

# Eth2Vec: Learning Contract-Wide Code Representations for Vulnerability Detection on Ethereum Smart Contracts

Nami Ashizawa  
Osaka University

Naoto Yanai  
Osaka University

Jason Paul Cruz  
Osaka University

Shingo Okamura  
National Institute of Technology, Nara College

**Abstract**—Ethereum smart contracts are programs that run on the Ethereum blockchain, and many smart contract vulnerabilities have been discovered in the past decade. Many security analysis tools have been created to detect such vulnerabilities, but their performance decreases drastically when codes to be analyzed are being rewritten. In this paper, we propose Eth2Vec, a machine-learning-based static analysis tool for vulnerability detection, with robustness against code rewrites in smart contracts. Existing machine-learning-based static analysis tools for vulnerability detection need features, which analysts create manually, as inputs. In contrast, Eth2Vec automatically learns features of vulnerable Ethereum Virtual Machine (EVM) bytecodes with tacit knowledge through a neural network for natural language processing. Therefore, Eth2Vec can detect vulnerabilities in smart contracts by comparing the code similarity between target EVM bytecodes and the EVM bytecodes it already learned. We conducted experiments with existing open databases, such as Etherscan, and our results show that Eth2Vec outperforms the existing work in terms of well-known metrics, i.e., precision, recall, and F1-score. Moreover, Eth2Vec can detect vulnerabilities even in rewritten codes.

## I. INTRODUCTION

### A. Backgrounds

Ethereum [48] is the largest platform that provides an execution environment for smart contracts, and many distributed applications have been developed and deployed as smart contracts on the Ethereum blockchain. Ethereum smart contracts are programs that are stored on the Ethereum blockchain and are run by the Ethereum Virtual Machine (EVM) as EVM bytecodes<sup>1</sup>.

Given the transparent and decentralized nature of the Ethereum blockchain, the EVM bytecodes of smart contracts deployed on the Ethereum blockchain can be accessed and analyzed by anyone. Unfortunately, this also means that an adversary can abuse smart contracts [50] by analyzing their EVM bytecodes and looking for vulnerabilities. Consequently, attacks on vulnerable smart contracts can occur and possibly cause significant damage, especially when the attacked smart contracts handle assets. For example, the DAO attack is an infamous springboard attack where the attacker/s exploited a vulnerability in the DAO smart contract and stole more than 60 million USD worth of Ether, the cryptocurrency used in Ethereum.

According to literature [50], the security of smart contracts cannot be guaranteed because of the complexity (or lack of complexity) of the programming languages used for creating smart contracts, e.g., Solidity, which are relatively new languages, and the insufficient knowledge of programmers when creating smart contracts. To make matters worse, smart contracts that are successfully deployed on a blockchain cannot be modified, i.e., their source codes cannot be edited or deleted. In the past years, many attacks on deployed smart contracts have been reported [1], and thus making sure that the source codes of smart contracts are not vulnerable before they are deployed on a blockchain is desirable. To do this, many analysis tools for vulnerability detection in Ethereum smart contracts have been developed [7].

In this paper, we aim to develop a static analysis tool that can precisely identify various vulnerabilities in smart contract codes with high throughput by analyzing these codes. In static analysis, only a source code of a target to be analyzed is provided as input to identify if the code has vulnerabilities without executing the target itself. Therefore, ideally, static analysis can be used to prevent vulnerable codes from being deployed.

However, static analysis has two problems in general: (1) accuracy of its vulnerability detection is limited, and (2) its analysis time can be long. For instance, disassembly from EVM bytecodes [4], [33], [41], [49] does not have the capability to identify whether a program is vulnerable or not, and the early literature focuses on simply improving the readability of disassembled codes. In other words, disassembled codes often need to be analyzed manually, consequently increasing the number of false positives and false negatives. Moreover, symbolic execution [6], [26], [44], [47], which extracts control flow graphs (CFGs) from a target code, achieves high accuracy by automating the analysis, but the generation of CFGs needs to search all states such that the target code transits. Therefore, the analysis takes significant amounts of time [46].

A potential solution to the problems described above is machine learning. Static analysis based on machine learning infers whether a given code is vulnerable by learning features of codes. In doing so, the analysis also achieves a versatile analysis within a short time. However, static analysis based on machine learning has two inherent limitations: (1) code rewrites decrease analysis accuracy, and (2) appropriate features of smart contracts are indefinite. In the first limitation, for instance, CFGs with inlined functions are different from

<sup>1</sup>Hereafter, “Ethereum smart contract/s” and “smart contract/s” are used interchangeably but have the same meaning.

those of the original functions. Therefore, identifying the original function with the inlined function as the same codes is challenging. A pair of functions with the same semantics but different structures may generate different analysis results because the differences between structures of the codes strongly affect the analysis. In the second limitation, although features are manually extracted for machine learning, code samples and open knowledge about smart contracts are insufficient. Common knowledge about smart contract features has never been established [50]. Notably, kinds of features representing vulnerabilities in codes of smart contracts are unknown. Besides, features that are robust against differences in code structures are still not obvious.

## B. Contributions

In this paper, we propose *Eth2Vec*, a static analysis tool based on machine learning that identifies smart contract vulnerabilities by learning smart contract codes via their EVM bytecodes, assembly codes, and abstract syntax trees. *Eth2Vec* has high throughput, high accuracy, and robustness against code rewrites. *Eth2Vec* is an analysis tool based on a neural network for natural language processing, and it outputs the existence and kind of vulnerabilities in a target smart contract code only by taking the code as input. Using *Eth2Vec*, a user can analyze codes of smart contracts quickly even without expert knowledge on smart contract vulnerabilities. To achieve this, we also developed a parser for EVM bytecodes, including compilation of the Solidity language, which is a high-level language used for creating smart contracts. As a result, developers can analyze their smart contract codes directly even without expert knowledge of vulnerabilities and before deploying them onto the blockchain.

In terms of analysis of Ethereum smart contracts by using machine learning, a major contribution of this paper is the provision of a method that is robust against code rewrites. As described in the previous subsection, analysis via a typical machine learning algorithm [29] may output wrong results when codes to be analyzed are rewritten. Such situation happens because existing tools learn *only* patterns of code descriptions, i.e., the tools cannot learn the underlying features of the codes themselves. As another aspect of the limitation above, features that should be leveraged for analysis have never been established in existing works, to the best of our knowledge.

*Eth2Vec* overcomes the limitation described above by leveraging a neural network for natural language processing. While existing tools [29], [46] learn features that are given manually, *Eth2Vec* automatically learns features by entrusting the feature extraction phase to a neural network. In other words, using a neural network can isolate the feature extraction from the technical difficulty of analysis of Ethereum smart contracts. Indeed, a neural network for natural language processing has achieved high accuracy in processing of assembly codes [9], [51]. Nonetheless, *Eth2Vec* is novel for utilizing a neural network for the analysis of Ethereum smart contracts. To do this, we also developed a module that gives the EVM bytecodes to the neural network as inputs. Furthermore, we designed a learning methodology based on a neural network for natural language processing which takes the EVM bytecodes

through compilation of Solidity codes as inputs. *Eth2Vec* achieves robust analysis against code rewrites as well.

We conducted experiments to evaluate the performance of *Eth2Vec*. We used 5,000 files from Etherscan<sup>2</sup> as the dataset of contracts in *Eth2Vec*, and then executed the 10-fold cross-validation. The experimental results show that *Eth2Vec* can detect vulnerabilities within 1.2 seconds per contract with an average precision of 77.0%. Notably, reentrancy, whose severity is the highest among known vulnerabilities [43], can be detected with 86.6% precision. Our results also indicate that *Eth2Vec* outperforms the method by Momeni et al. [29], i.e., the use of support vector machine (SVM) based on their recommended features, as a naive method in terms of precision, recall, and F1-score of vulnerability detection. Besides, when we checked outputs by *Eth2Vec* in detail, we found examples of code clones with code rewrites and their vulnerabilities, which were not found by the SVM-based method in the current experiment.

We plan to release the source codes of *Eth2Vec* via GitHub for reproducibility and as reference for future works.

## II. PRELIMINARIES

In this section, we describe background knowledge to help readers understand our work.

