

ESCORT: Ethereum Smart CONtRaCTs Vulnerability Detection using Deep Neural Network and Transfer Learning

Oliver Lutz
University of Würzburg

Huili Chen
University of California, San-Diego

Hossein Fereidooni
Technical University of Darmstadt

Christoph Sendner
University of Würzburg

Alexandra Dmitrienko
University of Würzburg

Ahmad Reza Sadeghi
Technical University of Darmstadt

Farinaz Koushanfar
University of California, San Diego

ABSTRACT

Ethereum smart contracts are automated decentralized applications on the blockchain that describe the terms of the agreement between buyers and sellers, reducing the need for trusted intermediaries and arbitration. However, the deployment of smart contracts introduces new attack vectors into the cryptocurrency systems. In particular, programming flaws in smart contracts can be and have already been exploited to gain enormous financial profits. It is thus an emerging yet crucial issue to detect vulnerabilities of different classes (e.g., reentrancy or multiple send bugs) in contracts in an effective and efficient manner. Existing machine learning-based vulnerability detection methods are limited and only inspect whether the smart contract is vulnerable, or train individual classifiers for each specific vulnerability, or demonstrate multi-class vulnerability detection without extensibility consideration. To overcome the scalability and generalization limitations of existing works, we propose ESCORT, the first Deep Neural Network (DNN)-based vulnerability detection framework for Ethereum smart contracts that supports lightweight *transfer learning* on unseen security vulnerabilities, thus is *extensible* and *generalizable*. ESCORT leverages a *multi-output* neural network architecture that consists of two parts: (i) A common feature extractor that learns the semantics of the input smart contract; (ii) Multiple branch structures where each branch learns a specific vulnerability type based on features obtained from the feature extractor. We perform a comprehensive evaluation of ESCORT on various smart contracts. Experimental results show that ESCORT achieves an average F1 score of 95% on six vulnerability types and the detection time is 0.02 seconds per contract. When extended to new vulnerability types, ESCORT yields an average F1 score of 93%. To the best of our knowledge, ESCORT is the first framework that enables transfer learning on new vulnerability types with minimal modification of the DNN model architecture and re-training overhead.

KEYWORDS

Blockchain, Ethereum Smart Contracts, Vulnerability Detection, Deep Neural Network, Transfer Learning

1 INTRODUCTION

The success of Bitcoin [64] fueled the interest in the cryptocurrency platforms. As a result, next generation blockchain-powered application platforms emerged, such as Ethereum [6] and Hyperledger [9].

Beyond financial transactions, these platforms provide smart contracts that are automated decentralized applications describing the terms of agreements and the transaction flow between the buyers and the sellers.

Generally, blockchains are append-only distributed and replicated databases, which maintain an ever-growing list of immutable and tamper-resistant data records (e.g., financial transactions and smart contracts). Immutability and tamper-resistance of blockchains stem from their append-only property and are paramount to the security of blockchain applications. In the context of digital currency and payments, they ensure that all the involved parties have access to a single history of payment transactions in a distributed setting and that such a history cannot be manipulated. In smart contract systems, immutability and tamper-resistance properties enforce "code is law" principle, meaning that conditions recorded in a smart contract are not to be modified since they have been written and published.

However, these properties bear their own security risks and challenges: First, smart contracts are written in error-prone programming languages such as Solidity [22], and can contain exploitable programming errors/bugs that are often overlooked or detected only after deployment on the blockchain, and cannot be simply fixed. Second, Ethereum operates on open networks where everyone can join without trusted third parties while smart contracts are often in control of significant financial assets. Hence, smart contracts are attractive and easy attack targets for adversaries to gain financial profits [30]. The consequences of bug exploitation may have global effects on the entire underlying blockchain platform, far beyond the boundaries of individual contracts. For instance, vulnerabilities in a single smart contract, the DAO [49], affected the entire Ethereum network, when in June 2016 the attacker exploited a reentrancy bug and had withdrawn most of its funds worth about 60 million US dollars [35, 78]. In the aftermath, the value of Ether, Ethereum's cryptocurrency, dropped dramatically [72], and the postulated "code is law" principle was undermined through the deployment of a hard fork – a manual intervention orchestrated by a notable minority, the team of Ethereum core developers. The Ethereum blockchain was thus split into two versions, Ethereum and Ethereum Classic, which are maintained in parallel since then. In another case, a critical bug accidentally triggered in 2017 resulted in the freezing of more than \$280M worth of Ether in the Parity multisig wallet [79].

