gas consumed by this particular transaction. The `block_timestamp` and `block_number` identify the time and number of the block that includes the transaction, while `block_hash` serves as the block's unique identifier. Lastly, the `from_scam` and `to_scam` fields indicate whether the sender's and receiver's addresses are associated with scams (0 for no, 1 for yes). The dataset also includes categorical data, `from_category` and `to_category`, which classify the nature of the participants, such as Phishing, Scamming, or Fake ICO. In this dataset, Benign had 57,000 transactions (79.28%), Scamming had 11,143 transactions (15.51%), Phishing had 3,106 transactions (4.32%), and Fake ICO had only 1 transaction. Dataset 2 includes 71,250 transactions.

To ensure the reliability of our datasets, we implemented a comprehensive data preprocessing phase. Maintaining the quality and consistency of data was paramount, so we tackled issues like missing values in critical fields by assigning reasonable default values—for instance, substituting missing entries in the `to_address` field with 'Unknown' and filling absent input fields with '0x'. We also addressed outliers by applying normalization techniques, which helped maintain uniformity across the datasets. Additionally, features with high cardinality such as the From and To addresses were encoded to manage their complexity, and temporal features were standardized to Unix epoch time. These preprocessing steps were crucial in preparing the data for the development of robust machine learning models.

After preprocessing, the next step is feature selection, which involves identifying key content-based and network-based features for phishing detection. Content-based features like Value, Gas, and Gas_Price describe transaction details, while network-based features such as From, To, `from_category`, and `to_category` capture transaction behavior. These features helped refine the dataset for analysis. Feature engineering was applied to improve data representation. High-cardinality features were encoded, and `from_category` and `to_category` were merged into `combined_category` to summarize transaction types. The input feature, stored as hexadecimal data, was converted into a numerical format for machine learning. These steps ensured that the data was ready for analysis.

4.2 Experimental Procedures

We utilized the two datasets to examine how simple AEs impact the classification accuracy of RF, DT, and KNN classifiers.

Minimal Manipulations. AEs were crafted by manipulating specific features. For the first dataset, we conducted the following simple manipulations.

- ① *Timestamp Manipulation (TimeStamp):* We apply predefined intervals to simulate future events, evaluating model robustness against temporal shifts. These intervals introduce different time variations, allowing us to assess how well classifiers maintain accuracy across changing time frames.
- ② *Value Manipulation (Value):* Transaction values were adjusted using two methods: adding a fixed percentage uniformly or applying a random percentage change based on the original value. These adjustments tested the models' sensitivity to value changes and their ability to detect phishing activities despite variations.
- ③ *Receiver Address Manipulation (To):* The receiver's address was randomly reassigned to other addresses in the dataset to simulate phishing transactions with different destinations. This approach tested the classifiers' ability to detect phishing attempts despite changes in transaction recipients.

- ④ *Sender Address Manipulation (From)*: The sender's address was replaced with different addresses to simulate phishing transactions from various sources. This adjustment assessed the classifiers' ability to detect phishing attempts despite changes in the transaction origin.

The two datasets were divided into training and testing sets using an 80-20 ratio to ensure a balanced foundation for model training and evaluation. We trained the RF, DT, and KNN classifiers on the training subsets. We assessed their performance using overall accuracy and class-specific accuracy metrics for each label: benign, phishing, and scamming.

We applied both targeted and untargeted adversarial attacks to the second dataset [9]. AEs were created by modifying key transaction features, enabling the simulation of realistic attacks and the identification of potential vulnerabilities.

Untargeted Attacks. The study of untargeted attacks involved generating AEs through broad, random perturbations across the feature space. Initially, all features were manipulated to assess the model's overall resilience. Results were compiled into detailed tables, highlighting performance across different perturbations and offering a comprehensive view of feature influence on robustness. Next, individual features were targeted for a more precise evaluation. Modifying *from_address* and *to_address* introduced unseen addresses, testing the model's ability to adapt to origin and destination changes. Adjusting *value*, *gas*, and *gas_price* simulated fluctuations, assessing sensitivity to transaction cost variations. Manipulating *block_timestamp* and *block_number* examined the model's response to sequence and timing changes. Altering *input*, *receipt_cumulative_gas_used*, and *receipt_gas_used* evaluated the impact of content modifications.

Targeted Attacks. We conducted targeted attacks focusing on benign, phishing, and scamming scenarios. We employed two methods for generating these targeted attacks: rule-based and gradient-based (using the Fast Gradient Sign Method).

Rule-based Modifications This approach adjusts key features, such as transaction value and timestamps, following predefined rules to introduce controlled variations that may influence classification. Additionally, random hexadecimal strings were generated to replace existing *from_address* and *to_address*, creating unseen addresses. This process evaluated the model's ability to recognize phishing attempts when transaction origins and destinations changed and examined its effectiveness in detecting complex fraud attempts.

- **Benign Targeted Attack:** This scenario evaluated the model's ability to maintain a benign classification despite manipulations, mimicking tactics used to disguise malicious activities. Artificial benign transactions were generated by slightly adjusting features such as *transaction value* and *block_timestamp*. For instance, minor modifications to transaction values tested whether small financial changes could impact classification accuracy, potentially leading to misclassifications.
- **Phishing Targeted Attack:** This attack altered phishing-labeled transactions to make them appear legitimate. Adjustments were made to *from_address*, *to_address*, and *value* to mimic attempts to bypass security checks. This method tested whether the model could still identify phishing transactions after changes to key transaction details.
- **Scamming Targeted Attack:** Scam-labeled transactions were manipulated to test whether they could be misclassified as benign or other category. Adjustments to *gas* and *gas_price* simulated efforts to obscure scams by modifying transaction costs. This analysis provided insights into the model's sensitivity to economic factors and its ability to detect scams despite cost-based manipulations.

FGSM. FGSM introduced small, calculated changes to influence the model's predictions while keeping feature values within a realistic range. Unlike the rule-based method, it preserved the

dataset's original distribution and adjusted features based on gradients to increase the likelihood of misclassification. This method tested the model's sensitivity to adversarial modifications. By keeping transactions similar to legitimate ones, FGSM created slight variations that resulted in incorrect classifications. This approach helped identify specific weaknesses in the model's ability to detect phishing transactions and provided insights into areas for improvement.

- **FGSM Details:** FGSM calculates perturbations that align with the gradient direction of the loss function. The perturbations were applied using the formula:

$$x' = x + \epsilon \cdot \text{sign}(\nabla_x J(\theta, x, y)),$$

where x is the original feature, x' is the perturbed feature (AE), ϵ is a small scalar controlling the perturbation's magnitude, ∇_x represents the gradient of the loss function J for x , and $J(\theta, x, y)$ is the loss function dependent on the model parameters θ , input x , and true label y . The term $\text{sign}(\nabla_x J(\theta, x, y))$ provides the direction in which to perturb the feature vector to maximize the loss function.

- **Features:** FGSM modifications were applied to transaction value, gas, gas_price, and block_timestamp. These features play a key role in determining transaction type, and even small changes can influence classification. The method introduced minimal but targeted adjustments, following the gradient direction to increase the likelihood of prediction errors.

Evaluation Metrics. The performance of the model was assessed using precision, recall, accuracy, and F1-score. The precision is defined as $\text{Precision} = \frac{TP}{TP+FP}$, $\text{Recall} = \frac{TP}{TP+FN}$, $\text{Accuracy} = \frac{TP+TN}{TP+TN+FP+FN}$, and $\text{F1-score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$ where TP (True Positives) refers to correctly predicted positive cases, TN (True Negatives) refers to correctly predicted negative cases, FP (False Positives) refers to incorrectly predicted positive cases and FN (False Negatives) refers to incorrectly predicted negative cases.

5 RESULTS AND ANALYSIS

5.1 Preliminary Results

We evaluate how different perturbations on the first dataset affect the accuracy and resilience of RF, DT, and KNN models. We examine how uniform and proportional value manipulations and address change impact model performance.

Timestamp Manipulations. We evaluated how timestamp changes affected classifier accuracy by applying shifts of +24 hours, +1 hour, +30 minutes, +15 minutes, and +5 minutes. These shifts represented transaction delays or intentional modifications to evade detection. The RF classifier showed the least impact, with accuracy dropping slightly from 0.98 to 0.95 for a one-day shift and to 0.97 for a one-hour shift. DT experienced a larger decline, with accuracy decreasing from 0.98 to 0.94 under a one-day shift. KNN was the most affected, with accuracy falling from 0.94 to 0.83 for the same change. These results indicate that RF is more resilient to timestamp variations, making it a strong option for phishing detection when transaction timings are altered (Table 10).

Value Manipulations. When transaction values were changed uniformly, RF and DT accuracy dropped to 0.69, making phishing detection less effective. Attackers who consistently modify transaction values could exploit this weakness. DT's recall for phishing was especially affected, falling to just 0.01%. In contrast, proportional changes had little effect, with accuracy and other measures staying nearly the same. This suggests that models adapt better to these variations, which reflect normal transaction patterns. KNN performed consistently under both conditions, showing it was less impacted by value changes (Table 3).

Table 2. Performance of Random Forest (RF), Decision Tree (DT), and K-Nearest Neighbors (KNN) models under timestamp manipulations across different increments. Metrics include accuracy, precision, recall, F1-score, and count for benign and phishing labels. Increments: original (O), +24 hours (+24), +1 hour (+1), +30 minutes (+30), +15 minutes (+15), and +5 minutes (+5). The baseline dataset represents normal conditions, while the adversarial dataset contains manipulated timestamps.