### A. Ethereum Smart Contracts

In Ethereum, there are two kinds of accounts, namely, an externally owned account (EOA) and a contract account. EOAs have a private key that can be used to access the corresponding Ether or contracts. A contract account has smart contract code, which an EOA cannot have, and it does not have a private key. Instead, it is owned and controlled by the logic of its smart contract code. In Ethereum, a smart contract refers to an *immutable* computer program that is deployed on the blockchain and runs *deterministically* in the context of the EVM. The immutability property indicates that, similar to any data published on a general blockchain, smart contract codes can be considered as trustworthy, i.e., once deployed, they cannot be changed or deleted. The deterministic property indicates that the execution of the coded functions of smart contracts will produce the same result for anyone who runs them. Once deployed on the blockchain, a contract is self-enforcing and managed by the peers in the network, i.e., its functions are executed when the conditions in the contract are met. A smart contract is given an identity in terms of a contract address. Using this address, it can receive Ether and its functions can be executed. A contract is invoked when its contract address is the destination of a transaction, which is a signed message originating from an EOA, transmitted by the network, and recorded on the blockchain. Such transaction causes a contract to run in the EVM using the transaction (and transaction’s data) as input. The data indicate which specific function in the contract to run and what parameters to pass to that function. To incentivize peers to execute contract functions, Ethereum relies on *gas*, which is paid in Ether, to “fuel computations”. The amount of gas needed to execute a transaction is relative to the complexity of the computations, thus also preventing infinite loops.

<sup>2</sup><https://etherscan.io>

Smart contracts are typically written in a high-level language, such as Solidity [10]. The source code is then compiled to low-level bytecode that runs in the EVM. The EVM is a simple stack-based architecture. Its instruction set is kept minimal to avoid incorrect implementations that could cause consensus problems. The EVM is a global singleton, i.e., it operates like a global, single-instance computer that runs in all peers in the network. Each peer runs a local copy of the EVM to validate the execution of contract functions, and the processed transactions and smart contracts are recorded on the blockchain.

## B. Machine Learning

Machine learning consists of two algorithms, i.e., *training* and *inference*. The training algorithm takes data as input to learn their features and optimize parameters inside the model for an objective function. On the other hand, the inference algorithm takes unseen data as input and infers a similar set of features to that of the training data. When each data is unlabeled, a learning algorithm is called unsupervised learning. The most popular approach to machine learning in the recent years is deep learning, which is based on neural networks and can extract features in a black-box manner. In this paper, we aim to develop a model for learning vulnerable smart contracts to detect vulnerabilities in unlearned smart contracts.

## III. ANALYSIS OF SMART CONTRACTS VIA MACHINE LEARNING

This section describes the analysis target as the problem setting and its technical difficulty to be tackled in this paper.

### A. Analysis Target

In this paper, we focus on security analysis of the Solidity language as a target of static analysis of Ethereum smart contracts. In particular, we aim to identify the existence of vulnerabilities and classify the kinds of vulnerabilities in the codes to be analyzed. This means that, for example, a developer uses a tool to analyze the smart contracts he/she is developing in local. Such tool potentially needs to convince a developer that a smart contract being developed does not have any vulnerability even if the developer does not have sufficient knowledge about smart contract vulnerabilities. Therefore, we aim to develop a tool that can identify the existence of vulnerabilities (if there are any) in smart contracts even if only the codes of the smart contracts are given as input. Such a specification for analyzing smart contracts is preferable because standardized knowledge about Ethereum smart contracts is insufficient compared to general programming languages such as C and Java [50]. Meanwhile, strong obfuscation, i.e., the use of encryption, on Solidity codes is out of the scope of this paper because such an obfuscated code of smart contracts does not exist as far as we know.

Hereafter, we refer to codes written in Solidity as a contract and a file of codes consisting of more than a single contract as a contract file. We also call contracts to evaluate a vulnerability as test contracts and those to learn the vulnerability as training contracts. We call codes obtained from the compilation of the contracts as EVM bytecodes. The largest unit in each contract is a function, and a library function is also identical to a

function. Finally, “the blockchain” will be used to refer to the Ethereum blockchain unless otherwise specified. The problem setting in this paper is then formalized as follows:

**Problem Formulation:** We formalize our approach for analysis of smart contracts as follows. Each contract  $c_i \in \mathcal{C}$  includes vulnerabilities  $V_i = \{v_1^i, \dots, v_l^i\} \in \mathcal{V}^l$ , where  $\mathcal{C}$  denotes a set of contracts,  $\mathcal{V}$  denotes a set of vulnerabilities independent of each other, and  $l$  denotes any number. Given any integer  $n \in \mathbb{N}$ , a combination of a contract and vulnerabilities  $CV = \{(c_1, V_1), \dots, (c_n, V_n)\}$  and a test contract  $c_t \in \mathcal{C}$  are inputs of a model  $M$ . Let  $\{\epsilon_i^{c_t}\}_{i \in [1, d]} \subseteq \mathbb{R}^d$  denote the output of the model  $M$  which has  $d = |\mathcal{V}|$  elements, where  $|\mathcal{V}|$  denotes the size of  $\mathcal{V}$  and  $\epsilon_i^{c_t}$  denotes a probability about vulnerabilities in  $\mathcal{V}$ . Our goal is to develop a tool that optimizes  $M(CV, c_t) \rightarrow \{\epsilon_i^{c_t}\}$ .

### B. Technical Difficulty

Although several machine-learning-based static analysis tools for smart contracts have been proposed in literature [24], [25], [29], [40], [46], accuracy of their vulnerability detection is limited even on known vulnerabilities. In other words, accuracy significantly decreases when codes are rewritten from the original codes. Namely, on the problem formulation described above, the limitation described in Section I-A as our motivation is formalized as, even if  $c_t$  with a vulnerability  $V_i$  which is included in  $CV$  for any  $i \in [1, n]$ ,  $\epsilon_i^{c_t} \ll \epsilon_i^{c_i}$  holds for  $c_t \notin CV$ . We call such a situation non-robust.

The limitation related to the accuracy of vulnerability detection described above is caused by insufficient extraction of features of the existing tools. In general, a machine learning model needs features as inputs, which are manually extracted, and then learns the features explicitly. However, features for representing smart contracts to be analyzed are non-obvious because the history of Ethereum smart contracts is shorter than other general languages such as C and Java [50]. Intuitively, the limitation about the non-robustness is denoted as that required of a model  $M$  is unknown. Another reason for the insufficient extraction of features is the lack of code samples of smart contracts [50].

## IV. DESIGN OF ETH2VEC

In this section, we present Eth2Vec. We first describe the design concept to overcome the technical difficulty described in the previous section and then present the tool overview and its building blocks. Finally, we present the Eth2Vec model with its objective function.

### A. Design Concept

We aim to solve the technical difficulty by leveraging neural networks for natural language processing. Loosely speaking, neural networks handle feature extraction in a black-box manner, and thus the extraction of features can be isolated from the technical difficulty of the training. Therefore, Eth2Vec can analyze vulnerabilities of codes even if essential features of the vulnerabilities are unclear.

The term natural language processing described above means that the computation of the code similarity so that each word and paragraph are vectorized by inputting text data to

neural networks, e.g., Word2Vec. When the security of codes is analyzed, a model learns vulnerable codes and it can then identify the vulnerabilities via the similarity of codes to be analyzed with the learned codes.

To incorporate the natural language processing, Eth2Vec utilizes the *PV-DM* model [21] as neural networks to deal with EVM bytecodes. *PV-DM* model learns document representation based on tokens in the document. However, according to Ding et al. [9], a document is sequentially laid out, which is different from program codes. In particular, program codes can be represented as a graph and has a specific syntax. On the design of Eth2Vec, toward the analysis of EVM bytecodes leveraging the *PV-DM* model, we developed a new module named *EVM Extractor* to represent an abstract syntax tree of the codes.

### B. Tool Overview

Eth2Vec consists of two modules, i.e., *PV-DM* model [21] for neural networks to deal with paragraphs and *EVM Extractor* to create inputs of the *PV-DM* model from Solidity source codes.

The overview of Eth2Vec is shown in Figure 1. First, a *PV-DM* model is utilized as neural networks to deal with bytecodes. In particular, the *PV-DM* model executes unsupervised learning by taking JSON files generated from bytecodes as input and then computes the code similarity for each contract. To do this, we developed the *EVM Extractor* as a module to create inputs of the *PV-DM* model from EVM bytecodes because the *PV-DM* model cannot deal with EVM bytecodes initially. More specifically, the *EVM Extractor* analyzes EVM bytecodes syntactically and creates JSON files for instruction-level, block-level, function-level, and contract-level.

As a result, Eth2Vec takes EVM bytecodes to be analyzed as input, and then returns lists of code clones and their vulnerabilities for contract-level from a user’s standpoint. Meanwhile, vulnerabilities of contracts for training data are identified in advance by the use of existing tools [26], [43]. Vulnerabilities for test data are then evaluated by the code similarity with the vulnerable contracts. We show an output example of Eth2Vec in Appendix A.

### C. Building Blocks

Eth2Vec utilizes the *PV-DM* model [21] and *EVM Extractor* as building blocks.

