

An Explainable AI Model for the Detecting Malicious Smart Contracts Based on EVM Opcode Based Features

*Note: Sub-titles are not captured for <https://ieeexplore.ieee.org> and should not be used

1st Roopak Surendran
Security Researcher

Abstract—Hackers may create malicious solidity programs and deploy it in the Ethereum block chain. These malicious smart contracts try to attack legitimate programs by exploiting its vulnerabilities such as reentrancy, tx.origin attack, bad randomness, delegatecall and so on. This may lead to drain of the funds, denial of service and so on. Hence, it is necessary to identify and prevent the malicious smart contract before deploying it into the blockchain. In this paper, we propose an ML based malicious smart contract detection mechanism by analyzing the EVM opcodes. After balancing the opcode frequency dataset with SMOTE algorithm, we transformed opcode frequencies to the binary values (0,1) using an entropy based supervised binning method. Then, an explainable AI model is trained with the proposed binary opcode based features. From the implementations, we found that the proposed mechanism can detect 99% of malicious smart contracts with a false positive rate of only 0.01. Finally, we incorporated LIME algorithm in our classifier to justify its predictions. We found that, LIME algorithm can explain why a particular smart contract app is declared as malicious by our ML classifier based on the binary value of EVM opcodes.

Index Terms—component, formatting, style, styling, insert.

I. INTRODUCTION

Smart contracts are the programs deployed in a block chain network that runs only when some predefined conditions or rules met [1], [2]. It enables to automate the execution of the business transactions without the involvement of a third party. Solidity is one of the most popular language for implementing smart contracts which can be deployed in the Ethereum [3], a decentralized block chain network. Solidity is an object-oriented programming language adopted from the concepts of JavaScript, C++ and python [4]. The programs developed using solidity language run in a virtual environment called EVM (Ethereum Virtual Machine).

Like other programming languages, security vulnerabilities may occur in solidity programming language also [3]. An attacker can create and deploy a malicious contract that exploits the vulnerabilities in another deployed legitimate contacts. This might lead to the attacks such as drainage of funds, DoS/DDoS and so on. Hence, it is required to detect the malicious contract and block it before deploying it into the block chain network. There are more than ten kinds of

vulnerabilities existing in a solidity program. Some of the major vulnerabilities are:

- **Bad Randomness:** With this vulnerability, an attacker can predict the output of solidity program functions that relies on randomness;
- **Denial of Service:** A smart contract can deny the functions of other legitimate contracts by performing DoS attack.;
- **Tx.Origin Attack:** It is a kind of phishing attack which deceives the contract owner to perform the privileged operations;
- **Reentrancy Attack:** This is a recursive process which allows unauthorized transfer of funds from the contract of legitimate user to the malicious contract;
- **Forced Ether Reception:** With this vulnerability, an attacker can forcefully send ether/s to any other vulnerable contracts;
- **Hiding Malicious Code:** This vulnerability allows to hide a malicious code inside another contact.

There are several tools such as mythril, slither and so on are introduced for analyzing vulnerabilities in a solidity program [5], [6]. These existing tools can only identify known vulnerabilities in a solidity application. It cannot discover unknown or zero-day vulnerabilities in an application. In the proposed mechanism, we use ML algorithms to identify whether a smart contract program tries to attack other applications or not. Further, the ML classifier can justify its decision using explainable AI algorithms.

Existing mechanisms uses solidity program code or EVM opcode-based features for identifying the vulnerabilities. While deploying the solidity program in EVM, it is compiled to the low level program called EVM byte code. After deploying, a developers can remove their actual high level solidity program from the Ethereum network. Hence, the solidity program codes are not always available for analysis. Hence, in this paper, we used EVM opcode-based features for identifying whether a contact is malicious or not.

In our collected smart contract dataset, the percentage of malicious contacts are below 1%. That is, the dataset is highly imbalanced. Training an ML classifier with a highly

imbalanced dataset may lead to the biased prediction towards the majority class. In order to overcome this limitation, we use data oversampling techniques such as SMOTE to balance the dataset by adding more malicious contracts. After balancing the opcode frequency dataset with SMOTE algorithm, we transformed opcode frequencies to the binary values (0 or 1) using an entropy based supervised binning method. Then, an explainable AI model is trained with the proposed binary opcode based features. From the implementations, we get 0.99 recall in detecting malicious contract with a false positive rate of only less than 1%. Further, we used explainable AI algorithm such as LIME to justify the predictions of the sample malicious and legitimate contract. We found that, our model can explain its predictions based on the values of key EVM opcode level features.