Increment	Model	Accuracy	Precision		Recall		F1-score		Count	
			Benign	Phishing	Benign	Phishing	Benign	Phishing	Benign	Phishing
Baseline Dataset										
O	RF	0.98	1.00	0.96	0.98	1.00	0.99	0.98	15989	7483
	DT	0.98	1.00	0.95	0.98	1.00	0.99	0.97	15989	7483
	KNN	0.94	1.00	0.85	0.92	1.00	0.96	0.92	15989	7483
Adversarial Dataset										
+24	RF	0.95	0.94	0.98	0.99	0.86	0.96	0.92	15498	7974
	DT	0.94	0.95	0.94	0.97	0.88	0.96	0.91	15601	7871
	KNN	0.83	0.84	0.82	0.94	0.62	0.88	0.70	14376	9096
+1	RF	0.97	0.97	0.97	0.99	0.94	0.98	0.96	15498	7974
	DT	0.96	0.97	0.95	0.98	0.94	0.97	0.94	15601	7871
	KNN	0.89	0.93	0.84	0.93	0.84	0.93	0.84	14376	9096
+30	RF	0.97	0.98	0.97	0.99	0.96	0.98	0.96	15498	7974
	DT	0.96	0.97	0.95	0.98	0.94	0.97	0.94	15601	7871
	KNN	0.91	0.95	0.85	0.92	0.90	0.94	0.87	14376	9096
+15	RF	0.98	0.99	0.97	0.99	0.97	0.99	0.97	15498	7974
	DT	0.96	0.98	0.95	0.98	0.95	0.98	0.95	15601	7871
	KNN	0.91	0.95	0.85	0.93	0.9	0.94	0.87	14376	9096
+5	RF	0.97	0.99	0.97	0.98	0.97	0.98	0.97	15498	7974
	DT	0.97	0.98	0.95	0.98	0.97	0.98	0.96	15601	7871
	KNN	0.93	0.98	0.85	0.92	0.96	0.95	0.9	14376	9096

Table 3. Performance evaluation of Random Forest (RF), Decision Tree (DT), and K-Nearest Neighbors (KNN) models subjected to 1% uniform (U) and proportional (P) value manipulation strategies compared to the original (O). Metrics include accuracy, precision, recall, F1-score, and count for benign and phishing labels.

Model	Strategy	Accuracy	Precision		Recall		F1-score		Count		
			Benign	Phishing	Benign	Phishing	Benign	Phishing	Benign	Phishing	
Random Forest (RF)											
O	Original	0.99	0.98	1.00	1.00	0.99	0.99	0.99	0.99	15803	7669
U	Uniform	0.69	0.96	0.68	0.02	1.00	0.03	0.81	0.03	23353	119
P	Proportional	0.99	0.98	1.00	1.00	0.99	0.99	0.99	0.99	15813	7659
Decision Tree (DT)											
O	Original	0.98	0.95	1.00	1.00	0.98	0.97	0.99	0.97	15601	7871
U	Uniform	0.69	0.75	0.68	0.02	1.00	0.03	0.81	0.03	23294	178
P	Proportional	0.98	0.95	1.00	1.00	0.98	0.97	0.99	0.97	15619	7853
K-Nearest Neighbors (KNN)											
O	Original	0.96	0.89	0.99	0.98	0.94	0.93	0.97	0.93	15192	8280
U	Uniform	0.96	0.89	0.99	0.98	0.94	0.93	0.97	0.93	15193	8279
P	Proportional	0.96	0.89	0.99	0.98	0.94	0.93	0.97	0.93	15192	8280

Address Manipulations. The models were tested against AEs by modifying the From and To address features in 5,000, 10,000, and 23,472 transactions. Address changes are common in phishing attacks, where attackers alter source or destination addresses to hide fraudulent transactions. For the RF model, accuracy dropped to 0.87 when the From address was changed (Table 4) and to 0.84 for the To address (Table 5). Precision and recall for phishing transactions also declined, leading to a lower F1 score. While RF remained stable, these results show that address changes can still affect its performance, which could be a concern in real world. The DT model experienced a smaller accuracy reduction, reaching 0.92 for From address changes (Table 4) and 0.93 for To (Table 5).

Table 4. Performance evaluation of Random Forest (RF), Decision Tree (DT), and K-Nearest Neighbors (KNN) models under manipulations of the From feature in Ethereum transaction datasets. Metrics include accuracy, precision, recall, F1-score, and count for benign and phishing labels. The baseline dataset (Ba) represents normal conditions, while other values correspond to adversarial manipulations.

Model	Manipulation	Accuracy	Precision		Recall		F1-score		Count	
			Benign	Phishing	Benign	Phishing	Benign	Phishing	Benign	Phishing
Random Forest (RF)										
—	Baseline	0.99	1.00	0.96	0.98	1.00	0.99	0.98	15708	7764
5000	Adversarial	0.96	0.96	0.97	0.98	0.91	0.97	0.94	16386	7086
10000	Adversarial	0.94	0.93	0.97	0.99	0.83	0.96	0.89	17027	6445
23472	Adversarial	0.87	0.84	0.97	0.99	0.60	0.91	0.74	18877	4595
Decision Tree (DT)										
—	Baseline	0.98	1.00	0.95	0.98	1.00	0.99	0.97	15601	7871
5000	Adversarial	0.97	0.98	0.95	0.98	0.96	0.98	0.96	15935	7537
10000	Adversarial	0.96	0.96	0.95	0.98	0.91	0.97	0.93	16272	7200
23472	Adversarial	0.92	0.91	0.94	0.98	0.79	0.94	0.86	17207	6265
K-Nearest Neighbors (KNN)										
—	Baseline	0.94	1.00	0.85	0.92	1.00	0.96	0.92	14706	8766
5000	Adversarial	0.93	0.97	0.85	0.92	0.93	0.94	0.89	15209	8263
10000	Adversarial	0.91	0.94	0.84	0.92	0.87	0.93	0.85	15767	7705
23472	Adversarial	0.85	0.86	0.81	0.93	0.67	0.89	0.74	17288	6184

Table 5. Performance evaluation of Random Forest (RF), Decision Tree (DT), and K-Nearest Neighbors (KNN) models under manipulations of the To feature in Ethereum transaction datasets. Metrics include accuracy, precision, recall, F1-score, and count for benign and phishing labels. The baseline dataset (Ba) represents normal conditions, while other values correspond to adversarial manipulations.

Model	Manipulation	Accuracy	Precision		Recall		F1-score		Count	
			Benign	Phishing	Benign	Phishing	Benign	Phishing	Benign	Phishing
Random Forest (RF)										
—	Baseline	0.99	1.00	0.96	0.98	1.00	0.99	0.98	15708	7764
5000	Adversarial	0.96	0.95	0.97	0.98	0.89	0.97	0.93	16558	6914
10000	Adversarial	0.92	0.91	0.97	0.99	0.79	0.95	0.87	17377	6095
23472	Adversarial	0.84	0.81	0.97	0.99	0.51	0.89	0.67	19537	3935
Decision Tree (DT)										
—	Baseline	0.98	1.00	0.95	0.98	1.00	0.99	0.97	15601	7871
5000	Adversarial	0.97	0.98	0.95	0.98	0.96	0.98	0.96	15884	7588
10000	Adversarial	0.96	0.97	0.95	0.98	0.93	0.97	0.94	16142	7330
23472	Adversarial	0.93	0.92	0.94	0.98	0.83	0.95	0.88	16872	6600
K-Nearest Neighbors (KNN)										
—	Baseline	0.94	1.00	0.85	0.92	1.00	0.96	0.92	14706	8766
5000	Adversarial	0.94	0.99	0.85	0.92	0.98	0.95	0.91	14849	8623
10000	Adversarial	0.94	0.98	0.85	0.92	0.97	0.95	0.91	14988	8484
23472	Adversarial	0.93	0.96	0.85	0.93	0.93	0.94	0.89	15351	8121

However, DT's ability to handle address modifications better may be due to its decision-making process, which can manage categorical data more effectively. KNN was the most affected, with accuracy falling to 0.85 for From and 0.93 for To changes. This led to sharp declines in precision and recall for phishing transactions (Tables 4 and 5). KNN's high sensitivity to address modifications suggests it may struggle in environments where attackers frequently change transaction addresses, making it less reliable in detecting phishing.

RF and DT handled timestamp and value changes better than KNN. To make ML models more secure and reliable, it is important to refine their design and include adversarial training. The differences in model performance show why thorough testing is necessary. Next, the analysis will examine the second dataset from Al-Emari *et al.* [9], which is used for multi-class classification.

Table 6. Impact of adversarial attacks on benign, phishing, and scamming detection for Random Forest (RF), Decision Tree (DT), and K-Nearest Neighbors (KNN). The table presents accuracy changes before and after attacks, instance count shifts, and label changes for phishing (P), benign (B), scamming (S), and fake ICO (F).

Metric	Random Forest (RF)		Decision Tree (DT)		K-Nearest Neighbors (KNN)	
	Pre	Post	Pre	Post	Pre	Post
Accuracy Changes Due to Attacks						
Benign Detection	1.00	0.84	0.99	0.84	0.97	0.90
Phishing Detection	0.96	0.01	0.96	0.01	0.41	0.02
Scamming Detection	0.99	0.14	0.98	0.14	0.67	0.07
Instance Counts Before and After Attacks						
Benign	11,431	12,026	11,431	12,063	11,431	12,861
Phishing	629	187	629	168	629	306
Scamming	2,189	2,033	2,189	2,018	2,189	1,083
Label Changes Due to Attack						
Benign → [P, S, F]	[0, 5, 0]		[0, 5, 0]		[60, 97, 0]	
Phishing → [B, S, F]	[325, 117, 0]		[325, 146, 0]		[372, 27, 0]	
Scamming → [B, P, F]	[274, 0, 0]		[311, 10, 0]		[1,214, 16, 0]	

5.2 Results of Targeted Attacks

Rule-based Modifications. This study first examines rule-based adversarial modifications that target specific transaction features to assess how simple changes can affect model performance. These modifications reflect real-world cases where attackers alter some features to avoid detection. Table 6 shows how these changes impact model performance under adversarial attacks.