1) *PV-DM model*: The *PV-DM* model is an extension of Word2Vec that treats text data for paragraph-level. Intuitively, it learns vector representations for each word and each paragraph. More concretely, given a text paragraph consisting of multiple sentences, a *PV-DM* model applies a sliding window over each sentence. The sliding window starts from the beginning of the sentence and moves forward a single word at each step. In doing so, the *PV-DM* model executes a multi-class prediction task such that it maps the current paragraph into a vector and each word in the context into a vector. More precisely, the model averages these vectors and infers the target word from the vocabulary via the softmax function. Formally, given a text that contains a list of paragraphs  $p \in T$ , each paragraph  $p$  contains a list of sentences  $s \in p$ , and each

sentence is a sequence of  $|s|$  words  $w_t \in s$ . Then, the *PV-DM* model maximizes the log probability as follows:

$$\sum_p \sum_s \sum_{t=k}^{|s|-k} \log \mathbf{P}(w_t | p, w_{t-k}, \dots, w_{t+k}). \quad (1)$$

*PV-DM* model is initially designed for text data that is sequentially laid out while it can also be for analysis on an assembly code by extending vector representation of language [9]. In this paper, we present a contract-wide representation learning model by leveraging the *PD-DM* model on the syntax of EVM bytecodes.

2) *EVM Extractor*: *EVM Extractor* is a module that syntactically analyzes EVM bytecodes for *PV-DM* model. In particular, *EVM Extractor* parses Solidity files as in the following hierarchical structure:

```

data: dictionary
├─ name: file name
├─ md5: md5 hash
├─ functions: list
│   └─ name: function name
│   └─ sea: start point
│   └─ see: end point
│   └─ id: function id
│   └─ call: list of callee functions
│   └─ blocks: list
│       └─ name: block name
│       └─ bytes: bytecodes
│       └─ sea: start point
│       └─ eea: end point
│       └─ id: block id
│       └─ call: list of callee blocks
│       └─ src: assembly instructions

```

The top-level, *data*, represents a contract file to be analyzed. The second level indicates file-dependent information. The third level indicates function-dependent information. The fourth level indicates a block consisting of multiple instructions. In each level, *sea* means the start address of the current function and that of the current block. Likewise, *see* and *eea* mean the start address of the next function and that of the next block, respectively. In contrast, *call* represents a callee function (or a callee block) from the current function (or block). Finally, at the bottom level, instructions of EVM bytecodes are stored in *src* together with their addresses. Meanwhile, a library function is treated as an individual contract.

### D. Design of Objective Function

Eth2Vec is a model that learns and infers codes of Ethereum smart contracts through unsupervised learning. Intuitively, it can be regarded as an extension of Word2Vec with a wider vector representation of codes given its ability to target Ethereum smart contracts.

The objective function of Eth2Vec is designed by extending that of the original *PV-DM* model described in Equation (1) along with the syntax shown in Section IV-C2 in the same argument in [9], which targets on an assembly language. Given a file of an Ethereum smart contract to be analyzed, Eth2Vec vectorizes the codes for instruction-level and then computes the

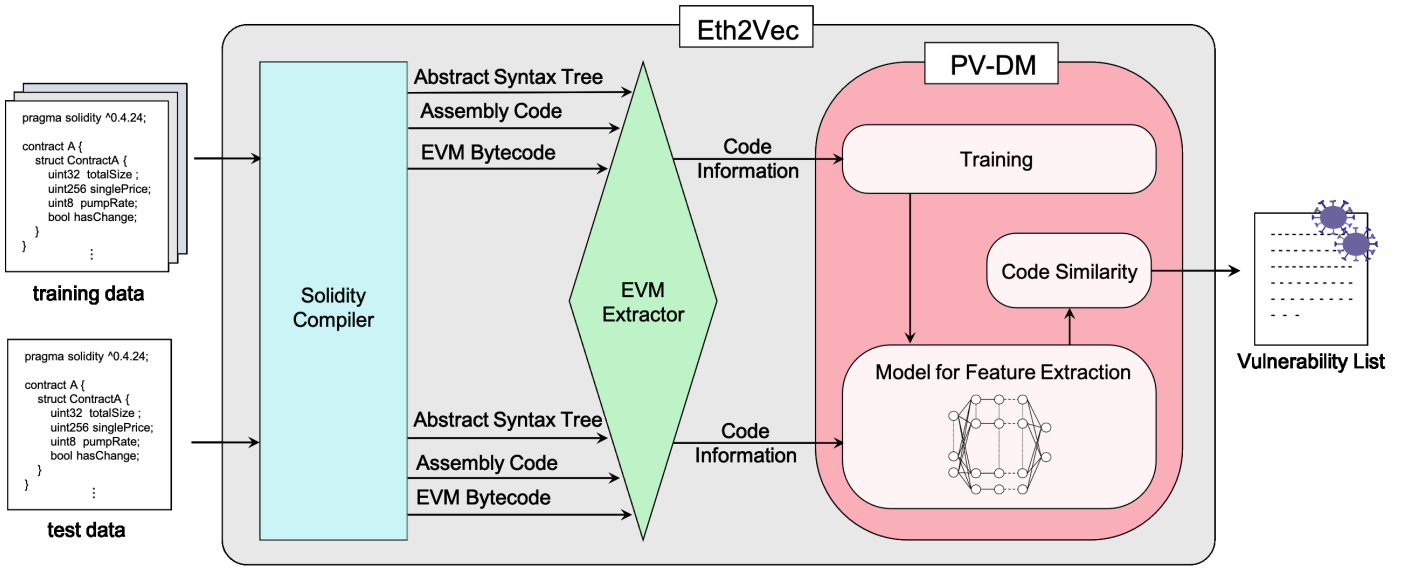


Fig. 1: Overview of Eth2Vec

code similarity. The process above is executed for block-level related to the instructions, function-level, and contract-level recursively, and therefore Eth2Vec can identify codes, whose distribution is similar to that of the training data, as a clone. Formally, an objective function is defined as follows:

$$\sum_{C_i \in \text{Dict}} \sum_{f_s} \sum_{seq_i} \sum_{in_j} \sum_{t_c} \log \mathbf{P}(t_c | C_i, in_{j-1}, in_{j+1}), \quad (2)$$

where we denote by  $\text{Dict}$  a contract file, by  $C_i$  a contract, by  $f_s$  a function, by  $seq_i$  a block consisting of plural instructions, by  $in_j$  each instruction, and by  $t_c$  a token with respect to the current instruction. Likewise, we denote by  $\mathcal{U}(C_i)$  plural functions, by  $\mathcal{S}(f_s)$  plural blocks, by  $\mathcal{I}(seq_i)$  a list of instructions, and by  $\mathcal{T}(in_j)$  a list of tokens.

Here, the first summation term, i.e.,  $\sum_{C_i \in \text{Dict}}$ , is included in the second level shown in the previous section, i.e., `file` name. Equation (2) is given contracts and their instructions, and then it maximizes the log probability for the current token  $t_c$ . Intuitively, the lexical meaning for each contract is computed through the current instruction and its neighbor instructions. Moreover, codes consisting of plural contracts can also be analyzed by extracting features for each contract.

We describe the objective function to represent the aforementioned intuition in detail below. Given a contract file  $Dict$ , a function  $f_s$  for each contract  $C_i$  is vectorized and then we denote by  $\vec{\theta}_{f_s}$  the vector representation of  $f_s$ . Furthermore, we denote by  $\mathcal{CT}(in)$  average of the vector representation of neighbor instructions of  $in$ . We then define a concatenation of the vector representation of the instruction itself and that of its operand as follows:

$$\mathcal{CT}(in) = \vec{v}_{\mathcal{P}(in)} \parallel \frac{1}{\mathcal{A}(in)} \sum_t \vec{v}_t, \quad (3)$$

where  $\mathcal{P}(in)$  is one operation with respect to  $in$ ,  $\mathcal{A}(in)$  is a list of operands with respect to  $in$ , and  $\parallel$  is a concatenation

of strings. In doing so, for any contract  $C_i$ , by averaging  $f_s$  and the  $j$ -th instruction  $in_j$  with  $\mathcal{CT}(in_{j-1})$  and  $\mathcal{CT}(in_{j+1})$ , a function  $\delta(in_j, f_s)$  to evaluate the joint memory of neighbor instructions in  $f_s$  is defined as follows:

$$\delta(in_j, f_s) = \frac{1}{3} \left( \vec{\theta}_{f_s} + \mathcal{CT}(in_{j-1}) + \mathcal{CT}(in_{j+1}) \right). \quad (4)$$