Once it was understood that real-world attacks and even innocent mistakes can lead to fatal economic loss, the problem of detecting smart contract vulnerabilities became very appealing to the research community. A wide range of automated tools were developed to help find vulnerabilities in smart contracts using various techniques, such as symbolic execution [59, 84], satisfiability modulo theories (SMT) solving [29], data flow analysis [39], runtime monitoring [37], and fuzzing [40], to name a few examples. However, present methods provide a limited trade-off between detection effectiveness and efficiency. For instance, symbolic execution-based vulnerability detection is slow since all program paths need to be examined [17], whereas data flow analysis has limited coverage [39]. Moreover, individual tools normally cover a limited set of vulnerabilities, hence to achieve sound testing, one would normally need to apply several tools.

Recently, Machine Learning (ML) has attracted the attention of security researchers due to its capability to learn the hidden representation from the abundant data [28]. Prior works have shown the effectiveness of ML techniques for detecting vulnerabilities in software written in C and C++ languages [74, 85]. There are also first attempts to apply ML techniques for the detection of smart contract vulnerabilities [48, 80, 87]. However, as we elaborate in details in Section 9, the existing solutions suffer from either of the following shortcomings: (i) They are inherently *unscalable* and *inflexible*, as the inclusion of any new vulnerability types would require training of new models; or (ii) They only distinguish between vulnerable and non-vulnerable smart contracts (i.e., binary classification), without the ability to detect vulnerability types; or (iii) The tools require source code of the smart contract, which limits their applicability scope.

Our goal and contributions. In this paper, we aim to address the deficiencies and limitations of existing solutions and propose ESCORT, the first Deep Neural Network (DNN) based framework for smart contract vulnerability detection that has the following properties: (i) Operates on bytecode of smart contracts and does not require access to the source code; (ii) Distinguishes safe contracts from vulnerable ones with one or more vulnerabilities; (iii) Automatically identifies the vulnerability types of detected vulnerabilities; and (iv) Presents an innovative multi-output model architecture that enables fast model extension to new vulnerabilities using the concept of *transfer learning*.

In particular, we make the following contributions:

- ESCORT enables *efficient* and *scalable* multi-vulnerability detection of smart contracts. It employs a *multi-output* architecture where the feature extractor learns the program semantics and each branch of the DNN captures the semantics of a specific vulnerability class. The defender only needs to train a single DNN based on smart contract bytecode to detect multiple vulnerability types.
- ESCORT is the first DNN-based framework that supports lightweight transfer learning on new vulnerability types, thus is *extensible* and *generalizable*. Our multi-output architecture can be easily expanded by concatenating a new classification branch to the feature extractor. Only the new branch needs to be trained on the new dataset during transfer learning.

- ESCORT is automated for inspecting vulnerabilities in smart contracts prior to their deployment and demonstrates superior vulnerability detection performance while incurring low overhead. We perform a comprehensive evaluation of ESCORT across various Ethereum smart contracts to corroborate its effectiveness, efficiency, and extensibility. Our framework can detect vulnerabilities in 0.02 seconds per smart contract and yields an average F1 score of 95% across all evaluated classes.
- We evaluate ESCORT using a dataset consisting of 93,497 smart contracts downloaded from Ethereum blockchain and labeled them using three existing vulnerability detection tools: Oyente [17], Mythril [14], and Dedaub [2]. To construct and label such a dataset, we developed a toolchain ContractScraper that automates the acquisition and labeling process. ContractScraper’s modular structure enables one to easily integrate other tools for labeling. We will open source ContractScraper and our dataset to encourage further research on smart contract security.