① *Benign Class.* RF and DT originally classified benign transactions with high accuracy on the test set. However, when exposed to adversarial changes, their accuracy dropped to 0.84, as shown in Table 6, a decline of over 0.15. This suggests that even small changes can make it harder for the models to correctly classify transactions, creating a risk where attackers can make fraudulent activities appear legitimate. In contrast, KNN performed better under these conditions, maintaining an accuracy of 0.90. Its method, which groups transactions based on similarity, may help it resist some of these changes. The number of benign transactions affected by adversarial modifications increased: from 11,431 to 12,026 for RF, 12,063 for DT, and 12,861 for KNN, as shown in Table 6.

② *Phishing Class.* RF and DT initially had a phishing detection accuracy of 0.96, as shown in Table 6. However, after adversarial changes, their accuracy dropped to 0.01, a decrease of more than 0.95. This shows that both models struggle to detect phishing when key transaction details are modified, making them vulnerable to targeted attacks. KNN, which started with a lower accuracy of 0.41, dropped to 0.02. Even though it had a weaker starting point, its near-total failure under attack highlights how difficult it is to maintain phishing detection when attackers exploit specific weaknesses (see Table 6). The number of correctly identified phishing transactions dropped sharply: from 629 to 187 for RF, 168 for DT, and 306 for KNN. Meanwhile, more phishing transactions were misclassified as benign, increasing to 325 for RF and DT and 372 for KNN, as shown in Table 6.

③ *Scamming Class* RF and DT originally had high accuracy in detecting scamming transactions, at 0.99 and 0.98, as shown in Table 6. After adversarial changes, their accuracy dropped to 0.14, a decrease of more than 0.85. This shows that both models are highly affected by small modifications, which could weaken their ability to detect fraud. KNN started with an accuracy of 0.67 but fell to 0.07. While it performed well in general classification, it struggled with scam detection, making it less reliable in such cases. The number of correctly classified scamming transactions dropped from 2,189 to 2,033 for RF, 2,018 for DT, and 1,083 for KNN. Misclassifications as benign were highest for KNN, with 1,214 cases, compared to 274 for RF and 311 for DT, as shown in Table 6.

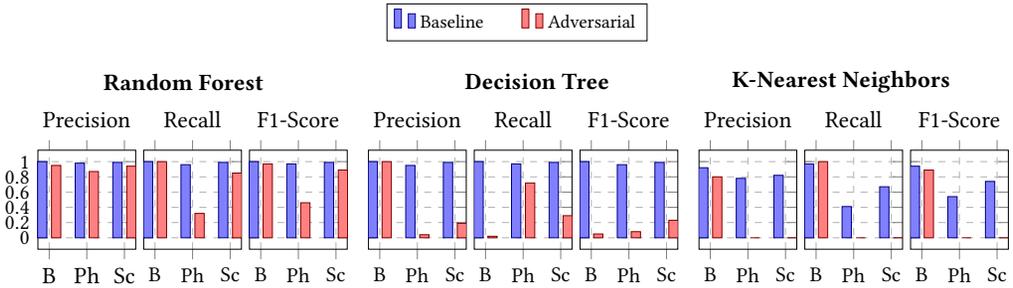


Fig. 2. Performance comparison of RF, DT, and KNN under baseline and adversarial conditions using the FGSM on Benign class. The figure illustrates the Precision, Recall, and F1 Score metrics across three classes: Benign, Phishing, and Scamming.

Takeaway

Targeted adversarial attacks undermine the effectiveness of machine learning models in detecting phishing and scamming activities, revealing a critical weakness in their robustness across different algorithms.

5.3 Gradient-based Approach Using FGSM

① *Benign Class* The overall accuracy of the RF model dropped from 0.99 to 0.94, with a notable decline in phishing detection. However, it maintained a high accuracy of 0.99 for benign transactions. This suggests that while RF can still recognize benign transactions under adversarial conditions, its ability to detect phishing attempts is significantly weakened, creating a trade-off between precision and security. In contrast, DT's accuracy fell sharply from 0.99 to 0.09, with benign recall dropping to 0.02, indicating high vulnerability. This sharp decline suggests that DT is particularly sensitive to gradient-based attacks, which can exploit its decision boundaries more easily than those of other models. KNN retained a perfect benign accuracy of 1.00 but saw its overall accuracy decrease from 0.90 to 0.80, showing its difficulty in identifying phishing and scamming transactions. The model behaved differently under adversarial conditions, as shown in Figures 2 and 7a. This decline highlights KNN's vulnerability to adversarial attacks, especially when compared to RF and DT. The results remained consistent across multiple metrics, including precision, recall, and F1 scores, with clear performance drops in phishing and scamming detection when exposed to adversarial changes such as FGSM.

② *Phishing Class* The RF model's phishing detection accuracy dropped from 0.96 to 0.47, with a significant decline in the F1 score. While RF still identified some phishing attempts under attack, its performance was far from reliable, highlighting the need for stronger defense mechanisms. DT's overall accuracy fell to 0.10, with phishing recall decreasing to 0.75, showing high sensitivity to adversarial attacks. This suggests a structural weakness in DT, making it easier for attacks to manipulate their predictions, particularly for phishing detection. KNN failed entirely to detect phishing transactions under adversarial conditions, with all related metrics dropping to zero. This result shows that KNN becomes ineffective against phishing threats when adversarial modifications are introduced. Figures 3 and 7b provide a detailed view of how FGSM impacts phishing detection across different models.

③ *Scamming Class* The overall accuracy of the RF model dropped from 0.99 to 0.81, with scamming accuracy decreasing from 0.99 to 0.76. While RF still outperforms other models in this scenario, the significant drop indicates a need for improved resistance to adversarial attacks in scam detection.

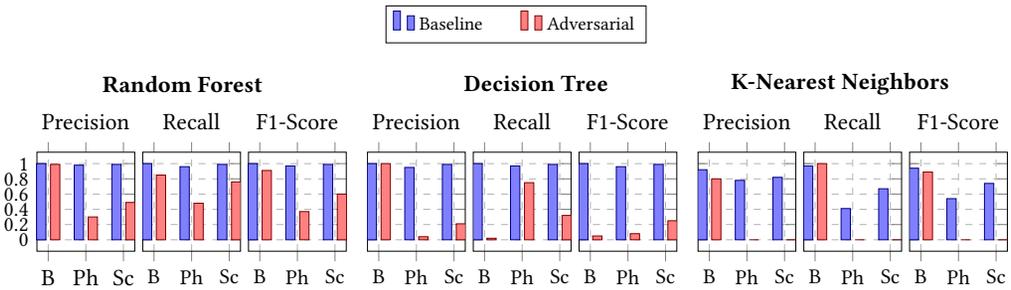


Fig. 3. Evaluation of RF, DT, and KNN performance with FGSM on the Phishing class. Metrics cover precision, recall, and F1 scores for the Benign, Phishing, and Scamming classes.

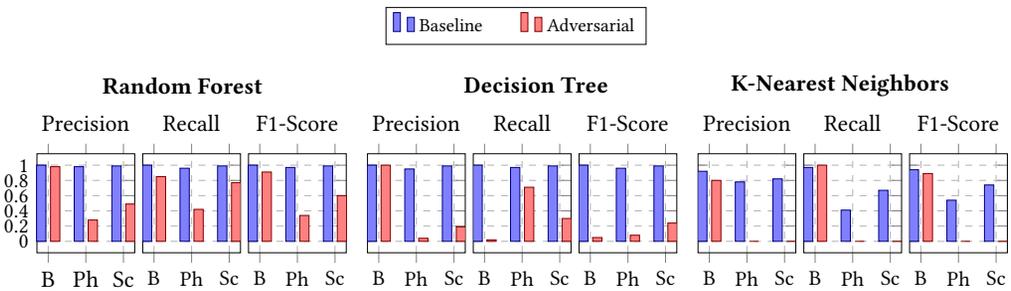


Fig. 4. Assessment of RF, DT, and KNN performance with FGSM on the Scamming class. Metrics include overall and scamming class accuracy, as well as precision, recall, and F1 scores for Benign, Phishing, and Scamming across baseline and adversarial conditions.

The DT’s overall accuracy fell to 0.09, with scamming recall drastically reducing to 0.30. The DT model’s performance suggests it is highly vulnerable to adversarial manipulations in scam-related features, making it less reliable for detecting such fraudulent activities. The KNN model’s scamming detection metrics also dropped to zero. Figures 4 and 7c depicts a comparison between the baseline and adversarial performance using FGSM on the scamming class across different models.

Takeaway

Gradient-based adversarial attacks, such as FGSM, severely compromise the effectiveness of machine learning models in detecting phishing and scamming activities, particularly for Decision Tree and KNN models, which become nearly ineffective under such conditions.

5.4 Results of Untargeted Attacks

We employed untargeted adversarial attacks to evaluate the models’ performance in adversarial scenarios. This method involved modifying different features to examine which models could sustain their detection capabilities when exposed to manipulated data. Rather than concentrating on specific targeted attacks, the primary objective was to assess the overall robustness of each model against a diverse range of perturbations.