The rest of the paper is organized as follows. In Section 2, literature review has been given. In Section 3, our proposed mechanism has been given. The results and discussions are given in Section 4. In Section 5, we conclude our research and future directions are given.

II. LITERATURE REVIEW

In this section, we discuss about the existing solidity vulnerability analyzers in detail. The existing vulnerability analyzers inspect either the solidity program code or EVM opcode (or a combination of both) for detecting the vulnerabilities.

In Contractward [7], the relationship among solidity code, opcode bigram features are used for identifying the vulnerabilities. In Eth2Vec [8], evm opcodes are used as features of ML classifier to identify whether a contract is vulnerable or not. Grieco et al. [9] developed a tool called Echidna to detect vulnerabilities by analyzing the solidity code. Ethainter [10] can detect composite vulnerabilities by checking information flow with data sanitization in solidity code. Ethver [11], a formal verification-based vulnerability analyzer which takes solidity codes as input and translate it into MDP (Markov Decision Process) models. Then, it is verified using PRISM model checker. Yang et al. [12] proposed a mechanism called FEther for identifying smart contract vulnerabilities by combining symbolic analysis and higher order logic theorem-proving. In GasGuage [13], fuzz testing is used to detect Ethereum gas related vulnerabilities. In Musc [14], the mutants of smart contracts are generated from AST (Abstract Syntax Tree) and these AST generated mutants are converted to solidity code for testing. Gupta et al. [15], suggested a deep learning based mechanism to detect vulnerable contracts from the EVM opcodes.

All of the existing mechanisms discussed the problem of detecting vulnerabilities in smart contracts. They did not make an attempt to detect the malicious smart contract which tries to exploit vulnerabilities. Moreover, the classifier decisions has not been justified in all these methods. In order to overcome this limitation, in this work, we propose a mechanism to transform opcode frequency based features to binary features for interpreting classifier decisions using explainable AI algorithm. From the implementations, we found that our model can

effectively interpret why a particular smart contract is declared as malicious by analyzing binary value of EVM opcodes.

III. METHODOLOGY

In this section, we discuss about our malicious contract detection mechanism in detail. The proposed mechanisms consist of three steps. In the first step, we extract EVM opcodes of malware and goodware apps in our dataset. The percentage of malicious contacts may be very lesser than the legitimate contracts. Hence, in the second step, we use oversampling technique such as SMOTE, to balance the dataset by adding more malicious contracts. Then, we train an ML classifier and tested with the features of unknown smart contracts. Further, explainable AI is used to justify the classifier predictions. The architecture of our mechanism is given in Figure 1.

A. Dataset Oversampling Phase

The total number of malicious contracts is very lower as compared to benign contracts. That is, there are only below 1% of malicious contracts in a dataset. Training with a highly imbalanced dataset may lead to the performance bias towards the majority class. Hence, in order to balance the training dataset, we used SMOTE algorithm [16]. The steps in data oversampling using SMOTE algorithm are given below.

- 1) Step 1: Extract the frequency values of n EVM opcodes X_i for $i = 1, 2, \dots, n$ in the malware and legitimate contracts in the dataset.
- 2) Step 2: Build a minority class set X , for each $x \in X$, the nearest neighbor of X is calculated using Euclidean distance formula. The Euclidean distance $D(A, B)$ of any given two points $A = (a_1, a_2, \dots, a_n)$ and $B = (b_1, b_2, \dots, b_n)$ is calculated as,

$$D(A, B) = \sqrt{\sum_{i=1}^n (a_i - b_i)^2}.$$

- 3) Step 3: Set the sampling rate N according to the imbalanced population, for each $x \in X$, N samples are selected for k nearest neighbors to construct a new set Y .
- 4) Step 4: For each sample $y_k \in Y$ for $k = 1, 2, \dots, N$, new samples are generated by the following formula $X + t \cdot |x - y_k|$, where $0 \leq t \leq 1$.