Therefore, the log probability in Equation (2) can be replaced with  $\mathbf{P}(t_c | C_i, in_{j-1}, in_{j+1}) = \mathbf{P}(t_c | \delta(in_j, f_s))$ . According to the literature [9], by letting  $X = (\vec{v}_{t_c})^T \times \delta(in_j, f_s)$  and utilizing a sigmoid function  $\sigma(X) = \frac{1}{1+e^{-X}}$ , the above computation can be approximated by the  $k$ -negative sampling [21], [28] as follows:

$$\begin{aligned} & \sum_{C_i} \sum_{f_s} \sum_{seq_i} \sum_{in_j} \sum_{t_c} J(\theta) \\ &= \sum_{C_i} \sum_{f_s} \sum_{seq_i} \sum_{in_j} \sum_{t_c} \log \mathbf{P}(t_c | \delta(in_j, f_s)), \end{aligned} \quad (5)$$

$$J(\theta) \approx \log(\sigma(X)) + \sum_{i=1}^k \mathbb{E}_{t_d \sim P_n(t_c)} (\llbracket t_d \neq t_c \rrbracket \log \sigma(X)), \quad (6)$$

where  $\llbracket t_d \neq t_c \rrbracket$  is an identity function which returns 1 if the expression inside the function is true or 0 otherwise. On the other hand,  $\mathbb{E}_{t_d \sim P_n(t_c)}$  is a sampling function that samples a token  $t_d$  in accordance with the noise distribution  $P_n(t_c)$  from  $t_c$ . We finally utilize Equation (6) as the objective function of Eth2Vec. Intuitively, the function maximizes the probability of a token  $t_c$  of the current instruction and decreases that of a token  $t_d$  of the other instructions. Then, we can compute the

gradient with respect to  $\theta_{f_s}$  through the following derivative:

$$\begin{aligned} \frac{\partial J(\theta)}{\partial \theta_{f_s}^{\rightarrow}} &= \frac{\partial J(\theta)}{\partial X} \frac{\partial X}{\partial \theta_{f_s}^{\rightarrow}} \\ &= \frac{\vec{v}_{t_c}}{3} (1 - \sigma(X)) \\ &+ \frac{\vec{v}_{t_c}}{3} \sum_i^k \mathbb{E}_{t_d \sim P_n(t_c)} (\llbracket t_d \neq t_c \rrbracket (1 - \sigma(X))) \end{aligned}$$

Moreover, we can also utilize a more intuitive gradient by approximating the aforementioned gradient computation as follows:

$$\frac{\partial J(\theta)}{\partial \theta_{f_s}^{\rightarrow}} \approx \frac{\vec{v}_{t_c}}{3} \sum_i^k \mathbb{E}_{t \sim P_n(t_c)} (\llbracket t = t_c \rrbracket - \sigma(X)). \quad (7)$$

Intuitively, Equation (7) approximates the original derivative described above in a manner that it moves a gradient to the positive if a token is identical to the current instruction or to the negative otherwise. Likewise, the gradient with respect to the vector representation  $\vec{v}_{t_c}$  of a token on the current instruction can be approximated. Moreover, for the vector representation  $\vec{v}_{\mathcal{P}(in)}$  of a token of some operation  $\mathcal{P}(in)$  and that  $\vec{v}_{\mathcal{A}(in)}$  of operands on the current instruction, their gradients can be approximated in a similar manner. These gradients are computed as follows although we omit their derivation:

$$\frac{\partial J(\theta)}{\partial \vec{v}_{t_c}} \approx (\llbracket t = t_c \rrbracket - \sigma(X)) \cdot \delta(in_j, f_s) \quad (8)$$

$$\frac{\partial J(\theta)}{\partial \vec{v}_{\mathcal{P}(j+1)}} \approx (\llbracket t = t_c \rrbracket - \sigma(X)) \cdot \frac{1}{\mathcal{A}(in)} \sum_t^{\mathcal{A}(in)} \vec{v}_t^{\rightarrow} \quad (9)$$

$$\frac{\partial J(\theta)}{\partial \vec{v}_{\mathcal{A}(in_{j+1})}} \approx (\llbracket t = t_c \rrbracket - \sigma(X)) \cdot (\vec{v}_{\mathcal{P}(in_{j+1})}^{\rightarrow}) \quad (10)$$

Similar equations hold for the previous instruction  $in_{j-1}$  by replacing  $in_{j+1}$  with  $in_{j-1}$  although we omit the details due to space limitation.

### E. Training and Inference

We now present the training algorithm and the inference algorithm with the gradients described in the previous section. First, the training algorithm is shown in Algorithm 1. The goal of the training is to optimize vectors for each instruction belonging to the given contract file as input. Instructions with similar meaning are mapped to a similar position in the vector space through the training. These resulting vectors by the algorithm can be used as features for vector representation of contract-level analysis. The algorithm is an unsupervised training algorithm, and thus the procedure does not require a ground-truth mapping between equivalent contracts.

The aforementioned features can be utilized directly in the inference algorithm shown in Algorithm 2. In particular, for an unseen contract  $C_t \notin \text{Dict}$ , the algorithm first initializes  $\theta_{f_t}^{\rightarrow}$ , which is associated with any function  $f_t$  belonging to  $C_t$ . Then, the algorithm follows the same procedure as the training algorithm. However, all  $\vec{v}_t^{\rightarrow}$ 's in the trained model are kept, and  $\theta_{f_t}^{\rightarrow}$ 's are updated following their errors. At the end of the inference,  $\theta_{f_t}^{\rightarrow}$  is output whereas the vectors for all  $C_i \in \text{Dict}$

---

**Algorithm 1** Training algorithm TRAIN of Eth2Vec for one epoch

---

**Require:** Dict

**Ensure:** vector representations for tokens of any  $in_j$  and  $\theta_{f_s}^{\rightarrow}$

```

1: for all  $C_i \in \text{Dict}$  do
2:   for all  $f_s \in \mathcal{U}(C_i)$  do
3:     for all  $seq_i \in \mathcal{S}(f_s)$  do
4:       for all  $in_j \in \mathcal{I}(seq_i)$  do
5:         vectorize  $f_s$  as  $\theta_{f_s}^{\rightarrow}$ 
6:         compute  $\mathcal{CT}(in_{j+1})$  by Eq. (3)
7:         compute  $\mathcal{CT}(in_{j-1})$  by Eq. (3)
8:         compute  $\delta(in_j, f_s)$  by Eq. (4)
9:         for all  $t_c \in \mathcal{T}(in_j)$  do
10:           compute a gradient for  $\theta_{f_s}^{\rightarrow}$  by Eq. (7)
11:           compute a gradient for  $\vec{v}_{t_c}^{\rightarrow}$  by Eq. (8)
12:           compute a gradient for  $\vec{v}_{\mathcal{P}(j+1)}^{\rightarrow}$  by Eq. (9)
13:           compute a gradient for  $\vec{v}_{\mathcal{A}(in_{j+1})}^{\rightarrow}$  by Eq. (10)
14:           compute a gradient for  $\vec{v}_{\mathcal{P}(j-1)}^{\rightarrow}$  by Eq. (9)
15:           compute a gradient for  $\vec{v}_{\mathcal{A}(in_{j-1})}^{\rightarrow}$  by Eq. (10)
16:         end for
17:         update vectors for tokens of  $in_j$ 
18:         update  $\theta_{f_s}^{\rightarrow}$ 
19:       end for
20:     end for
21:   end for
22: end for
```

---

**Algorithm 2** Inference algorithm Infer of Eth2Vec for a query

---

**Require:** contract  $C_t \notin \text{Dict}$

**Ensure:** vector representation  $\theta_{f_t}^{\rightarrow}$  for any  $f_t \in \mathcal{U}(C_t)$

```

1: Initialize vector representation  $\theta_{f_t}^{\rightarrow}$  of  $f_t \in \mathcal{U}(C_t)$ 
2: for all  $f_t \in \mathcal{U}(C_t)$  do
3:   for all  $seq_i \in \mathcal{S}(f_t)$  do
4:     for all  $in_j \in \mathcal{I}(seq_i)$  do
5:       compute  $\mathcal{CT}(in_{j+1})$  by Eq. (3)
6:       compute  $\mathcal{CT}(in_{j-1})$  by Eq. (3)
7:       compute  $\delta(in_j, f_s)$  by Eq. (4)
8:       for all  $t_c \in \mathcal{T}(in_j)$  do
9:         compute a gradient for  $\theta_{f_t}^{\rightarrow}$  by Eq. (7)
10:        update  $\theta_{f_t}^{\rightarrow}$ 
11:      end for
12:    end for
13:  end for
14: end for
```

---

remain the same except for  $\theta_{f_t}^{\rightarrow}$ . Finally, to search contracts for a match with the given contract  $C_t$ , the resultant vectors are compared with those of the training contracts using a typical statistical method, e.g., the cosine similarity.