① *All Features* The RF model’s accuracy decreased to 0.95, DT’s to 0.91, with phishing detection severely impaired, and KNN maintained its baseline accuracy. These findings indicate that while RF and DT models exhibit robustness, they remain vulnerable to broad, untargeted adversarial

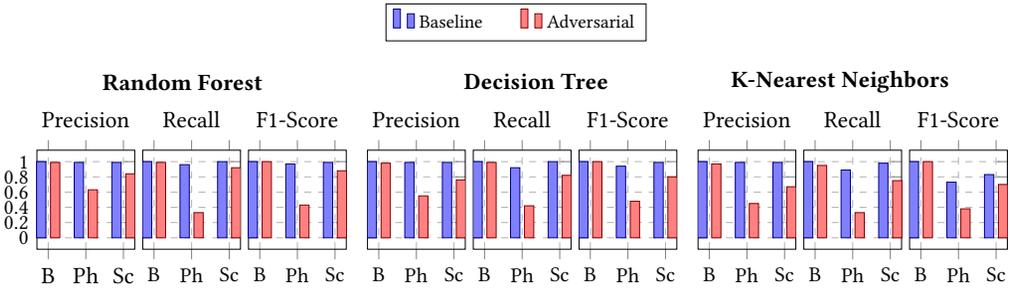


Fig. 5. Performance comparison of RF, DT, and KNN models under baseline and untargeted adversarial attacks. The plots illustrate Precision, Recall, and F1-Score across Benign, Phishing, and Scamming.

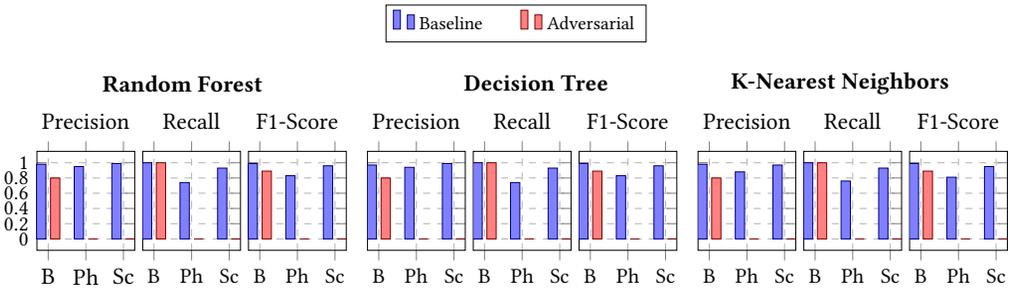


Fig. 6. RF, DT, and KNN performance with address feature manipulation. Metrics: accuracy, precision, recall, and F1 for Benign, Phishing, and Scamming classes.

modifications. Notably, the phishing precision for RF drops from 0.99 to 0.63, and for DT, from 0.96 to 0.15. Similarly, phishing recall decreases from 0.96 to 0.33 for RF and from 0.97 to 0.15 for DT. These reductions underscore the models' reduced ability to identify phishing attempts when under attack. The details are illustrated in Figures 5 and 7d, where performance metrics under these conditions are compared across all models. This significant drop in accuracy highlights the susceptibility of these models to broad adversarial manipulations.

② *Address Features* AEs focusing on the *from_address* and *to_address* features resulted in a decline in overall accuracy to 0.80 for all models. None of the models detected phishing or scamming, indicating a high sensitivity to address manipulations. This result highlights the critical importance of accurately processing and interpreting address features, as adversaries can easily exploit these weaknesses to evade detection. The complete drop in phishing and scamming detection metrics to zero across all models is alarming, considering the baseline phishing precision and recall were 0.95 and 0.74 for RF, 0.94 and 0.74 for DT, and 0.88 and 0.76 for KNN, respectively. These declines underscore the vulnerability of these models when address features are manipulated. This vulnerability is further illustrated in Figures 6 and 7e, demonstrating the impact of unseen address manipulations on model performance.

These findings highlight the importance of accurate processing, as attackers can easily exploit any manipulation to bypass detection. The fact that phishing and scamming detection metrics dropped to zero across all models is alarming. Before manipulation, the baseline phishing precision and recall were 0.95 and 0.74 for RF, 0.94 and 0.74 for DT, and 0.88 and 0.76 for KNN. The drop in performance underscores how vulnerable these models become when key features are altered. The results in Table 7, show the impact of unseen address manipulations on model performance.

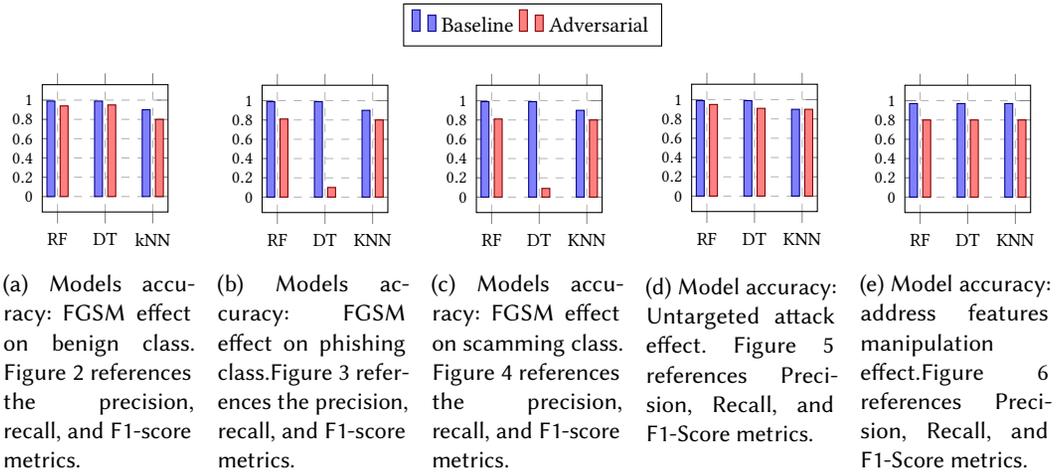


Fig. 7. Combined model accuracy across 2-6 figures.

Table 7. Impact of unseen address features on model performance for different models (RF, DT, KNN) using original and modified data. Metrics include accuracy, precision, recall, and F1-score for benign, phishing, and scamming classes.

Model	Data	Accuracy	Precision			Recall			F1-score		
			Benign	Phishing	Scam	Benign	Phishing	Scam	Benign	Phishing	Scam
RF	Original	0.97	0.98	0.95	0.99	1.00	0.74	0.93	0.99	0.83	0.96
	Modified	0.80	0.80	0.00	0.00	1.00	0.00	0.00	0.89	0.00	0.00
DT	Original	0.97	0.97	0.94	0.99	1.00	0.74	0.93	0.99	0.83	0.96
	Modified	0.80	0.80	0.00	0.00	1.00	0.00	0.00	0.89	0.00	0.00
KNN	Original	0.97	0.98	0.88	0.97	1.00	0.76	0.93	0.99	0.81	0.95
	Modified	0.80	0.80	0.00	0.00	1.00	0.00	0.00	0.89	0.00	0.00

③ *Financial Features AEs* targeting financial features (*value, gas, gas_price*) led to a reduced RF’s accuracy of 0.79. Similarly, the DT’s accuracy dropped to 0.79. The KNN model’s accuracy slightly decreased to 0.90. The impact of these adversarial manipulations is particularly evident in the phishing precision and recall metrics. For example, under value manipulation, RF’s phishing precision dropped from 0.69 to 0.24 and recall from 0.56 to 0.27. Similarly, DT’s phishing precision decreased from 0.59 to 0.20, with recall dropping from 0.59 to 0.31. KNN also showed a reduction in phishing precision from 0.73 to 0.58 and recall from 0.51 to 0.39 under the same conditions. These results are detailed in Table 8.

④ *Using Temporal Features* Adversarial manipulations of temporal features (*block_timestamp, block_number*) showed the RF model’s accuracy fell from 0.99 to 0.80, with phishing detection metrics nearly nullified. This significant drop indicates that temporal features are a major vulnerability for RF and likely for other models, which adversaries can exploit to severely disrupt the model’s ability to detect malicious activities. For instance, under timestamp manipulation, RF’s phishing recall dropped drastically from 0.91 to zero. The DT’s accuracy similarly declined to 0.80. The KNN model’s accuracy also dropped to 80.26%. These results are further detailed in Table 9, which compares model performance across different strategies, including original, timestamp-manipulated, block number-manipulated, and combined approaches, highlighting the models’ vulnerabilities to temporal feature perturbations.

Table 8. Impact of financial feature manipulation on model performance for different models (RF, DT, KNN) using original and manipulated data. Metrics include accuracy, precision, recall, and F1-score for benign, phishing, and scamming classes.

Model	Strategy	Accuracy	Precision			Recall			F1-score		
			Benign	Phishing	Scam	Benign	Phishing	Scam	Benign	Phishing	Scam
RF	Original	0.92	0.95	0.69	0.85	0.97	0.56	0.79	0.96	0.62	0.82
	Value	0.79	0.90	0.24	0.53	0.87	0.27	0.60	0.88	0.26	0.56
	Gas	0.90	0.93	0.69	0.83	0.97	0.45	0.73	0.95	0.54	0.78
	Gas Price	0.91	0.92	0.74	0.87	0.98	0.48	0.70	0.95	0.58	0.78
	Combined	0.78	0.87	0.25	0.48	0.88	0.19	0.48	0.87	0.21	0.48
DT	Original	0.91	0.95	0.59	0.81	0.95	0.59	0.80	0.95	0.59	0.80
	Value	0.79	0.92	0.20	0.53	0.85	0.31	0.66	0.88	0.24	0.59
	Gas	0.84	0.94	0.42	0.60	0.88	0.54	0.76	0.91	0.47	0.67
	Gas Price	0.88	0.93	0.52	0.75	0.94	0.51	0.70	0.93	0.51	0.72
	Combined	0.71	0.89	0.15	0.39	0.77	0.29	0.54	0.83	0.20	0.45
KNN	Original	0.92	0.94	0.73	0.85	0.97	0.51	0.76	0.96	0.60	0.80
	Value	0.90	0.93	0.58	0.80	0.96	0.39	0.74	0.94	0.46	0.77
	Gas	0.89	0.93	0.68	0.74	0.95	0.47	0.72	0.94	0.56	0.73
	Gas Price	0.89	0.92	0.67	0.80	0.96	0.43	0.68	0.94	0.52	0.73
	Combined	0.86	0.91	0.53	0.67	0.94	0.31	0.64	0.92	0.39	0.65

Table 9. Impact of temporal feature manipulation on model performance for different models (RF, DT, KNN) using original and manipulated data. Metrics include accuracy, precision, recall, and F1-score for benign, phishing, and scamming classes.