B. Feature Transformation Using Entropy Based Supervised Binning Algorithm

In both training and test dataset, we use binning technique to transform discrete frequency based opcode features to binary categorical features. Assume that, you have a train or test dataset with a continuous feature $X = \{X_i : i = 1, 2, 3, \dots, n\}$ and a corresponding target variable $Y_i \in \{0, 1\}$ denotes its label (1 represents malicious and 0 represents legitimate). Entropy is a measure of uncertainty or disorder

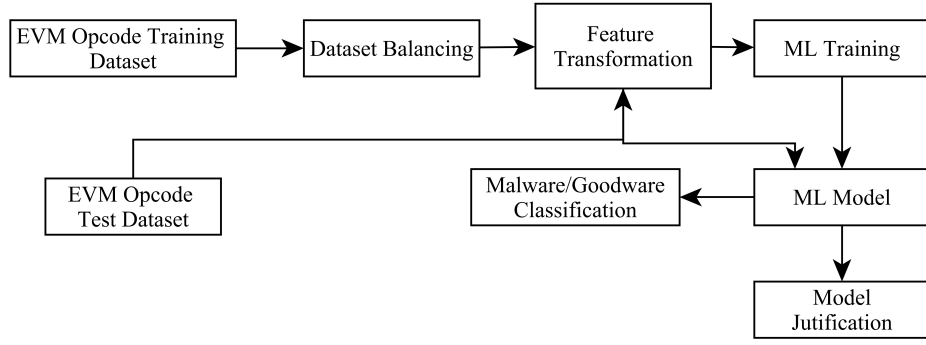


Fig. 1: Proposed Malicious Contract Detection Mechanism

in a set of data. For a given set of data S with k classes, the entropy $H(S)$ is calculated as:

$$H(S) = - \sum_{i=1}^k p_i \cdot \log_2(p_i),$$

where p_i is the proportion of data points in class i . Information Gain is used to measure the reduction in entropy achieved by partitioning a dataset based on a certain attribute (e.g., a continuous feature). The information gain is calculated as:

$$\text{Information Gain} = H(S) - \sum_j \frac{|S_j|}{|S|} \cdot H(S_j),$$

where S_j is the subset of S for which feature X falls into bin j . The supervised binning algorithm is given below:

- Start by sorting the feature X in ascending order;
- Calculate information gain for each possible split point;
- Choose the split point that maximizes the information gain, effectively dividing the data into two bins;

After identifying the splitting points $S(X_i)$ to the n frequency based opcode features X_i for $i = 1, 2, 3, \dots, n$ in the dataset, we transform opcode frequencies $O(X_i)$ of X_i to binary features $B(X_i)$ by comparing with the corresponding splitting points $S(X_i)$. That is,

$$B(X_i) = \begin{cases} 1 & \text{if } O(X_i) \geq S(X_i); \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

C. Interpreting Blackbox ML Models Using LIME Algorithm

In this section, we discuss about machine learning based malicious smart contract detection from opcode sequences. After constructing a balanced data set of binary opcode based features, we train an ML classifier. This ML classifier can detect whether an unknown contract is malicious or not. We incorporate explainable AI algorithm such as LIME [17] to justify the classifier predictions. The explanation $E(x)$ of a feature vector x produced by LIME is following:

$$E(x) = \arg \min_{g \in G} L(f, g, \pi_x) + \Omega(g),$$

where f is the actual model, g is the linear explanation model, π_x is the proximity measure between an instance of z to x and $\Omega(g)$ is the model complexity. π_x gives the weights w

to the perturbed instances z' based on their distance from x . $L(f, g, \pi_x)$ is calculated as,

$$L(f, g, \pi_x) = \sum_{z, z' \in Z} \pi_x(z) (f(z) - g(z'))^2,$$

where $\pi_x(z) = e^{(-D(x, z)^2 / (\sigma^2))}$, $g(z') = w \cdot g(z)$ and $D(x, z)$ is the distance function.