### F. Implementation

We implement Eth2Vec by utilizing Kam1n0<sup>3</sup> and py-solc-x<sup>4</sup>. First, the main module, PV-DM model, is

<sup>3</sup>Kam1n0 version 2.0.0: <https://github.com/McGill-DMaS/Kam1n0-Community>

<sup>4</sup>py-solc-x: <https://pypi.org/project/py-solc-x/>

implemented by `Kamln0`, which is a server system [8] utilized for on binary analysis [9]. We mainly modified the source codes in `DisassemblyFactoryIDA.java` and `ExtractBinaryViaIDA.py`. In particular, `ExtractBinaryViaIDA.py` initially generates a JSON file extracted from IDA by disassembling binary codes, and then `DisassemblyFactoryIDA.java` takes the file to store the binary codes within `Kamln0`. However, IDA cannot use EVM bytecodes for implementing `Eth2Vec`. Therefore, we changed `ExtractBinaryViaIDA.py`: For instance, the code information is obtained by compiling a Solidity file with `py-solc-x` without IDA, and then its resultant assembly codes, abstract syntax tree (AST), and binary codes are extracted. We plan to release the source codes of the whole implementation publicly via GitHub.

## V. EXPERIMENTS

In this section, we describe the experiments we conducted to evaluate `Eth2Vec`. First, we describe the purpose of the experiments. Then, we discuss the datasets and the training methodologies used for evaluation. Finally, we show the experimental results.

### A. Purpose of Experiments

To evaluate the performance of `Eth2Vec`, we try to identify vulnerabilities in codes to be analyzed through the training with the known vulnerable contracts. To do this, we first check if `Eth2Vec` appropriately represents the relationship between codes to be analyzed and codes learned in the training phase. We evaluate clone detection of codes written in the Solidity language to confirm whether `Eth2Vec` can appropriately extract features of the codes. In doing so, we also evaluate semantic clones based on the lexical-semantic relationship to confirm the robustness of `Eth2Vec` against code rewrites.

Next, we check whether all the learned vulnerabilities are identified by training the codes and their code similarity with a given code. In doing so, we also evaluate consistency of the vulnerability detection with the results of the clone detection described above. In particular, we check if an output of the vulnerability detection is identical to the vulnerabilities of code identified as clones and if the clones are contained in the clone detection’s output.

We confirm that `Eth2Vec` can extract features precisely and thus can detect vulnerabilities through the experiments mentioned above. We also compare the performance of `Eth2Vec` with that of the existing work by Momeni et al. [29], which extracts features manually, as a baseline.

### B. Experimental Setting

As mentioned above, experiments are conducted in two stages, i.e., clone detection and vulnerability detection of smart contracts. We describe the setting for each experiment below. We first describe datasets and the baseline in detail.

1) *Dataset*: We collect 5,000 contract files from Etherscan<sup>5</sup>, which is an open database of smart contracts, as a dataset utilized in the experiments. These 5,000 files are also

identical to files utilized in a recent work [14]. We then utilize only the files that can be compiled by solidity version -0.4.11 as a compiler. The dataset contains 95,152 contracts and 1,193,868 blocks.

2) *Baseline*: We compare the performance of `Eth2Vec` with that of the scheme based on support vector machine (SVM) by Momeni et al. [29] as a baseline. Although there are several versatile results based on machine learning [14], [25], [46], which can detect various vulnerabilities, their source codes are unpublished or we were unable to build their source codes in our environment. We thus adopt the SVM-based method by Momeni et al. which we were able to reproduce from scratch.

Momeni et al. extracted 16 features from an abstract syntax tree (AST) of an Ethereum smart contract. We utilize 15 of these features excluding hexadecimal addresses because hexadecimal addresses cannot be obtained from source codes, i.e., without deployment. Other features are described in Appendix B.

3) *Clone Detection*: We check if `Eth2Vec` can precisely identify clones of test contracts as input through training contracts. In particular, on the 10-fold cross-validation, 500 test contracts are randomly chosen, and the remaining 4,500 contracts are utilized as the training contracts. This process is iterated ten times. Meanwhile, a threshold of the code similarity to detect clones is 0.8.

In the setting mentioned above, precision is utilized as an evaluation metric. For clones for each function, i.e., function-level clones, which are output by a threshold more than 0.8, we check if a function in a contract, which has the highest similarity, is indeed a clone of the given input. To compare the performance of `Eth2Vec`, we also evaluate the SVM-based method described above by manually labeling each function. Meanwhile, clones in this experiment correspond to the type-I to type-IV clones in literature [36], and semantic clones are identical to the type-IV clones.

4) *Vulnerability Detection*: We check if `Eth2Vec` can precisely detect vulnerabilities of test contracts as input through learning the training contracts. In particular, we check whether vulnerabilities output by `Eth2Vec` with a threshold 0.8 about the test contracts are identical to true vulnerabilities of the contracts.

In this experiment, we adopt the 10-fold cross-validation similar to that in the clone detection experiment. We also confirm the vulnerabilities of the test contracts as ground truth of the experiment by utilizing Oyente [26] and SmartCheck [43]. The vulnerabilities are listed in Appendix C due to space limitation. In this experiment, we adopt well-known metrics, i.e., accuracy, precision, recall, and F1-score.

### C. Results on Clone Detection

The results of the clone detection by `Eth2Vec` are shown in Table I. According to the table, `Eth2Vec` can implicitly extract more features than the SVM-based method [29]. For instance, the values of the standard deviation by `Eth2Vec` are significantly smaller than those by the SVM-based method. Intuitively, this means that feature extraction by `Eth2Vec` is more stable, and therefore `Eth2Vec` is able to represent features on various codes. When the clone detection performance is

<sup>5</sup><https://etherscan.io/>



TABLE I: Precision of Clone Detection of Eth2Vec. This table shows the average and standard deviation of 10 executions of precision measurement on 10-fold cross-validation. Furthermore, the row of SVM w/o few clones represents result of the SVM-based setting computed by removing few clones. The numbers are truncated to one decimal place.

	Eth2Vec	SVM [29]	SVM w/o few clones
Average	74.9%	34.6%	42.7%
Standard Deviation	0.9	34.6	43.6

```

1 function _transfer(address _from, address _to, uint
  _value) internal {
2   require(_to != 0x0);
3   require(balanceOf[_from] >= _value);
4   require(balanceOf[_to] + _value > balanceOf[_to
  ]);
5   uint previousBalances = balanceOf[_from] +
  balanceOf[_to];
6   balanceOf[_from] -= _value;
7   balanceOf[_to] += _value;
8   Transfer(_from, _to, _value);
9   assert(balanceOf[_from] + balanceOf[_to] ==
  previousBalances);
10 }

```

Listing 1: Original function of Transfer

```

1 function _transfer(address _from, address _to, uint
  _value) internal {
2   require(_to != 0x0);
3   require(balanceOf[_from] >= _value);
4   require(balanceOf[_to] + _value > balanceOf[_to
  ]);
5   balanceOf[_from] -= _value;
6   balanceOf[_to] += _value;
7   Transfer(_from, _to, _value);
8 }

```

Listing 2: Clone of the original function of Transfer found by Eth2Vec

compared to EClone [25] as a reference, the detection rate of EClone with a threshold of 0.4 was 58.2%. Although we did not implement EClone by ourselves as described in Section V-B2, we believe that Eth2Vec can detect clones of Ethereum smart contracts better than EClone because of its more sophisticated feature extraction.

When we checked an output of Eth2Vec in detail, it contained semantic clones. We show a concrete example in Listing 1 and Listing 2. The output of Eth2Vec, i.e., Listing 2, is a clone of the original function shown in Listing 1 excluding lines 5 and 9. Listing 2 is identical to a type-III clone of Listing 1 according to the definitions in [36]. This means that Eth2Vec can detect clones precisely, even those of rewritten codes.

Meanwhile, the precision of the SVM-based method is small because the test datasets contain few clones in several executions. The SVM-based method infers negatives for the given test data under such situation. Therefore, the precision, the recall, and the F1-score become 0, resulting in a low average and a more extensive standard deviation. Nonetheless, when we re-compute the precision by removing the results with

the few clones among ground-truth, the precision is updated to 42.7%.

#### D. Results on Vulnerability Detection

The results of the vulnerability detection by Eth2Vec are shown in Table II.

First, the standard deviations of Eth2Vec are lower than those of the SVM-based method. This means that Eth2Vec diminishes the effect of difference between features of vulnerabilities because Eth2Vec can detect the vulnerabilities more stably than the SVM-based method without Eth2Vec can precisely extract features and thus, it can take into account various vulnerabilities.

Second, we confirm that Eth2Vec is also robust against the dispersion of data. In particular, the appearance of vulnerabilities with high severity, i.e., reentrancy and time dependency, is low in the current dataset. Consequently, the number of the corresponding labels is less, and the recall of the SVM-based method became low. By contrast, the recall of Eth2Vec is more stable than that of the SVM-based method even against vulnerabilities, whose frequency of appearance is low, in comparison with the other vulnerabilities.

Moreover, when we checked an output of Eth2Vec in detail, we found a concrete example, as shown in Listing 3. The code is output as a clone of Listing 1 and the code similarity of Listing 1 with Listing 3 is lower than that with Listing 2. As difference between these codes, Listing 1 and Listing 2 do not have vulnerabilities while Listing 3 contains the integer overflow vulnerability. Specifically, line 4 in Listing 3 is different from Listing 1 with respect to an expression statement, and the statement in Listing 3 is identical to a typical form of integer overflow<sup>6</sup>.