Model	Strategy	Accuracy	Precision			Recall			F1-score		
			Benign	Phishing	Scam	Benign	Phishing	Scam	Benign	Phishing	Scam
RF	Original	0.99	1.00	0.88	0.97	1.00	0.91	0.96	1.00	0.89	0.97
	Timestamp	0.80	0.80	1.00	1.00	1.00	0.00	0.00	0.89	0.01	0.00
	Block Number	0.95	1.00	0.46	0.80	1.00	0.19	0.94	1.00	0.27	0.86
	Combined	0.80	0.80	0.00	0.89	1.00	0.00	0.01	0.89	0.00	0.01
DT	Original	0.99	1.00	0.87	0.97	1.00	0.91	0.96	1.00	0.89	0.97
	Timestamp	0.80	0.80	1.00	1.00	1.00	0.00	0.00	0.89	0.01	0.00
	Block Number	0.95	1.00	0.46	0.80	1.00	0.19	0.94	1.00	0.27	0.86
	Combined	0.80	0.80	0.00	0.89	1.00	0.00	0.01	0.89	0.00	0.01
KNN	Original	0.99	1.00	0.92	0.96	1.00	0.87	0.98	1.00	0.89	0.97
	Timestamp	0.80	0.80	1.00	1.00	1.00	0.00	0.00	0.89	0.01	0.00
	Block Number	0.94	1.00	0.52	0.77	0.99	0.20	0.95	1.00	0.29	0.85
	Combined	0.80	0.80	0.00	0.85	1.00	0.00	0.01	0.89	0.00	0.02

Takeaway

These results underscore the need for robust defensive mechanisms against AEs. The significant decline in performance metrics under simple conditions highlight the need for a more reliable classification of transactions.

6 DEFENSE MECHANISM

In this section, we propose a defense mechanism to mitigate the impact of adversarial attacks on Ethereum phishing transaction detection models. Given the vulnerability of machine learning models to adversarial perturbations, our approach centers around adversarial training [28, 55], which enhances model robustness by incorporating AEs into the training process and allowing the model to learn to resist manipulations. We note that it is important to distinguish adversarial training from traditional data poisoning. While data poisoning aims to compromise model performance by injecting misleading data, adversarial training is a defensive measure designed to improve model resilience to real-world adversarial scenarios.

6.1 Adversarial Training Approach

We applied adversarial training to defend against targeted adversarial attacks across all classes—benign, phishing, and scamming—which were individually targeted. We also focused on specific feature manipulations, particularly the *TimeStamp* and *Value* features. To implement the adversarial training approach, we followed these steps.

① **Generating Adversarial Examples.** We employed existing adversarial example generation functions to create targeted samples by manipulating key features. This involved generating AEs with minimal perturbations to the *TimeStamp* and *Value* features. Moreover, we used FGSM to craft targeted AEs by perturbing the input features associated with benign, phishing, and scamming labels, aiming to challenge the models with realistic adversarial scenarios [16]. These perturbations were carefully controlled to remain within realistic bounds, ensuring that the examples would simulate genuine attack scenarios while minimizing potential performance compromises.

② **Models Retraining.** The RF, DT, and KNN models were retrained using an augmented dataset that included both original and adversarially modified data. During this process, model parameters were adjusted to improve resilience against adversarial attacks. Unlike poisoning, which reduces model performance, the AEs used in this study were intended to identify weaknesses in key feature spaces that influence classification accuracy. This method aimed to help the models generalize to new attacks without becoming overly adapted to AEs. By introducing potentially misleading patterns, the goal was to enhance the models' ability to classify transactions, particularly those in categories such as benign and phishing.

③ **Models Evaluation.** After training with AEs, the models were tested on both the original and modified sets to evaluate their ability to recognize manipulated. The assessment focused on key performance metrics, including accuracy, precision, recall, and F1 score, while also examining misclassification patterns to understand how the models responded to adversarial attacks.

7 POST RETRAINING RESULTS

7.1 Preliminary Results

After applying timestamp and value manipulations, we assess the impact of adversarial retraining on the RF, DT, and KNN models. We analyze how retraining enhances model accuracy and resilience, focusing on restoring precision, recall, and F1 scores across all models.

Timestamp. Before retraining, timestamp manipulations caused significant accuracy reductions, with RF at 0.95, DT at 0.94, and KNN at 0.83 (Table 10). After retraining, these models showed enhanced resilience, maintaining accuracies of 0.98 for RF and DT, and improving to 0.94 for KNN (Table 10). Stability in precision, recall, and F1 scores across all models indicates that adversarial training effectively mitigated the adverse effects of timestamp manipulations. The differential impact on KNN can be attributed to its inherent sensitivity to the input space distribution, which makes timestamp manipulations disruptive. We note that while we observed slight variances in the non-adversarial test sets after retraining, these were negligible compared to the significant gains in adversarial robustness.

Value. Value manipulations significantly affected the RF and DT models, reducing their accuracy from 0.99 and 0.98 to 0.69, respectively (Table 3). After adversarial retraining (Table 11), both models regained their original accuracy, while KNN remained stable. This suggests that adversarial training helped restore classification performance under value manipulations. The differences in how RF and DT responded to these changes reflect variations in their handling of feature splits and decision thresholds, with DT being more sensitive to abrupt shifts in feature values than RF.

Table 10. Performance of RF, DT, and KNN after adversarial training on timestamp manipulations for different models, increments, and datasets. Metrics include accuracy, precision, recall, F1-score, and instance count for benign and phishing classes.

Increment	Model	Dataset	Accuracy	Precision		Recall		F1-score		Count	
				Benign	Phish	Benign	Phish	Benign	Phish	Benign	Phish
Original	RF	Baseline	0.98	1.00	0.97	0.98	1.00	0.99	0.98	15,989	7,483
	DT	Baseline	0.98	1.00	0.95	0.98	1.00	0.99	0.97	15,989	7,483
	KNN	Baseline	0.95	1.00	0.88	0.94	1.00	0.97	0.94	15,989	7,483
+24 Hours	RF	Adversarial	0.98	1.00	0.97	0.98	1.00	0.99	0.98	15,989	7,483
	DT	Adversarial	0.98	1.00	0.95	0.98	1.00	0.99	0.97	15,989	7,483
	KNN	Adversarial	0.94	1.00	0.86	0.92	1.00	0.96	0.92	15,989	7,483
+1 Hour	RF	Adversarial	0.98	1.00	0.97	0.98	1.00	0.99	0.98	15,989	7,483
	DT	Adversarial	0.98	1.00	0.95	0.98	1.00	0.99	0.97	15,989	7,483
	KNN	Adversarial	0.95	1.00	0.88	0.94	1.00	0.97	0.94	15,989	7,483
+30 Minutes	RF	Adversarial	0.98	1.00	0.97	0.98	1.00	0.99	0.98	15,989	7,483
	DT	Adversarial	0.98	1.00	0.95	0.98	1.00	0.99	0.97	15,989	7,483
	KNN	Adversarial	0.95	1.00	0.88	0.94	1.00	0.97	0.94	15,989	7,483
+15 Minutes	RF	Adversarial	0.98	1.00	0.97	0.98	1.00	0.99	0.98	15,989	7,483
	DT	Adversarial	0.98	1.00	0.95	0.98	1.00	0.99	0.97	15,989	7,483
	KNN	Adversarial	0.95	1.00	0.89	0.94	1.00	0.97	0.94	15,989	7,483
+5 Minutes	RF	Adversarial	0.98	1.00	0.97	0.98	1.00	0.99	0.98	15,989	7,483
	DT	Adversarial	0.98	1.00	0.95	0.98	1.00	0.99	0.97	15,989	7,483
	KNN	Adversarial	0.95	1.00	0.88	0.94	1.00	0.97	0.94	15,989	7,483

Table 11. Performance evaluation of RF, DT, and KNN models after adversarial training on value manipulation, subjected to 1% uniform (U) and proportional (P) value manipulation strategies compared to the original (O). Metrics include accuracy, precision, recall, F1-score, and instance count for benign and phishing classes.

Model	Strategy	Accuracy	Precision		Recall		F1-score		Count	
			Benign	Phish	Benign	Phish	Benign	Phish	Benign	Phish
RF	Original	0.99	0.98	1.00	1.00	0.99	0.99	0.99	7,483	15,989
	Uniform	0.99	0.98	1.00	1.00	0.99	0.99	0.99	7,661	15,811
	Proportional	0.99	0.98	1.00	1.00	0.99	0.99	0.99	7,652	15,820
DT	Original	0.98	0.95	1.00	1.00	0.97	0.97	0.99	7,483	15,989
	Uniform	0.98	0.95	1.00	1.00	0.97	0.97	0.99	7,888	15,584
	Proportional	0.98	0.95	1.00	1.00	0.97	0.97	0.99	7,884	15,588
KNN	Original	0.98	0.93	1.00	1.00	0.96	0.96	0.98	8,064	15,408
	Uniform	0.98	0.93	1.00	1.00	0.96	0.96	0.98	8,065	15,407
	Proportional	0.98	0.93	1.00	1.00	0.96	0.96	0.98	8,064	15,408

7.2 Post Retraining Results of Targeted Attacks

After retraining, the RF, DT, and KNN models were tested to see how they responded to targeted attacks on different classes. The results showed improvements in their ability to correctly classify benign, phishing, and scamming cases. The KNN model had a slightly higher number of misclassifications, likely because it relies on local neighborhood data, which can be less reliable when class boundaries are altered. Despite this, the overall findings show that adversarial training helped improve model performance.