IV. RESULTS AND DISCUSSIONS

In this section, we discuss the performance results of our ML based malicious smart contract detection mechanism and their explainability. We collected the opcode sequences of malicious and benign contracts from the forta dataset ¹. The collected dataset is highly imbalanced because there exists only a very few malicious contracts in the wild. The collected dataset consists of 139451 legitimate and 143 malicious contracts. In their mechanism, a logistic regression (LR) classifier was trained with the previously detected malicious opcode patterns of to detect malicious behavior in unseen smart contracts. Because of the high dataset imbalance, the classifier can detect only 59% of malicious smart contracts.

The proposed mechanism was implemented in a windows 11 PC with 128GB of memory. After collecting the dataset, we calculated the frequencies of opcodes in each malicious and legitimate contracts. The, we add a label '1' at the end of malicious contract and '0' at the end of legitimate contract. Finally, a CSV dataset file contain malicious and legitimate smart contract opcode frequency based feature vectors was created. The opcodes found in the smart contract programs are given in Table 1.

It is known that, an ML model will be biased to the features of majority class if we train it with a highly imbalanced dataset. Initially, we transformed frequency based opcode features to binary features using the entropy based supervised binning algorithm. Then, we used 70% of malicious smart contract (93 apps) and 138247 legitimate contract for training ML classifier. The training dataset has been balanced by adding equal number of synthetically generated malicious contract features using SMOTE algorithm. Then, we train an ML classifier in the SMOTE balanced dataset of transformed

¹<https://huggingface.co/datasets/forta/malicious-smart-contract-dataset>

SELFBALANCE	BASEFEE	MLOAD	MSTORE
SLOAD	SSTORE	JUMP	JUMPI
PC	MSIZE	GAS	RETURN
EXTCODEHASH	BLOCKHASH	COINBASE	TIMESTAMP
NUMBER	PREVRANDAO	GASLIMIT	CHAINID
BALANCE	ORIGIN	CALLER	CALLVALUE
CALLDATALOAD	CALLDATACOPY	CODESIZE	CODECOPY
GASPRICE	EXTCODESIZE	EXTCODECOPY	SLT
SGT	EQ	ISZERO	AND
OR	XOR	NOT	BYTE
SHL	SHR	SAR	KECCAK256
ADDRESS	CALLCODE	RETURN	DELEGATECALL
CREATE	STATICCALL	REVERT	SELFDESTRUCT
LOG	SWAP	DUP	PUSH
POP	STOP	ADD	MUL
SUB	CALL	DIV	SDIV
MOD	SMOD	ADDMOD	MULMOD
EXP	SIGNEXTEND	LT	GT
CREATE			

TABLE I: List of Opcodes in Solidity Program

binary opcode based features corresponding to 70% of malicious smart contracts and tested those of remaining 30% of malicious smart contracts for evaluation. The classifiers [18] used for evaluation are:

- Naive Bayes;
- Logistic Regression;
- K- Nearest Neighbour;
- Decision Tree.

The performance results in different machine learning classifiers are given in Table 2. From the Table 2, we can see that the proposed mechanism can detect 99% of malicious smart contracts with a false positive rate of only below 1%. The comparison results against the mechanism used in Forta is given in Table 3. From Table 3, we can see that our mechanism outperforms Forta with a recall rate of 0.99. In order to find out the relevant opcode based features, we used extra trees feature selection method in the opcode frequency dataset [19]. The relevant opcode features X_i , its frequencies $O(X_i)$, splitting points $S(X_i)$ and assigned binary labels $B(X_i)$ of a sample smart contract application is given in Table 4. In rest of our work, we use only these relevant opcodes for justifying the classifier predictions.

A. Illustrating Explainability of ML with a Malicious and Legitimate Contract

In this section, we illustrate the explainability of our ML classifier in a sample malicious and legitimate contract. We fed the binary opcode based features of the malicious and

legitimate contracts in a trained explainable AI model. Here, we made justification on the basis of top features those have a contribution value greater than zero. The classifier justification for a malicious contract is given in Table 5. From Table 5, we can see that, the majority of features supports malicious class label. In Figure 2, we can see that the total contributions of features support malicious contract is higher than that of legitimate. Hence, it is clear that the features corresponding to the smart contract is malicious. The classifier justification for legitimate contract is given in Table 6. From Table 6, we can see that, the majority of features supports legitimate class label. In Figure 3, we can see that the total contributions of features support legitimate contract is higher than that of malicious. Hence, it is clear that the features corresponding to the smart contract is legitimate.