The aforementioned example indicates that the vulnerability detection of Eth2Vec is robust against code rewrites. Although Listing 1 and Listing 3 seem to be more similar than Listing 1 and Listing 2, Listing 3 is different from a “precise” clone of Listing 1 because Listing 3 contains the integer overflow vulnerability, which is not included in Listing 1. Therefore, Eth2Vec returns a higher code similarity for Listing 2 and Listing 1 than for Listing 3 and Listing 1. We thus confirm that Eth2Vec has robust vulnerability detection.

#### E. Throughput for Inference

The throughput of the inference by Eth2Vec is shown in Figure 2. The processing time for vulnerability detection depends on the number of functions and that for displaying the detection results, i.e., summarizing by Kam1n0, depends on the number of detected clones. According to the current measurement, the detection throughput of Eth2Vec is faster than that of the SVM-based method except for the process of EVM Extractor. Eth2Vec can analyze within about 0.371 seconds per contract in comparison with about 0.397 seconds per contract by the SVM-based method. Meanwhile, the processing time for displaying the detection results is longer than the detection itself. A long analysis needed 5.6 seconds for detection, 0.015 seconds for saving, and 3.5 seconds for summarizing, where

<sup>6</sup><https://github.com/ConsenSys/mythril/wiki/Integer-Overflow>



TABLE II: Vulnerability Detection Results on Eth2Vec: In the following table, we measured well-known metrics for both Eth2Vec and the SVM-based method [29]. The values of the averages and standard deviations are computed by the values of each row. The numbers are truncated to one decimal place.

Vulnerability	Severity	Eth2Vec			SVM [29]		
		Precision [%]	Recall [%]	F1-Score [%]	Precision [%]	Recall [%]	F1-Score [%]
Reentrancy	3	86.6	54.8	61.5	30.0	7.8	12.3
Time Dependency	2	75.2	17.0	27.3	55.0	2.8	5.3
ERC-20 Transfer	1	95.6	58.4	72.4	89.0	95.3	92.0
Gas Consumption	1	48.0	29.0	32.4	10.0	3.1	4.7
Implicit Visibility	1	68.9	82.0	74.8	71.5	77.5	73.8
Integer Overflow	1	89.9	57.6	70.1	84.9	73.1	78.3
Integer Underflow	1	74.6	56.0	63.7	75.1	39.2	50.0
Average		77.0	50.7	57.5	59.3	42.7	45.2
Standard Deviation		14.7	19.7	18.0	27.3	36.4	34.7

```

1 function _transfer(address _from, address _to, uint
  _value) internal {
2   require(_to != 0x0);
3   require(balanceOf[_from] >= _value);
4   require(balanceOf[_to] + _value >= balanceOf[_to
  ]);
5   uint previousBalances = balanceOf[_from] +
  balanceOf[_to];
6   balanceOf[_from] -= _value;
7   balanceOf[_to] += _value;
8   emit Transfer(_from, _to, _value);
9   assert(balanceOf[_from] + balanceOf[_to] ==
  previousBalances);
10 }

```

Listing 3: Vulnerability example by Eth2Vec to the original function

the number of contracts is nine and the number of functions is 94.

### F. Limitations

In this section, we describe several limitations of Eth2Vec, which will be improved in future works. First, the current construction is in the unsupervised setting, and the supervised setting is not implemented. According to Hill et al. [19], the supervised setting seems to be more suitable for the classification of natural language processing than the unsupervised setting. Intuitively, evaluating a kind of vulnerability contained in codes is a classification problem about vulnerabilities. Namely, the performance of Eth2Vec can be improved by utilizing the supervised setting.

Second, the current work only detected vulnerabilities of contracts included in the training dataset utilized in our experiment. There are potentially many vulnerable contracts that have already been deployed, and thus it would be ideal if all deployed contracts are analyzed to determine potential vulnerabilities.

Finally, the current construction does not support inter-contract analysis where multiple contracts are interconnected. For instance, Rodler et al. [35] presented a new kind of vulnerability by utilizing CALL and CREATE instructions, which requires inter-contract analysis [6], [47]. Thus, the attacks by Rodler et al. are out of the scope of Eth2Vec currently. Our future work will aim to extend the current construction of Eth2Vec and overcome the limitations described above.

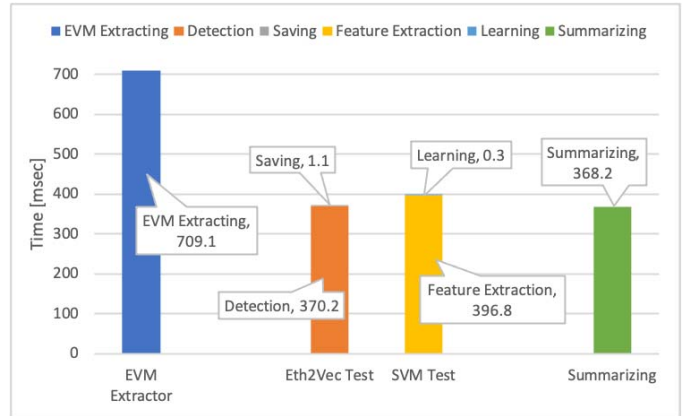


Fig. 2: Throughput of Eth2Vec: We measured throughput of Eth2Vec for each process. The item of Eth2Vec Test is identical to the vulnerability detection, which contains the detection itself and saving detection results. The item of Summarizing represents the processing time for the display process of the saved data. The entire process until a user obtains an output is a summation of EVM Extractor, Eth2Vec Test, and Summarizing. On this measurement, we randomly picked up 20 contract files, where a threshold is 0.8 and five clones, i.e., candidates of vulnerabilities, at maximum for each function are output.

## VI. RELATED WORKS

In this section, we describe related works in terms of security analysis by machine learning.

### A. Security Analysis of Ethereum by Machine Learning

As machine learning-based analysis, ContractWard [46] and the tool by Momeni et al. [29] are based on support vector machine (SVM) and random forest, whose features are extracted manually by an analyst. Hence, identifiable vulnerabilities are limited due to insufficient feature extraction as described in Section III. As research on automated feature extraction, VulDeeSmartContract [34] has been proposed by combining Word2Vec with long-short term memory (LSTM), but it specializes only in reentrancy vulnerability. Eth2Vec can be considered as a tool dealing with more versatile vulnerabilities in a similar approach.

Next, EClone [24] is a tool that detects code clones of Ethereum smart contracts through the computation of code similarity. Although EClone does not utilize neural networks, a symbolic analysis tool suitable for vector computations of codes has also been proposed. The detectability of Eth2Vec can be improved potentially by combining it with EClone.

In the neural network-based approach, there is ILF [17] which automatically generates an input of fuzzing test via neural networks. Loosely speaking, ILF learns codes by extracting features via symbolic execution. Since ILF learns outputs by symbolic execution, it is an entirely different tool from Eth2Vec, which learns codes themselves.

Finally, SmartEmbed [14] is a tool that identifies bugs of smart contract codes by leveraging FastText, which represents the codes as vectors. SmartEmbed is the closest work to Eth2Vec. However, SmartEmbed did not discuss an objective function and training algorithm explicitly and just vectorizes codes in accordance with Word2Vec and FastText. Consequently, there are several code samples which are unidentified as clones, as shown in the paper [14]<sup>7</sup>. In contrast, Eth2Vec can precisely extract features of codes with a vast vocabulary and thus can detect vulnerabilities even after codes are rewritten. Nevertheless, SmartEmbed is an elegant work that also discussed code repairing as well as versatile bugs. Interested readers are advised to read the SmartEmbed paper.

### B. Other Analysis Tools for Ethereum

To the best of our knowledge, symbolic execution [26] is the principal approach for analysis of Ethereum smart contracts. Since symbolic execution deals with unknown variables as symbolic variables, it is potentially suitable for analysis of smart contracts, which utilizes information outside codes [6], i.e., blockchain. Hence, many tools have been proposed so far [5], [12], [23], [30]–[32], [44]. The primary motivation of recent works aims to extend analysis areas, e.g., inter-contract analysis and contract creation. Although symbolic execution consumes a longer time in comparison with a machine learning-based approach as described in Section I-A, the vulnerability detection of Eth2Vec may be improved by combining it with such tools.

Another approach for smart contract analysis is formal verification [3] which deduces whether a program satisfies a specification via predicate logic. In general, formal verification can provide precise analysis by representing the security in a mathematical way. However, the verification itself is a rigid and challenging work. Many works just formalized a specification of Ethereum and did not achieve the analysis of vulnerabilities [15], [16], [16], [18], [20], [45]. As the latest work, Grishchenko et al. presented a formal verification tool named eThor [37] that can analyze vulnerabilities in a versatile way. We consider eThor as the most elegant tool in formal verification, and thus readers interested in formal verification are suggested to read the eThor paper.