Benign Class. Benign class accuracy was severely impacted by targeted adversarial attacks, with RF and DT accuracies dropping from 1.00 and 0.99 to 0.84, and KNN from 0.97 to 0.90—indicating significant performance declines of over 15% and 7%, respectively (Table 6). These attacks caused notable misclassifications, especially in the KNN model, where 60 benign instances were incorrectly classified as phishing and 97 as scamming. However, after adversarial training, the models showed substantial recovery (Table 12). Both RF and DT regained their high accuracy levels of 0.99, and KNN improved to 0.98, reflecting a strong restoration of performance. Misclassifications were drastically reduced, with RF and DT eliminating nearly all errors, while KNN minimized its misclassification rates to just 0.22% for phishing and 0.03% for scamming. Overall, adversarial training proved highly

Table 12. Impact of adversarial training using FGSM on RF, DT, and KNN models for benign, phishing, and scamming classes. The table presents accuracy before (Pre) and after (Post) retraining, as well as misclassification rates for different instance types. Cells of repeated values are merged.

Metric	Random Forest (RF)		Decision Tree (DT)		K-Nearest Neighbors (KNN)	
	Pre	Post	Pre	Post	Pre	Post
Accuracy Before and After Adversarial Training						
Benign Detection	0.99				0.98	
Phishing Detection	0.99				0.98	
Scamming Detection	0.99				0.98	
Instance Counts Before Adversarial Training						
Benign	Phishing		Scamming		Fake ICO	
11,431	629		2,189		2,189	
Misclassification Rates After Adversarial Training						
Benign Misclassifications						
Benign → Phishing	0.00%				0.22%	
Benign → Scamming	0.00%				0.03%	
Phishing → Benign	0.00%				8.43%	
Scamming → Benign	0.00%		0.14%		0.00%	
Phishing Misclassifications						
Phishing → Benign	0.00%		1.75%		9.70%	
Phishing → Scamming	0.00%				8.11%	
Scamming → Phishing	0.32%		0.55%		0.73%	
Scamming Misclassifications						
Scamming → Benign	0.00%				2.69%	
Scamming → Phishing	0.00%				1.14%	
Phishing → Scamming	4.77%		6.20%		7.79%	
Fake ICO → Scamming	0.00%				100%	

effective in restoring and even enhancing the robustness of these models against adversarial attacks, particularly for the benign class, though KNN still showed slight residual vulnerabilities.

Phishing Class. Phishing class accuracy was drastically affected by targeted attacks, with RF and DT dropping to 0.01 and KNN to 0.02—reflecting a more than 95% decrease in performance (Table 6). However, after adversarial training, these models showed significant resilience. Both RF and DT rebounded to an accuracy of 0.99, and KNN improved to 0.98, indicating a near-complete recovery (Table 12). Misclassifications were nearly eliminated, with RF achieving 0% errors, and DT showing minimal errors. KNN, while much improved, still displayed some vulnerability with small but noticeable misclassification rates. Overall, adversarial training has proven highly effective, particularly for RF and DT, though KNN may benefit from further refinement to achieve similar levels of robustness.

Scamming Class. Scamming class accuracy was significantly impacted by targeted adversarial attacks, with RF and DT accuracies dropping from 0.99 and 0.98 to 0.14, and KNN plummeting from 0.67 to 0.07—representing drops of over 85% and 90%, respectively (Table 6). These attacks led to a substantial number of scamming instances being misclassified, particularly in the KNN model, where 1,214 instances were incorrectly classified as benign. However, after applying adversarial training, the models demonstrated significant recovery (Table 12). RF and DT accuracies returned to 0.99, and KNN improved to 0.98, showing a strong restoration of performance. Misclassifications were greatly reduced, with RF achieving a 0% misclassification rate for scamming instances, and DT and KNN reducing errors to minimal levels. Despite these improvements, KNN still showed some residual vulnerability, with 2.69% of scamming instances misclassified as benign and 1.14% as phishing. Overall, adversarial training proved highly effective for most models, particularly RF and DT, though KNN may still benefit from further refinement to achieve comparable robustness.

8 DISCUSSION

Feature Selection for Optimal Classification. This study examines how feature selection affects ML models in transaction classification. The choice of features impacts a model's ability to classify transactions and resist manipulation. Among the tested features, **timestamp** and **value** impacted model performance. Their stability under manipulation suggests they capture transaction details that differentiate legitimate from fraudulent activities. Timestamp changes affected all models, with RF handling them better than DT and KNN. This indicates that time-related features contribute to fraud detection, as transaction timing can reveal patterns linked to fraudulent behavior. The accuracy results show that transaction time and date help identify fraud.

Value manipulations, including uniform and proportional changes, affected model performance, especially under uniform. RF and DT models showed accuracy drops, while KNN remained stable. After adversarial training, RF and DT regained accuracy, showing that adversarial training helps reduce the impact of value changes. The results indicate that transaction value is an important feature in classification. The ability to manipulate transaction value without detection suggests a potential risk, as it may enable bypassing of fraud detection systems.

Most Resistant Features to Adversarial Attacks. The results show that **address features**, specifically the From and To addresses, are less affected by adversarial attacks than other features. This may be because address features capture transaction patterns that are harder to change. When these features were manipulated, DT and RF models showed some accuracy reduction, while KNN was more affected. The differences across models suggest that how addresses are represented and how each model processes them influence their resistance to attacks. Address features may hold relationship patterns between transactions that remain stable even when altered.

In addition, the analysis showed that while important for accuracy, temporal features were also resistant to manipulations. This suggests that temporal data capture transaction patterns that are harder to alter without making significant changes. Shifting timestamps led to accuracy declines, but the impact was smaller than value manipulations. This indicates that temporal features play a key role in classification and are less vulnerable to adversarial attacks due to the complexity and uniqueness of timestamp data.

Adversarial training is an effective way to improve model robustness against targeted attacks. Exposing models to AEs during training made them more resistant to critical features such as **timestamp** and **value** disruptions. This approach helped compare the effects of such attacks, showing that adversarial training can be a practical method for reinforcing phishing detection models against exploitable weaknesses.

Best Combinations. For a more resilient classification, combining **temporal and address features** has proven to be effective. This integration takes advantage of both temporal patterns and relational data, creating models that are less vulnerable to adversarial interference. These features provide a layer of protection. Temporal features capture transaction timing and distribution patterns, while address features offer a stable relational structure that is harder to manipulate. This approach strengthens the model's accuracy, even when adversarial attacks attempt to exploit either feature.

Similarly, combining **temporal features with financial features**, such as transaction value and gas price, further enhances robustness. Although financial features can be uniformly manipulated, combining them with temporal data provides additional context that improves the model's resilience. The temporal features help to contextualize the financial data, mitigating the impact of adversarial value manipulations.

The following recommendations can be drawn for the effective and robust classification of transactions, especially in adversarial attacks. ① **Focus on Temporal and Address Features:**

Timestamp and address data could be key features in classification models, as they offer resilience against adversarial manipulations and play a crucial role in maintaining accuracy. ② **Integrate Financial Features with Temporal Data:** Use financial transaction data with temporal features to improve robustness and provide a transactional context that helps counteract adversarial manipulations. ③ **Adopt a Multi-Feature Approach:** Utilize a combination of diverse feature types to leverage their respective strengths and ensure a balanced, resilient classification model capable of withstanding various adversarial manipulations.

9 CONCLUSION

This study examines how machine learning models can be affected by adversarial attacks in detecting fraudulent transactions. The results show that different classifiers have varying levels of vulnerability. RF is more resistant, while DT and KNN are more easily influenced by these attacks. This highlights the need for careful feature selection and adversarial training to improve model robustness. These techniques help reduce the impact of attacks and improve classification accuracy. Using temporal and address features along with financial data strengthens model defenses and supports reliable classification under adversarial conditions. The adversarial training approach tested in this study also helps reduce the effects of targeted attacks, making it a valuable defense strategy. Future research will explore more complex attack scenarios and apply these methods to other financial platforms to evaluate their applicability and effectiveness.

ACKNOWLEDGEMENT

An earlier version of this work appeared in proceedings of WISA 2024 as “Simple Perturbations Subvert Ethereum Phishing Transactions Detection” [11].