V. CONCLUSION AND FUTURE SCOPE

In this section, we discuss about the limitations and future directions of our work. Here, we build an explainable ML classifier which is capable to detect malicious contracts and make justifications on its prediction. Our ML classifier can able to detect malicious smart contracts with a precision and F1 score of 0.99. With the help of explainable AI algorithm, we found that the frequencies of certain opcodes are very high in malicious contracts.

It is possible for an adversary to defeat our ML classifier using adversarial attacks [20]. In future, we will build an adversarial resistant ML classifier to detect adversarial malicious

Algorithm	TPR	FPR	Precision	Accuracy	F1Score
Naive Bayes	0.96	0.04	0.99	0.99	0.98
LR	0.99	0.01	0.99	0.99	0.99
DT	0.99	0.01	0.99	0.99	0.99
KNN	0.90	0.10	0.92	0.99	0.91

TABLE II: Performance of ML Classifiers in Our Dataset

Method	TPR	FPR	Precision	Accuracy	F1Score
Proposed Method	0.99	0.01	0.99	0.99	0.99
Forta	0.59	N/A	0.88	N/A	N/A

TABLE III: Performance of Our Mechanism against Forta

EVM Opcode	Opcode Frequency	Splitting Point	Binary Value
SSTORE	10	17	0
RETURNDATACOPY	1	17	0
SLT	0	0	0
EQ	16	32	0
OR	129	163	0
RETURN	17	75	0
DELEGATECALL	1	4	0
LOG	3	5	0
SUB	57	37	1
SLOAD	21	63	0

TABLE IV: Details of Key EVM Opcodes and Conditions for Binary Feature Transformation

contracts. For that, we will synthetically generate malicious contracts similar to the real world scenario using conditional GAN and train the adversarial resistant classifier with it [21]. Hence, in future, the proposed mechanism will be capable to detect unknown adversarial malicious contract attacks.

Our mechanism justifies the predictions on the basis of frequency values of certain opcodes. With respect to the change in evolving malicious contracts, these frequency range value of malicious and legitimate contracts may change in future. In this situation, explainable AI algorithm will not work [22]. In order to overcome this limitations, we will combine semantic analysis algorithms with explainable AI to make more reliable explanations for evolving smart contracts in future.

REFERENCES

- [1] B. K. Mohanta, S. S. Panda, D. Jena, An overview of smart contract and use cases in blockchain technology, in: 2018 9th international conference on computing, communication and networking technologies (ICCCNT), IEEE, 2018, pp. 1–4.
- [2] Z. Zheng, S. Xie, H. Dai, X. Chen, H. Wang, An overview of blockchain technology: Architecture, consensus, and future trends, in: 2017 IEEE international congress on big data (BigData congress), Ieee, 2017, pp. 557–564.
- [3] Z. A. Khan, A. S. Namin, Ethereum smart contracts: Vulnerabilities and their classifications, in: 2020 IEEE International Conference on Big Data (Big Data), IEEE, 2020, pp. 1–10.
- [4] M. Wohrer, U. Zdun, Smart contracts: security patterns in the ethereum ecosystem and solidity, in: 2018 International Workshop on Blockchain Oriented Software Engineering (IWBOSE), IEEE, 2018, pp. 2–8.
- [5] J. Feist, G. Grieco, A. Groce, Slither: a static analysis framework for smart contracts, in: 2019 IEEE/ACM 2nd International Workshop on Emerging Trends in Software Engineering for Blockchain (WETSEB), IEEE, 2019, pp. 8–15.
- [6] M. Di Angelo, G. Salzer, A survey of tools for analyzing ethereum smart contracts, in: 2019 IEEE international conference on decentralized applications and infrastructures (DAPPCON), IEEE, 2019, pp. 69–78.
- [7] W. Wang, J. Song, G. Xu, Y. Li, H. Wang, C. Su, Contractward: Automated vulnerability detection models for ethereum smart contracts, IEEE Transactions on Network Science and Engineering 8 (2) (2020) 1133–1144.
- [8] N. Ashizawa, N. Yanai, J. P. Cruz, S. Okamura, Eth2vec: learning contract-wide code representations for vulnerability detection on ethereum smart contracts, in: Proceedings of the 3rd ACM international symposium on blockchain and secure critical infrastructure, 2021, pp. 47–59.
- [9] G. Grieco, W. Song, A. Cygan, J. Feist, A. Groce, Echidna: effective, usable, and fast fuzzing for smart contracts, in: Proceedings of the 29th ACM SIGSOFT International Symposium on Software Testing and Analysis, 2020, pp. 557–560.
- [10] L. Brent, N. Grech, S. Lagouvardos, B. Scholz, Y. Smaragdakis, Ethainter: a smart contract security analyzer for composite vulnerabilities, in: Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation, 2020, pp. 454–469.
- [11] E. Mazurek, Ethver: Formal verification of randomized ethereum smart contracts, in: Financial Cryptography and Data Security. FC 2021 International Workshops: CoDecFin, DeFi, VOTING, and WTSC, Virtual Event, March 5, 2021, Revised Selected Papers 25, Springer, 2021, pp. 364–380.
- [12] Z. Yang, H. Lei, Fether: An extensible definitional interpreter for smart-