Overall, ESCORT is the first flexible and generalizable smart contract detection technique with superior vulnerability detection performance. **Outline.** The remaining part of the paper is organized as follows: In Section 2, we provide the necessary background information on smart contracts and their vulnerabilities, as well as give insights into deep learning. We define the threat model and identify the design challenges in Section 4. Section 5 sheds light into ESCORT’s design specifics, while Section 7 and 8 provide implementation details and evaluation results, respectively. After that, the relevant related work is surveyed in Section 9. We finally conclude in Section 10.

2 BACKGROUND

We introduce smart contracts and their vulnerabilities in Section 2.1, and the background on deep learning (DL) in Section 2.2.

2.1 Smart Contracts and Vulnerabilities

A smart contract is written in high-level programming languages such as Solidity [22] and is called by its address to run operations on the blockchain. Once compiled, the bytecode of the contract is generated and executed inside the Ethereum Virtual Machines (EVM). Since there is a one-to-one mapping between a blockchain operation and bytecode representation, it is feasible to analyze the control flow of a contract at the bytecode-level. When triggered, the execution of the smart contract is autonomous and enforceable for all participating parties [36].

The EVM itself is a stack-based machine with a word size of 256 bits and stack size of 1024 [88]. The memory applies a word-addressable model. Once a contract is deployed on the blockchain, it requires gas to function. Gas is the unit used to pay the computational cost of the miners running contracts or transactions and is paid in Ether.

Similar to any other software, smart contracts might suffer from vulnerabilities and programming bugs. The Smart Contract Weakness Classification (SWC) Registry [16] collects information about various vulnerabilities. We differentiate five categories of vulnerability types: External Calls, Programming Errors, Execution Cost, Influence by Miners, and Privacy.

External Calls. Any public function of a smart contract can be called by any other contract. A malicious user can then exploit public availability to attack vulnerable functions of smart contracts. A prominent example is the so-called reentrancy bug (SWC-107 [16]). Here, an attacker can call a contract’s function multiple times before the initial call is terminated. If the internal contract state is not securely updated, the attacker can drain Ether from the contract by recursively calling the function.

Programming Errors. Some of the programming errors in smart contracts are very similar to those found in traditional programs, such as missing input validation, typecast bugs, use of untrusted inputs in security operations, unhandled exception, exception disorder, and Integer overflow and underflow vulnerabilities. In another example, an *assert* function used in tests and not removed by the programmer in the release version may lead to its misuse by an attacker, which can result in exploitable error conditions (SWC-110 [16]). Other vulnerabilities can be specific to smart contracts.

Examples are greedy contracts that lock Ether indefinitely, gas-less send bug that does not provide sufficient gas to execute the code of the smart contract, Ether lost in transfer if sent to unknown recipients, etc. Further examples are *callstack depth* limit reached exception bug and unprotected *selfdestruct* instructions (see SWC-106 [16]), where an attacker can call a smart contract’s public function containing a *selfdestruct* to terminate the smart contract, or he can fill up the stack to reach the stack size limit. Both attacks result in a Denial of Service (DoS) of vulnerable smart contracts.

Execution Cost. Every transaction on the Ethereum network costs gas. However, every block has a spendable gas limitation. An attacker can use this limit to induce a DoS of a vulnerable contract. For example, if the execution time of a function is dependent on input from the caller, a malicious caller can expand the execution time of the smart contract over the gas limit (SWC-128 [16]). Thereby, execution is terminated by exceeding the gas limit before it is finished. Another way an attacker can misuse the gas limit per block is to induce an error on a *send* call. If a programmer bundles *multiple sends* in one function of the smart contract, the attacker can then prevent the execution of other *send* calls in the function.

Influence by Miners. Miners are entities that actually execute transactions on the blockchain. They can decide which transactions to execute, in what order, and also able to influence environment variables (e.g., timestamps). To illustrate the problem, let us assume a scenario where a smart contract is instructed to send Ether to the first user that solved the puzzle. If two users commit a transaction with the solution at the same time, a miner decides who will be first and therefore will be getting the Ether (SWC-114 [16]). This vulnerability type is generally referred to as Transaction Order Dependence (TOD).