REFERENCES

- [1] Ahmed Abusnaina, Mohammed Abuhamad, Hisham Alasmay, Afsah Anwar, Rhongho Jang, Saeed Salem, DaeHun Nyang, and David Mohaisen. 2022. DL-FHMC: Deep Learning-Based Fine-Grained Hierarchical Learning Approach for Robust Malware Classification. *IEEE Trans. Dependable Secur. Comput.* 19, 5 (2022), 3432–3447. <https://doi.org/10.1109/TDSC.2021.3097296>
- [2] Ahmed Abusnaina, Afsah Anwar, Sultan Alshamrani, Abdulrahman Alabduljabbar, RhongHo Jang, DaeHun Nyang, and David Mohaisen. 2022. Systematically Evaluating the Robustness of ML-based IoT Malware Detection Systems. In *25th International Symposium on Research in Attacks, Intrusions and Defenses, RAID*. ACM, 308–320. <https://doi.org/10.1145/3545948.3545960>
- [3] Ahmed Abusnaina, Rhongho Jang, Aminollah Khormali, DaeHun Nyang, and David Mohaisen. 2020. DFD: Adversarial Learning-based Approach to Defend Against Website Fingerprinting. In *39th IEEE Conference on Computer Communications, INFOCOM*. IEEE, 2459–2468. <https://doi.org/10.1109/INFOCOM41043.2020.9155465>
- [4] Ahmed Abusnaina, Aminollah Khormali, Hisham Alasmay, Jeman Park, Afsah Anwar, and Aziz Mohaisen. 2019. Adversarial Learning Attacks on Graph-based IoT Malware Detection Systems. In *39th IEEE International Conference on Distributed Computing Systems, ICDCS*. IEEE, 1296–1305. <https://doi.org/10.1109/ICDCS.2019.00130>
- [5] Ahmed Abusnaina, Yuhang Wu, Sunpreet S. Arora, Yizhen Wang, Fei Wang, Hao Yang, and David Mohaisen. 2021. Adversarial Example Detection Using Latent Neighborhood Graph. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV*. IEEE, 7667–7676. <https://doi.org/10.1109/ICCV48922.2021.00759>
- [6] Ayodeji Adeniran, Mohammed Alkinoon, and David Mohaisen. 2023. Understanding the Utilization of Cryptocurrency in the Metaverse and Security Implications. In *Computational Data and Social Networks - 12th International Conference, CSoNet (Lecture Notes in Computer Science, Vol. 14479)*. Springer, 268–281. https://doi.org/10.1007/978-981-97-0669-3_25
- [7] Ayodeji Adeniran, Kieran Human, and David Mohaisen. 2024. Dissecting the Infrastructure Used in Web-based Cryptojacking: A Measurement Perspective. In *International Conference Information Security Applications, WISA*. <https://doi.org/10.48550/arXiv.2408.03426>
- [8] Rachit Agarwal, Tanmay Thapliyal, and Sandeep K. Shukla. 2022. Analyzing Malicious Activities and Detecting Adversarial Behavior in Cryptocurrency based Permissionless Blockchains: An Ethereum Usecase. *Distributed Ledger Technol. Res. Pract.* 1, 2 (2022), 1–21. <https://doi.org/10.1145/3549527>

- [9] Salam Al-Emari, Mohammed Anbar, Yousef K. Sanjalawe, and Selvakumar Manickam. 2020. A Labeled Transactions-Based Dataset on the Ethereum Network. In *Advances in Cyber Security - Second International Conference, ACeS (Communications in Computer and Information Science, Vol. 1347)*. Springer, 61–79. https://doi.org/10.1007/978-981-33-6835-4_5
- [10] Hisham Alasmary, Ahmed Abusnaina, Rhongho Jang, Mohammed Abuhamad, Afsah Anwar, DaeHun Nyang, and David Mohaisen. 2020. Soteria: Detecting Adversarial Examples in Control Flow Graph-based Malware Classifiers. In *40th IEEE International Conference on Distributed Computing Systems, ICDCS*. IEEE, 888–898. <https://doi.org/10.1109/ICDCS47774.2020.00089>
- [11] Ahod Alghuried and David Mohaisen. 2024. Simple Perturbations Subvert Ethereum Phishing Transactions Detection: An Empirical Analysis. *CoRR* abs/2408.03441 (2024). <https://doi.org/10.48550/ARXIV.2408.03441>
- [12] Hung Ba. 2019. Improving Detection of Credit Card Fraudulent Transactions using Generative Adversarial Networks. *CoRR* abs/1907.03355 (2019). <http://arxiv.org/abs/1907.03355>
- [13] Arjun Nitin Bhagoji, Warren He, Bo Li, and Dawn Song. 2018. Practical Black-Box Attacks on Deep Neural Networks Using Efficient Query Mechanisms. In *Computer Vision - ECCV (Lecture Notes in Computer Science, Vol. 11216)*. Springer, 158–174. https://doi.org/10.1007/978-3-030-01258-8_10
- [14] Yair Carmon, Aditi Raghunathan, Ludwig Schmidt, John C. Duchi, and Percy Liang. 2019. Unlabeled Data Improves Adversarial Robustness. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems*. 11190–11201. <https://proceedings.neurips.cc/paper/2019/hash/32e0bd1497aa43e02a42f47d9d6515ad-Abstract.html>
- [15] Francesco Cartella, Orlando Anunciação, Yuki Funabiki, Daisuke Yamaguchi, Toru Akishita, and Olivier Elshocht. 2021. Adversarial Attacks for Tabular Data: Application to Fraud Detection and Imbalanced Data. In *Proceedings of the Workshop on Artificial Intelligence Safety 2021 (SafeAI 2021) co-located with the Thirty-Fifth AAAI Conference on Artificial Intelligence (AAAI, Vol. 2808)*. CEUR-WS.org. https://ceur-ws.org/Vol-2808/Paper_4.pdf
- [16] Yuanyuan Chen, Jing Qiu, Xiaojiang Du, Lihua Yin, and Zhihong Tian. 2020. Security of Mobile Multimedia Data: The Adversarial Examples for Spatio-temporal Data. *Comput. Networks* 181 (2020), 107432. <https://doi.org/10.1016/J.COMNET.2020.107432>
- [17] Zhen Chen, Sheng-Zheng Liu, Jia Huang, Yu-Han Xiu, Hao Zhang, and Haixia Long. 2024. Ethereum Phishing Scam Detection Based on Data Augmentation Method and Hybrid Graph Neural Network Model. *Sensors* 24, 12 (2024), 4022. <https://doi.org/10.3390/S24124022>
- [18] Francesco Croce, Maksym Andriushchenko, Vikash Sehwal, Edoardo DeBenedetti, Nicolas Flammarion, Mung Chiang, Prateek Mittal, and Matthias Hein. 2021. RobustBench: a standardized adversarial robustness benchmark. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks*. <https://arxiv.org/pdf/2010.09670>
- [19] Pablo de Juan Fidalgo, Carmen Camara, and Pedro Peris-Lopez. 2022. Generation and Classification of Illicit Bitcoin Transactions. In *Proceedings of the International Conference on Ubiquitous Computing & Ambient Intelligence, UCAmI (Lecture Notes in Networks and Systems, Vol. 594)*. Springer, 1086–1097. https://doi.org/10.1007/978-3-031-21333-5_108
- [20] Yifan Ding, Liqiang Wang, Huan Zhang, Jinfeng Yi, Deliang Fan, and Boqing Gong. 2019. Defending Against Adversarial Attacks Using Random Forest. In *IEEE Conference on Computer Vision and Pattern Recognition Workshops, CVPR*. Computer Vision Foundation / IEEE, 105–114. <https://doi.org/10.1109/CVPRW.2019.00019>
- [21] Ivan Fursov, Matvey Morozov, Nina Kaplounkhaya, Elizaveta Kovtun, Rodrigo Rivera-Castro, Gleb Gusev, Dmitry Babaev, Ivan Kireev, Alexey Zaytsev, and Evgeny Burnaev. 2021. Adversarial Attacks on Deep Models for Financial Transaction Records. In *KDD '21: The 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. ACM, 2868–2878. <https://doi.org/10.1145/3447548.3467145>
- [22] Daniel Gibert, Luca Demetrio, Giulio Zizzo, Quan Le, Jordi Planes, and Battista Biggio. 2024. Certified Adversarial Robustness of Machine Learning-based Malware Detectors via (De)Randomized Smoothing. *CoRR* abs/2405.00392. <https://doi.org/10.48550/ARXIV.2405.00392>
- [23] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. 2015. Explaining and Harnessing Adversarial Examples. In *3rd International Conference on Learning Representations, ICLR*, Yoshua Bengio and Yann LeCun (Eds.). <http://arxiv.org/abs/1412.6572>
- [24] Qingyu Guo, Zhao Li, Bo An, Pengrui Hui, Jiaming Huang, Long Zhang, and Mengchen Zhao. 2019. Securing the Deep Fraud Detector in Large-Scale E-Commerce Platform via Adversarial Machine Learning Approach. In *The World Wide Web Conference, WWW*. ACM, 616–626. <https://doi.org/10.1145/3308558.3313533>
- [25] Arkan Hammoodi Hasan Kabla, Mohammed Anbar, Selvakumar Manickam, Taief Alaa Alamiedy, Peterson Bernabe Cruspe, Ahmed K. Al-Ani, and Shankar Karuppayah. 2022. Applicability of Intrusion Detection System on Ethereum Attacks: A Comprehensive Review. *IEEE Access* 10 (2022), 71632–71655. <https://doi.org/10.1109/ACCESS.2022.3188637>
- [26] Arkan Hammoodi Hasan Kabla, Mohammed Anbar, Selvakumar Manickam, and Shankar Karuppayah. 2022. Eth-PSD: A Machine Learning-Based Phishing Scam Detection Approach in Ethereum. *IEEE Access* 10 (2022), 118043–118057.