Feature	Actual Label Value	Supported Class Label	Contribution
RETURNDATACOPY	1	Malicious	0.25
RETURN	1	Malicious	0.21
SLT	1	Malicious	0.21
EQ	1	Legitimate	0.19
LOG	1	Legitimate	0.18
DELEGATECALL	0	Malicious	16
OR	0	Malicious	10
SSTORE	0	Malicious	0.03
SLOAD	0	Malicious	0.03
SUB	0	Malicious	0.03

TABLE V: LIME Algorithm Based Justification for Classifier Prediction in Malicious Contract

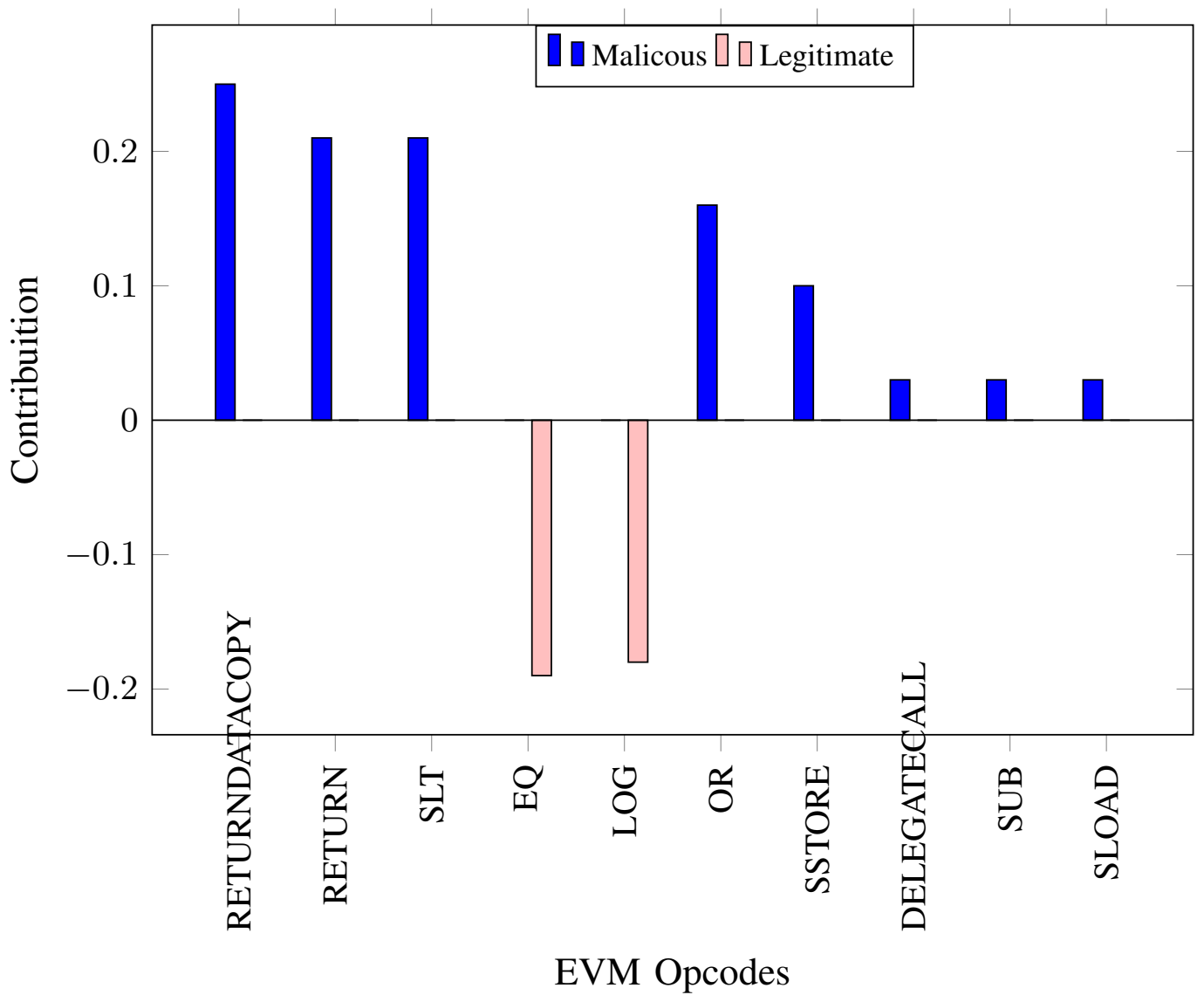


Fig. 2: Contributions of EVM Opcode Features in Malicious Contract

Feature	Actual Label Value	Supported Class Label	Contribution
RETURNDATACOPY	0	Legitimate	0.24
RETURN	0	Legitimate	0.21
SLT	1	Malicious	0.20
EQ	0	Malicious	0.18
LOG	1	Legitimate	0.16
OR	1	Legitimate	0.08
SSTORE	1	Legitimate	0.08
DELEGATECALL	0	Malicious	0.08
SUB	1	Malicious	0.03
SLOAD	1	Legitimate	0.02

TABLE VI: LIME Algorithm Based Justification for Classifier Prediction in Legitimate Contract