Finally, there are several tools on dynamic analysis which execute codes themselves [13], [35]. However, using dynamic analysis, an analyst needs to implement and execute attack

patterns by him-/herself to prevent vulnerabilities in advance. Although a universal attack pattern that can capture various attacks was proposed [13], an analyst still needs to have knowledge about attacks on smart contracts. Therefore, analysts with the ability of utilizing a dynamic analysis are limited, and thus static analysis is more reasonable.

### C. Security Applications of Natural Language Processing

We briefly describe natural language processing applications to the binary analysis and cybersecurity defense as further related works. After Shin et al. [39] proposed binary analysis based on neural networks, many works for binary analysis based on natural language processing, such as Word2Vec, have been proposed in recent years [9], [11], [27], [51]. These works are based on neural networks. The closest work to Eth2Vec is Asm2Vec [9], which focuses on assembly language. Eth2Vec can be seen as an extension of the implementation of Asm2Vec. The source code of Asm2Vec is publicly available and readers are advised to read that paper in conjunction with our work. We also note that the same construction can be implemented with other tools described above.

Likewise, natural language processing is another attractive area in cybersecurity. Walk2Friends, which performs a social relation inference attack [2], and DarkEmbed [42], which learns low dimensional distributed representations of darkweb/deepweb discussions, are examples of natural language processing. Other examples are Log2Vec [22] and Attack2Vec [38], which learn more generalized cybersecurity information. These tools apply natural language processing to cybersecurity areas, while Eth2Vec is an application to Ethereum smart contracts.

## VII. CONCLUSION

In this paper, we proposed Eth2Vec, a static analysis tool based on machine learning that detects vulnerabilities of Ethereum smart contracts. The most striking property of Eth2Vec is the automated feature extraction for each contract by leveraging neural networks for natural language processing. Consequently, by extracting features implicitly and incorporating lexical semantics between contracts, the vulnerabilities can be detected with 77.0% precision even after the codes are rewritten. Moreover, reentrancy which is one of the most important vulnerabilities can be detected with 86.6% precision. We also demonstrated that Eth2Vec outperforms the SVM-based method by Momeni et al. [29] in terms of precision, recall, and F1-score. We are preparing the release of the implementation of Eth2Vec via GitHub as well.

We are in the process of improving detection performance, i.e., precision, recall, and F1-score, for vulnerabilities and their underlying code clones. In particular, we intend to realize inter-contract analysis whereby multiple contracts affect each other. The performance of Eth2Vec will be improved significantly by introducing a function to obtain such information. Further studies, which take inter-contract analysis into account, will thus need to be undertaken.

## REFERENCES

- [1] N. Atzei, M. Bartoletti, and T. Cimoli. A survey of attacks on ethereum smart contracts (sok). In *Proc. of POST 2017*, volume 10204 of *LNCS*,

<sup>7</sup>Although the authors of [14] describe that the samples are not clones, these sample are type-IV clones following the definition in [36].

- pages 164–186. Springer, 2017.
- [2] M. Backes, M. Humbert, J. Pang, and Y. Zhang. Walk2friends: Inferring social links from mobility profiles. In *Proc. of CCS 2017*, page 1943–1957. ACM, 2017.
  - [3] K. Bhargavan, A. Delignat-Lavaud, C. Fournet, A. Gollamudi, G. Gonthier, N. Kobeissi, N. Kulatova, A. Rastogi, T. Sibut-Pinote, N. Swamy, et al. Formal verification of smart contracts: Short paper. In *Proc. of PLAS 2016*, pages 91–96. ACM, 2016.
  - [4] L. Brent, A. Jurisevic, M. Kong, E. Liu, F. Gauthier, V. Gramoli, R. Holz, and B. Scholz. Vandal: A scalable security analysis framework for smart contracts. *arXiv preprint arXiv:1809.03981*, 2018.
  - [5] T. Chen, X. Li, X. Luo, and X. Zhang. Under-optimized smart contracts devour your money. In *Proc. of SANER 2017*, pages 442–446. IEEE, 2017.
  - [6] Y. Chinen, N. Yanai, J. P. Cruz, and S. Okamura. Hunting for re-entrancy attacks in ethereum smart contracts via static analysis. *arXiv preprint arXiv:2007.01029*, 2020.
  - [7] M. Di Angelo and G. Salzer. A survey of tools for analyzing ethereum smart contracts. In *Proc. of DAPPCON 2019*, pages 69–78. IEEE, 2019.
  - [8] S. H. Ding, B. C. Fung, and P. Charland. KamIn0: Mapreduce-based assembly clone search for reverse engineering. In *Proc. of KDD 2016*, page 461–470. ACM, 2016.
  - [9] S. H. Ding, B. C. Fung, and P. Charland. Asm2vec: Boosting static representation robustness for binary clone search against code obfuscation and compiler optimization. In *Proc. of IEEE S&P 2019*, pages 472–489. IEEE, 2019.
  - [10] S. C. S. . documentation. <https://solidity.readthedocs.io/en/v0.5.11/security-considerations.html>.
  - [11] Y. Duan, X. Li, J. Wang, and H. Yin. Deepbindiff: Learning program-wide code representations for binary diffing. In *Proc. of NDSS 2020*. Internet Society, 2020.
  - [12] J. Feist, G. Grieco, and A. Groce. Slither: A static analysis framework for smart contracts. In *Proc. of WETSEB 2019*, pages 8–15. IEEE, 2019.
  - [13] C. Ferreira Torres, M. Steichen, R. Norvill, B. Fiz Pontiveros, and H. Jonker. Aegis: Shielding vulnerable smart contracts against attacks. In *Proc. of AsiaCCS 2020*, page 584–597. ACM, 2020.
  - [14] Z. Gao, L. Jiang, X. Xia, D. Lo, and J. Grundy. Checking smart contracts with structural code embedding. *IEEE Transactions on Software Engineering*, pages 1–1, 2020.
  - [15] I. Grishchenko, M. Maffei, and C. Schneidewind. Foundations and tools for the static analysis of ethereum smart contracts. In *Proc. of CAV 2018*, volume 10981 of *LNCS*, pages 51–78. Springer, 2018.
  - [16] I. Grishchenko, M. Maffei, and C. Schneidewind. A semantic framework for the security analysis of ethereum smart contracts. In *Proc. of POST 2018*, volume 10804 of *LNCS*, pages 243–269. Springer, 2018.
  - [17] J. He, M. Balunoviundefined, N. Ambroladze, P. Tsankov, and M. Vechev. Learning to fuzz from symbolic execution with application to smart contracts. In *Proc. of CCS 2019*, page 531–548. ACM, 2019.
  - [18] E. Hildenbrandt, M. Saxena, N. Rodrigues, X. Zhu, P. Daian, D. Guth, B. Moore, D. Park, Y. Zhang, A. Stefanescu, et al. Kevm: A complete formal semantics of the ethereum virtual machine. In *Proc. of CSF 2018*, pages 204–217. IEEE, 2018.
  - [19] F. Hill, K. Cho, and A. Korhonen. Learning distributed representations of sentences from unlabelled data. In *Proc. of NAACL HLT 2016*, pages 1367–1377. ACL, 2016.
  - [20] S. Kalra, S. Goel, M. Dhawan, and S. Sharma. Zeus: Analyzing safety of smart contracts. In *Proc. of NDSS 2018*. Internet Society, 2018.
  - [21] Q. Le and T. Mikolov. Distributed representations of sentences and documents. In *Proc. of ICML 2014*, pages 1188–1196, 2014.
  - [22] F. Liu, Y. Wen, D. Zhang, X. Jiang, X. Xing, and D. Meng. Log2vec: A heterogeneous graph embedding based approach for detecting cyber threats within enterprise. In *Proc. of CCS 2019*. ACM, 2019.
  - [23] H. Liu, C. Liu, W. Zhao, Y. Jiang, and J. Sun. S-gram: towards semantic-aware security auditing for ethereum smart contracts. In *Proc. of ASE 2018*, pages 814–819. ACM, 2018.
  - [24] H. Liu, Z. Yang, Y. Jiang, W. Zhao, and J. Sun. Enabling clone detection for ethereum via smart contract birthmarks. In *Proc. of ICPC 2019*, pages 105–115. IEEE, 2019.
  - [25] H. Liu, Z. Yang, C. Liu, Y. Jiang, W. Zhao, and J. Sun. Eclone: Detect semantic clones in ethereum via symbolic transaction sketch. In *Proc. of ESEC/FSE 2018*, page 900–903. ACM, 2018.
  - [26] L. Luu, D.-H. Chu, H. Olickel, P. Saxena, and A. Hobor. Making smart contracts smarter. In *Proc. of CCS 2016*, pages 254–269. ACM, 2016.
  - [27] L. Massarelli, G. A. Di Luna, F. Petroni, R. Baldoni, and L. Querzoni. Safe: Self-attentive function embeddings for binary similarity. In R. Perdisci, C. Maurice, G. Giacinto, and M. Almgren, editors, *Proc. of DIMVA 2019*, volume 11543 of *LNCS*, pages 309–329. Springer, 2019.
  - [28] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean. Distributed representations of words and phrases and their compositionality. In *Proc. of NIPS 2013*, pages 3111–3119, 2013.
  - [29] P. Momeni, Y. Wang, and R. Samavi. Machine learning model for smart contracts security analysis. In *Proc. of PST 2019*, pages 1–6. IEEE, 2019.
  - [30] M. Mossberg, F. Manzano, E. Hennenfent, A. Groce, G. Grieco, J. Feist, T. Brunson, and A. Dinaburg. Manticore: A user-friendly symbolic execution framework for binaries and smart contracts. *arXiv preprint arXiv:1907.03890*, 2019.
  - [31] B. Mueller. Smashing smart contracts. In *Proc. of HITBSECCONF 2018*, 2018.
  - [32] I. Nikolić, A. Kolluri, I. Sergey, P. Saxena, and A. Hobor. Finding the greedy, prodigal, and suicidal contracts at scale. In *Proc. of ACSAC 2018*, pages 653–663. ACM, 2018.
  - [33] R. Norvill, B. B. F. Pontiveros, R. State, and A. Cullen. Visual emulation for ethereum’s virtual machine. In *Proc. of NOMS 2018*, pages 1–4. IEEE, 2018.
  - [34] P. Qian, Z. Liu, Q. He, R. Zimmermann, and X. Wang. Towards automated reentrancy detection for smart contracts based on sequential models. *IEEE Access*, 8:19685–19695, 2020.
  - [35] M. Rodler, W. Li, G. O. Karame, and L. Davi. Sereum: Protecting existing smart contracts against re-entrancy attacks. In *Proc. of NDSS 2019*. Internet Society, 2019.
  - [36] C. K. Roy and J. R. Cordy. A survey on software clone detection research. *School of Computing TR 2007-541*, Queen’s University, 541(115):64–68, 2007.
  - [37] C. Schneidewind, I. Grishchenko, M. Scherer, and M. Maffei. Ethor: Practical and provably sound static analysis of ethereum smart contracts. In *Proc. of CCS 2020*, page 621–640. ACM, 2020.
  - [38] Y. Shen and G. Stringhini. Attack2vec: Leveraging temporal word embeddings to understand the evolution of cyberattacks. In *Proc. of USENIX Security 2019*, pages 905–921. USENIX Association, 2019.
  - [39] E. C. R. Shin, D. Song, and R. Moazzezi. Recognizing functions in binaries with neural networks. In *Proc. of USENIX Security 2015*, pages 611–626. USENIX Association, 2015.
  - [40] J. Song, H. He, Z. Lv, C. Su, G. Xu, and W. Wang. An efficient vulnerability detection model for ethereum smart contracts. In *Proc. of NSS 2019*, volume 11928 of *LNCS*, pages 433–442. Springer, 2019.
  - [41] M. Suiche. Porosity: A decompiler for blockchain-based smart contracts bytecode. *Proc. of DEFCON 2017*, 25:11, 2017.
  - [42] N. Tavabi, P. Goyal, M. Almukaynizi, P. Shakarian, and K. Lerman. Darkembed: Exploit prediction with neural language models. In *Proc. of AAAI 2018*, pages 7849–7854. AAAI Press, 2018.
  - [43] S. Tikhomirov, E. Voskresenskaya, I. Ivanitskiy, R. Takhaviev, E. Marchenko, and Y. Alexandrov. Smartcheck: Static analysis of ethereum smart contracts. In *Proc. of WETSEB 2018*, pages 9–16. ACM, 2018.
  - [44] C. F. Torres, J. Schütte, et al. Osiris: Hunting for integer bugs in ethereum smart contracts. In *Proc. of ACSAC 2018*, pages 664–676. ACM, 2018.
  - [45] P. Tsankov, A. Dan, D. Drachler-Cohen, A. Gervais, F. Buenzli, and M. Vechev. Securify: Practical security analysis of smart contracts. In *Proc. of CCS 2018*, pages 67–82. ACM, 2018.
  - [46] W. Wang, J. Song, G. Xu, Y. Li, H. Wang, and C. Su. Contractward: Automated vulnerability detection models for ethereum smart contracts. *IEEE Transactions on Network Science and Engineering*, pages 1–1 (Early Access), 2020.
  - [47] K. Weiss and J. Schütte. Annotary: A Concolic Execution System for