- <https://doi.org/10.1109/ACCESS.2022.3220780>
- [27] Dan Li, Dacheng Chen, Jonathan Goh, and See-Kiong Ng. 2018. Anomaly Detection with Generative Adversarial Networks for Multivariate Time Series. *CoRR* abs/1809.04758 (2018). <http://arxiv.org/abs/1809.04758>
- [28] Jie Li, Tianqing Zhu, Wei Ren, and Kim-Kwang Raymond Choo. 2023. Improve individual fairness in federated learning via adversarial training. *Comput. Secur.* 132 (2023), 103336. <https://doi.org/10.1016/J.COSE.2023.103336>
- [29] Xiaodan Li, Yuefeng Chen, Yuan He, and Hui Xue. 2019. AdvKnn: Adversarial Attacks On K-Nearest Neighbor Classifiers With Approximate Gradients. *CoRR* abs/1911.06591 (2019). <http://arxiv.org/abs/1911.06591>
- [30] Jintao Luo, Jiwei Qin, Ruijin Wang, and Lu Li. 2024. A Phishing Account Detection Model via Network Embedding for Ethereum. *IEEE Trans. Circuits Syst. II Express Briefs* 71, 2 (2024), 622–626. <https://doi.org/10.1109/TCSII.2023.3267822>
- [31] Haifeng Lv and Yong Ding. 2023. Phishing detection on Ethereum via transaction subgraphs embedding. *IET Blockchain* 3, 4 (2023), 194–203. <https://doi.org/10.1049/BLC2.12034>
- [32] Alberto Mozo, Ángel González-Prieto, Antonio Pastor Perales, Sandra Gómez Canaval, and Edgar Talavera. 2021. Synthetic flow-based cryptomining attack generation through Generative Adversarial Networks. *CoRR* abs/2107.14776 (2021). <https://arxiv.org/abs/2107.14776>
- [33] Nina Narodytska and Shiva Prasad Kasiviswanathan. 2017. Simple Black-Box Adversarial Attacks on Deep Neural Networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops, CVPR*. IEEE Computer Society, 1310–1318. <https://doi.org/10.1109/CVPRW.2017.172>
- [34] Cuong Phuc Ngo, Amadeus Aristo Winarto, Connie Khor Li Kou, Sojeong Park, Farhan Akram, and Hwee Kuan Lee. 2019. Fence GAN: Towards Better Anomaly Detection. In *31st IEEE International Conference on Tools with Artificial Intelligence, ICTAI*. IEEE, 141–148. <https://doi.org/10.1109/ICTAI.2019.00028>
- [35] Vinicius C. Oliveira, Júlia Almeida Valadares, José Eduardo de Azevedo Sousa, Alex Borges Vieira, Heder Soares Bernardino, Saulo Moraes Villela, and Glauber Dias Gonçalves. 2021. Analyzing Transaction Confirmation in Ethereum Using Machine Learning Techniques. *SIGMETRICS Perform. Evaluation Rev.* 48, 4 (2021), 12–15. <https://doi.org/10.1145/3466826.3466832>
- [36] Nicolas Papernot, Patrick D. McDaniel, Xi Wu, Somesh Jha, and Ananthram Swami. 2016. Distillation as a Defense to Adversarial Perturbations Against Deep Neural Networks. In *IEEE Symposium on Security and Privacy, SP*. IEEE Computer Society, 582–597. <https://doi.org/10.1109/SP.2016.41>
- [37] Elnaz Rabieinejad, Abbas Yazdinejad, Reza M. Parizi, and Ali Dehghantanha. 2023. Generative Adversarial Networks for Cyber Threat Hunting in Ethereum Blockchain. *Distributed Ledger Technol. Res. Pract.* 2, 2 (2023), 1–19. <https://doi.org/10.1145/3584666>
- [38] Vaishali Ravindranath, M. K. Nallakaruppan, M. Lawanya Shri, Balamurugan Balusamy, and Siddhartha Bhattacharyya. 2024. Evaluation of performance enhancement in Ethereum fraud detection using oversampling techniques. *Appl. Soft Comput.* 161 (2024), 111698. <https://doi.org/10.1016/J.ASOC.2024.111698>
- [39] Muhammad Saad and David Mohaisen. 2023. Analyzing In-browser Cryptojacking. *CoRR* abs/2304.13253 (2023). <https://doi.org/10.48550/ARXIV.2304.13253>
- [40] Muhammad Saad, Laurent Njilla, Charles A. Kamhoua, Joongheon Kim, DaeHun Nyang, and Aziz Mohaisen. 2019. Mempool optimization for Defending Against DDoS Attacks in PoW-based Blockchain Systems. In *IEEE International Conference on Blockchain and Cryptocurrency, ICBC*. IEEE, 285–292. <https://doi.org/10.1109/BLOC.2019.8751476>
- [41] Muhammad Saad, Jeffrey Spaulding, Laurent Njilla, Charles A. Kamhoua, Sachin Shetty, DaeHun Nyang, and David Mohaisen. 2020. Exploring the Attack Surface of Blockchain: A Comprehensive Survey. *IEEE Commun. Surv. Tutorials* 22, 3 (2020), 1977–2008. <https://doi.org/10.1109/COMST.2020.2975999>
- [42] Yousef K. Sanjalawe and Salam Al-Emari. 2023. Abnormal Transactions Detection in the Ethereum Network Using Semi-Supervised Generative Adversarial Networks. *IEEE Access* 11 (2023), 98516–98531. <https://doi.org/10.1109/ACCESS.2023.3313630>
- [43] Dule Shu, Nandi O. Leslie, Charles A. Kamhoua, and Conrad S. Tucker. 2020. Generative adversarial attacks against intrusion detection systems using active learning. In *Proceedings of the 2nd ACM Workshop on Wireless Security and Machine Learning, WiseML@WiSec*. ACM, 1–6. <https://doi.org/10.1145/3395352.3402618>
- [44] Samuel Henrique Silva and Peyman Najafirad. 2020. Opportunities and Challenges in Deep Learning Adversarial Robustness: A Survey. *CoRR* abs/2007.00753 (2020). <https://arxiv.org/abs/2007.00753>
- [45] Harsh Jot Singh and Abdelhakim Senhaji Hafid. 2019. Prediction of Transaction Confirmation Time in Ethereum Blockchain Using Machine Learning. In *Blockchain and Applications - International Congress, BLOCKCHAIN (Advances in Intelligent Systems and Computing, Vol. 1010)*. Springer, 126–133. https://doi.org/10.1007/978-3-030-23813-1_16
- [46] David Stutz, Matthias Hein, and Bernt Schiele. 2019. Disentangling Adversarial Robustness and Generalization. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*. Computer Vision Foundation / IEEE. <https://doi.org/10.1109/CVPR.2019.00714>
- [47] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. 2014. Intriguing properties of neural networks. In *2nd International Conference on Learning Representations, ICLR*.

<http://arxiv.org/abs/1312.6199>

- [48] Runnan Tan, Qingfeng Tan, Qin Zhang, Peng Zhang, Yushun Xie, and Zhao Li. 2023. Ethereum fraud behavior detection based on graph neural networks. *Computing* 105, 10 (2023), 2143–2170. <https://doi.org/10.1007/S00607-023-01177-7>
- [49] Daniël Vos and Sicco Verwer. 2021. Efficient Training of Robust Decision Trees Against Adversarial Examples. In *Proceedings of the 38th International Conference on Machine Learning, ICML (Proceedings of Machine Learning Research, Vol. 139)*. PMLR, 10586–10595. <http://proceedings.mlr.press/v139/vos21a.html>
- [50] Jiajing Wu, Qi Yuan, Dan Lin, Wei You, Weili Chen, Chuan Chen, and Zibin Zheng. 2022. Who Are the Phishers? Phishing Scam Detection on Ethereum via Network Embedding. *IEEE Trans. Syst. Man Cybern. Syst.* 52, 2 (2022), 1156–1166. <https://doi.org/10.1109/TSMC.2020.3016821>
- [51] Cihang Xie, Yuxin Wu, Laurens van der Maaten, Alan L. Yuille, and Kaiming He. 2019. Feature Denoising for Improving Adversarial Robustness. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*. Computer Vision Foundation / IEEE, 501–509. <https://doi.org/10.1109/CVPR.2019.00059>
- [52] Jin Yang, Tao Li, Gang Liang, Yunpeng Wang, Tianyu Gao, and Fangdong Zhu. 2020. Spam transaction attack detection model based on GRU and WGAN-div. *Comput. Commun.* 161 (2020), 172–182. <https://doi.org/10.1016/J.COMCOM.2020.07.031>
- [53] Xikang Yang, Biyu Zhou, Xuehai Tang, Xiaodan Zhang, Jizhong Han, and Songlin Hu. 2023. Translets: Toward Explainable Phishing Fraud Detection in Ethereum. In *IEEE International Conference on High Performance Computing & Communications, Data Science & Systems, Smart City & Dependability in Sensor, Cloud & Big Data Systems & Application, HPCC/DSS/SmartCity/DependSys*. IEEE, 582–591. <https://doi.org/10.1109/HPCC-DSS-SMARTCITY-DEPENDSYS60770.2023.00085>
- [54] Keting Yin and Binglong Ye. 2023. Phishing Scam Detection for Ethereum Based on Community Enhanced Graph Convolutional Networks. In *Neural Information Processing - 30th International Conference, ICONIP (Communications in Computer and Information Science, Vol. 1965)*. Springer, 191–206. https://doi.org/10.1007/978-981-99-8145-8_16
- [55] Liangheng Zhang, Congmei Jiang, Zhaosen Chai, and Yu He. 2024. Adversarial attack and training for deep neural network based power quality disturbance classification. *Eng. Appl. Artif. Intell.* 127, Part A (2024), 107245. <https://doi.org/10.1016/J.ENGAPPAL.2023.107245>
- [56] Francesco Zola, Jan Lukas Bruse, Xabier Etxeberria Barrio, Mikel Galar, and Raul Orduna Urrutia. 2020. Generative Adversarial Networks for Bitcoin Data Augmentation. In *2nd Conference on Blockchain Research & Applications for Innovative Networks and Services*. IEEE, 136–143. <https://doi.org/10.1109/BRAINS49436.2020.9223269>
- [57] Francesco Zola, Lander Seguro-Gil, Jan L. Bruse, Mikel Galar, and Raul Orduna Urrutia. 2022. Attacking Bitcoin anonymity: generative adversarial networks for improving Bitcoin entity classification. *Appl. Intell.* 52, 15 (2022), 17289–17314. <https://doi.org/10.1007/S10489-022-03378-7>