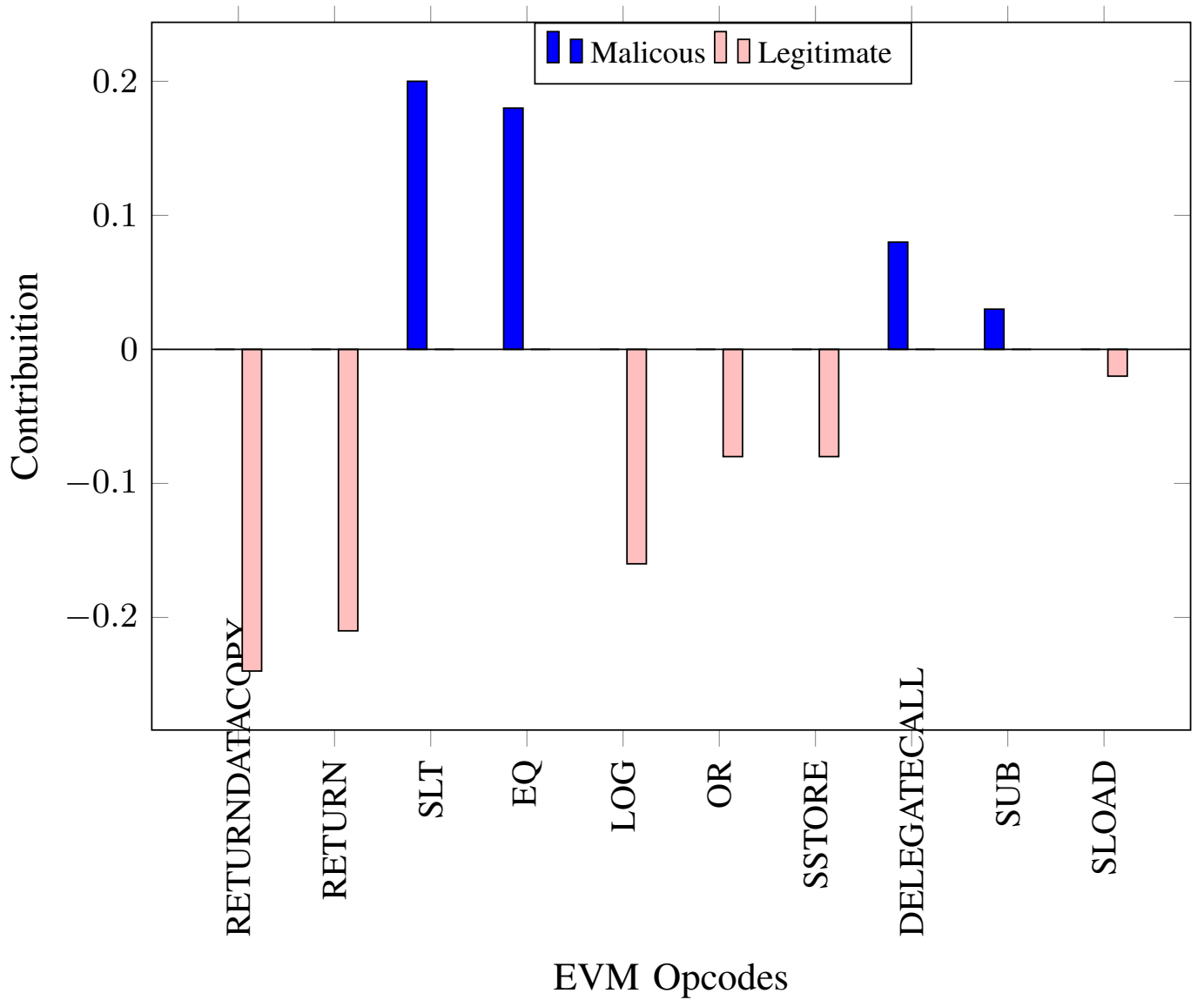


Fig. 3: Contributions of EVM Opcode Features in Legitimate Contract

contract verifications in coq, *ieee access* 7 (2019) 37770–37791.

- [13] B. Nassirzadeh, H. Sun, S. Banescu, V. Ganesh, Gas gauge: A security analysis tool for smart contract out-of-gas vulnerabilities, in: *The International Conference on Mathematical Research for Blockchain Economy*, Springer, 2022, pp. 143–167.
- [14] M. Barboni, A. Morichetta, A. Polini, Sumo: A mutation testing strategy for solidity smart contracts, in: *2021 IEEE/ACM International Conference on Automation of Software Test (AST)*, IEEE, 2021, pp. 50–59.
- [15] R. Gupta, M. M. Patel, A. Shukla, S. Tanwar, Deep learning-based malicious smart contract detection scheme for internet of things environment, *Computers & Electrical Engineering* 97 (2022) 107583.
- [16] S. Kotsiantis, D. Kanellopoulos, P. Pintelas, et al., Handling imbalanced datasets: A review, *GESTS international transactions on computer science and engineering* 30 (1) (2006) 25–36.
- [17] A. Gramegna, P. Giudici, Shap and lime: an evaluation of discriminative power in credit risk, *Frontiers in Artificial Intelligence* 4 (2021) 752558.
- [18] S. B. Kotsiantis, I. Zaharakis, P. Pintelas, et al., Supervised machine learning: A review of classification techniques, *Emerging artificial intelligence applications in computer engineering* 160 (1) (2007) 3–24.
- [19] A. R. Kharwar, D. V. Thakor, An ensemble approach for feature selection and classification in intrusion detection using extra-tree algorithm, *International Journal of Information Security and Privacy (IJISP)* 16 (1) (2022) 1–21.
- [20] O. Ibitoye, R. Abou-Khamis, A. Matrawy, M. O. Shafiq, The threat of adversarial attacks on machine learning in network security—a survey, *arXiv preprint arXiv:1911.02621* (2019).
- [21] L. Xu, M. Skoularidou, A. Cuesta-Infante, K. Veeramachaneni, Modeling tabular data using conditional gan, *Advances in neural information processing systems* 32 (2019).
- [22] H. de Bruijn, M. Warnier, M. Janssen, The perils and pitfalls of explainable ai: Strategies for explaining algorithmic decision-making, *Government information quarterly* 39 (2) (2022) 101666.