Privacy. Solidity [22] offers different visibility labels for functions and variables. Most notably, a programmer can define a function as private or public. The default setting for functions is public, which can be overlooked by the programmer (SWC-100 [16]). If a variable or function is set to private, it can’t be seen or changed by other contracts. However, even if it is set as private, an attacker can still parse the blockchain data, where those variables are stored in plaintext (SWC-136 [16]).

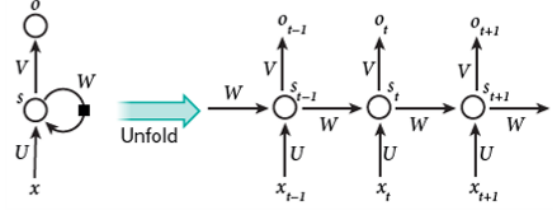


Figure 1: A recurrent neural network and its unfolding in time [56]. The computation of the forward pass is shown.

2.2 Deep Learning

ESCORT operates on the bytecode representation of the smart contracts, which can be considered as a special case of text data. We introduce background about text representation and recurrent neural networks below.

Text Representation. The text modality is typically transformed into numerical vectors for usage in ML algorithms. This transformation can be realized in different ways, such as a bag of words [86], n-gram language model [33], and embedding layer [90]. The numerical vectors converted from the text data are then used as the direct input to the DL models.

Recurrent Neural Network. An RNN is a category of DNNs where the connections between nodes (i.e., neurons) construct a direct computational graph along the temporal sequence [52]. A key property of the RNN is that the model can use its internal states as memory cells to store the knowledge about prior inputs, thus capturing the contextual information in the sequential input data (e.g., text document). As shown in Figure 1, the hidden states obtained from the previous input (s_t) affects the output in the current time step (o_{t+1}). The unfolded RNN diagram reveals the ‘parameter sharing’ mechanism of RNNs where the weight matrices (W , V , and U) are shared across different time steps. Parameter sharing makes RNNs generalizable to unseen sequences of different lengths.

Multi-label vs. Multi-class classification. Smart contracts vulnerability detection can be realized with two different paradigms. The first one is known as *Multi-class classification*, which refers to the case where the classification task has more than two classes that are mutually *exclusive*. In particular, each sample is assigned with one and only one class label. The second one, *Multi-label classification*, also involves multiple classes while a data sample can have more than one associated labels. This is because the classes in multi-label tasks describe non-exclusive attributes of the input (e.g., color and length).

Let us use an example to illustrate the difference between these two paradigms. Given a clothing dataset with three colors (black, blue, red) and four categories (jeans, dress, shirt, shoes), we want to train a model to predict these two clothing attributes simultaneously. Figure 2 shows the architecture of the multi-class and multi-label formulation of the clothing classification task. The multi-class design has only one set of dense layers (i.e., ‘heads’) at the bottom of the DNN where the last dense layer has $3 \times 4 = 12$ neurons. The network topology for multi-label classification has two sets of dense heads at the end of the DNN where the last dense layer in each branch has 4 and 3 neurons to learn the clothing category and color attribute, respectively. We call the network design of

Name	Method	Capability	Extensible	Required Input	Detection time per Contract (s)	F1-score
This work	ML (LSTM) + transfer learning	Multi-label	Yes	Bytecode	0.02	0.95
Oyente [59]	Symbolic execution	Multi-class	No	Source code and bytecode	350	0.38
Mythril [63]	Symbolic execution, taint analysis, and SMT	Multi-class	No	Bytecode	11.1	0.47
Dedaub [3]	Flow and loop analysis	Gas-focused vulnerability	No	Source code	20	NA
Securify [84]	Symbolic analysis	Binary decision	No	Bytecode	30	0.54
Vandal [32]	Logic-driven static program analysis	Multi-class	No	Bytecode	4.15	NA
ContractWard [87]	ML (bigram model)	Binary decision	No	Opcode	4	0.96
Towards Sequential [80]	ML (LSTM)	Binary decision	No	Opcode	2.2	0.86
NLP-inspired [42]	ML (AWD-LSTM)	Multi-class	No	Opcode	NA	0.9
Color-inspired [48]	ML (CNN)	Multi-label	No	Bytecode	1.5	0.94
Graph NN-based [92]	ML (GNN)	Multi-class	No	Source code	NA	0.77

Table 1: Qualitative comparison of ESCORT and existing smart contract vulnerability detection methods.

the *multi-label* classification with multiple sets of dense heads as ‘*multi-output*’ architecture throughout this paper. The ‘*stem-branch*’ topology makes the multi-output architecture extensible to learn new attributes. ESCORT leverages this observation to devise an efficient and extensible smart contract inspection solution.