Developing Secure Smart Contracts. In *Proc. of ESORICS 2019*, volume 11735 of *LNCS*, pages 747–766. Springer, 2019.

- [48] G. Wood. Ethereum: A secure decentralised generalised transaction ledger byzantium version. <https://ethereum.github.io/yellowpaper/paper.pdf>.
- [49] Y. Zhou, D. Kumar, S. Bakshi, J. Mason, A. Miller, and M. Bailey. Erays: reverse engineering ethereum’s opaque smart contracts. In *Proc. of USENIX Security 2018*, pages 1371–1385. Usenix Association, 2018.
- [50] W. Zou, D. Lo, P. S. Kochhar, X.-B. D. Le, X. Xia, Y. Feng, Z. Chen, and B. Xu. Smart contract development: Challenges and opportunities. *IEEE Transactions on Software Engineering*, pages 1–1, 2019.
- [51] F. Zuo, X. Li, P. Young, L. Luo, Q. Zeng, and Z. Zhang. Neural machine translation inspired binary code similarity comparison beyond function pairs. In *Proc. of NDSS 2019*. Internet Society, 2019.

## APPENDIX

### A. Output Example by Eth2Vec

We show an example of an output by Eth2Vec in Fig. 3. The interface of Eth2Vec follows that of Kam1n0 [8] as described in Section IV-F.

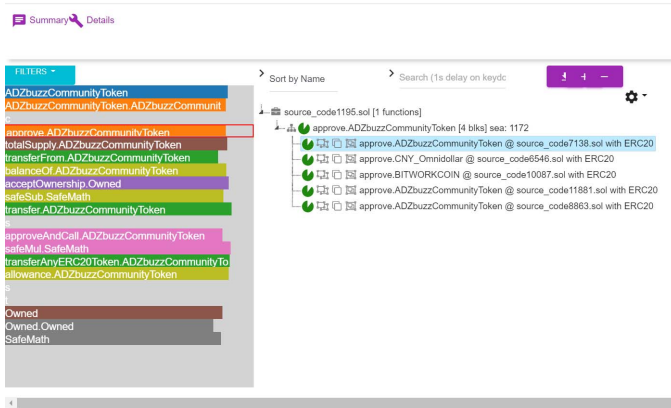


Fig. 3: Output example by Eth2Vec: The left-side on the windows displays a list of functions contained in each contract and the right-side shows a list of vulnerabilities included in the chosen function.

### B. Feature Selection by Momeni et al.

Momeni et al. [29] extracted 16 features from an abstract syntax tree (AST) of an Ethereum smart contract. Among them, we did not adopt Hexadecimal addresses as described in Section V. The whole list of vulnerabilities as follows, where all the information can be obtained from AST except for hexadecimal addresses:

- Lines of codes
- Contract definitions
- Function definitions
- Binary operations
- Function calls
- Blocks
- Expression statements
- Event definitions
- Bytes
- Elementary type addresses
- Modifier definitions
- Placeholder statements

- Modifier invocation
- Approve function definitions
- Constant values
- Hexadecimal addresses

### C. List of Vulnerabilities

TABLE III: List of vulnerabilities: In this paper, we target the following vulnerabilities for evaluation of Eth2Vec.

Name	Severity	Description
Reentrancy	3	External contracts should be called after all local state updates
Time Dependency	2	Miners can alter timestamps. Make critical code independent of the environment
ERC-20 Transfer	1	The contract throws where the ERC20 standard expects a bool. Return <code>false</code> instead
Gas Consumption	1	A transaction fails by exceeding an upper bound on the amount of gas that can be spent
Implicit Visibility	1	Functions are public by default. Avoid ambiguity: explicitly declare visibility level
Integer Overflow	1	The return value is not checked.
Integer Underflow	1	Always check return values of functions

In this paper, we target the vulnerabilities described in Table III for evaluation of Eth2Vec. The columns of Severity and Description follow descriptions in SmartCheck [43] except for gas consumption. Description about the gas consumption follows description in [https://consensys.github.io/smart-contract-best-practices/known\\_attacks/](https://consensys.github.io/smart-contract-best-practices/known_attacks/).