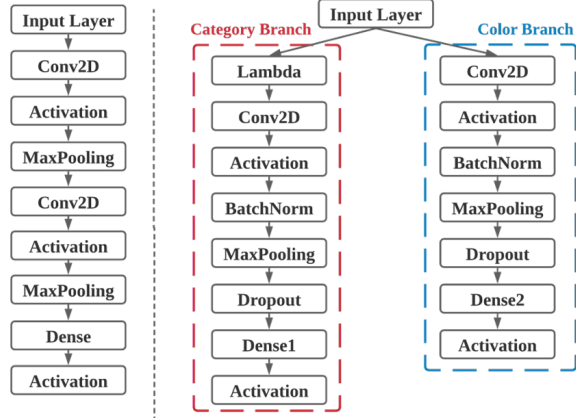


Figure 2: DNN architectures of multi-class (left) and multi-label (right) formulation of the clothing classification task [10].

3 MOTIVATION

Vulnerability detection of smart contracts is a challenging task. Table 1 summarizes the properties and performance metrics of previous detection tools. To the best of our knowledge, none of the existing works takes into account of the extensibility requirement of vulnerability detection in their design. This means that when new vulnerabilities of smart contracts are discovered and exploited by the adversary, the detection tools cannot quickly adapt to the new contracts. ESCORT is motivated to provide an efficient and extensible solution to concurrent detection of multiple vulnerability types using DL techniques. To this end, we first investigate two research questions (RQs) about smart contract inspection in this section and answer them in details in Section 5.

RQ1: How to build a single DL model for detecting multiple vulnerability types?

Prior works have attempted to use ML algorithms to inspect the vulnerabilities in smart contracts. An LSTM-based detection approach is proposed in [80] to distinguish vulnerable contracts from safe ones (i.e., binary classification), which does not meet the goal of RQ1. ContractWard [87] uses ‘One vs. Rest’ algorithms and trains individual ML models for each vulnerability class. Therefore, ContractWard [87] does not meet the model number constraints in RQ1. The paper [48] suggests to transform contract bytecode to RGB images and employs a Convolutional Neural Network (CNN) for binary classification. Their method can support multi-label classification by re-training the obtained CNN on the corresponding dataset, thus satisfying both requirements in RQ1.

RQ2: How to make the DL model extensible to new vulnerabilities?

ContractWard [87] can be adapted to detect a new vulnerability type by training a *new ML model* on the n-gram features from scratch. The color-inspired inspection method [48] does not consider the extensibility requirement of the detection. However, if applying the traditional transfer learning approach in the ML domain, the defender may replace the last dense layer of the pre-trained CNN with a new one. The new layer is then trained on the mixture of the old vulnerability data and the new one. Note that training a new ML model from scratch or fine-tuning a pre-trained ML model on a large dataset incurs high computational cost, thus is not desirable in practice.

4 THREAT MODEL AND CHALLENGES

We introduce our threat model in Section 4.1 and the challenges of developing an effective vulnerability detection technique for smart contracts in Section 4.2.

4.1 Threat Model

We consider a scenario where the attacker is a malicious party that can obtain knowledge from any public data structure in the blockchain and can upload his contract code to the Ethereum system. The exact attack requirements and actions that the adversary need to take for the exploitation of a smart contract are specific

to a vulnerability class, which we explain in Section 8.1. Note that ESCORT is generic and we demonstrate its effectiveness with eight common vulnerability types in this paper. The defender can be an Ethereum contract designer that aims to ensure the program is not exploitable by malicious adversaries during the code development time. This role can also be taken by the end user that intends to verify the security of the contract at runtime before sending any transactions to it. We assume the performance metrics reported in the previous papers [43, 59, 62] as well as the open-sourced implementation of existing detection tools [2, 14, 17] are reliable. This assumption is feasible since expert inspection has been performed to cross-validate the performance of proposed detection methods in the previous works. ESCORT relies on this assumption since our method is an instantiation of supervised learning paradigm.

ESCORT aims to provide the defenders with a holistic smart contract vulnerability detection solution. To this end, we formulate vulnerability inspection as a *supervised multi-label classification* problem where the input is the contract (can be represented in high-level language, opcodes, or bytecodes) and the output is the corresponding vulnerability types as introduced in Section 8.1.

4.2 Challenges

We identify the challenges to develop an effective DNN for automated smart contract vulnerability detection below.

(C1) Data Collection and Pre-processing. Supervised learning of a DNN requires a sufficiently large labeled dataset. Analyzing smart contracts using the source code is difficult since not that many smart contracts are open-sourced. On another hand, blockchain platforms typically host their smart contracts in a form of a bytecode, which is publicly available. However, performing analysis on the bytecode raises another challenge due to the large bytecode length of long contracts. This long sequence requires a large memory footprint during model training, suggesting that using the bytecode as direct input to the DNN model is impractical.

(C2) Feature Extraction. The second challenge concerns the problem of finding the proper feature representation of the smart contract program. On the one hand, manual design of contract features is time-consuming and has limited efficacy, since the contract bytecode is long and hard to interpret by human developers. On the other hand, current implementations of automated feature extraction [80, 87, 92]

(Section 9) mainly apply traditional software testing techniques on the smart contract without exploring its domain-specific properties, thus taking a long time to inspect a single contract.

(C3) Dataset Imbalance. The third challenge is a class imbalance in the training set. Prior works have shown that the number of contracts with specific vulnerability types is much lower than the one of non-vulnerable contracts [87]. Learning the characteristics of vulnerable contracts with a DNN is challenging since the stochastic gradient descent (SGD) based learning of DNN models is inherently biased towards the majority class, while our objective is to recognize the minority class (i.e., vulnerable contracts). As such, it is crucial to provide the DNN model with sufficient vulnerable contracts to ensure a high true positive rate.

(C4) Efficiency. The fourth challenge is to ensure the efficiency of DNN training and inference for concurrent detection of multiple

vulnerability types. Identifying diverse attacks with a single DNN detector is challenging since different vulnerabilities exploit distinct loopholes in the contract, which might be hard to capture with a conventional DNN. Efficiency is important for practical development of the contract scanner since devising individual classifiers for each vulnerability class, as done, e.g., in [87], is unscalable and incurs large computation overhead.

(C5) Extensibility/Generalizability. Finally, smart contract inspection should be capable of learning new vulnerability types quickly while preserving the knowledge of the known ones. We call this requirement ‘extensibility’. This property is important for both, researchers and practitioners, since new attacks on smart contracts are emerging at a fast speed. Augmenting an existing contract detector to new attacks is non-trivial since the new attack exploits the unseen and unpredictable susceptibilities of smart contracts compared to the previously known attacks. Training a new DNN from scratch to accommodate the new vulnerabilities consumes extensive resources and incurs additional engineering costs.

We show how ESCORT framework tackles each of the above challenges in Section 5.

5 ESCORT DESIGN

We propose ESCORT, the first extensible and transfer learning-friendly DNN-based framework for vulnerability detection of Ethereum smart contracts. The key innovation of ESCORT is that we decompose the task of vulnerability detection into two subtasks: (T1) Learning the bytecode features of general contracts (*attack-agnostic*); (T2) Learning to identify each particular vulnerability class (*attack-specific*). To achieve (T1), we design a *feature extractor* that captures the semantic and syntactic information of contract bytecode regardless of its vulnerabilities. To perform (T2), we devise an individual *vulnerability class branch* to characterize susceptibility given the bytecode features extracted in (T1). Our **divide-and-conquer** design is highly modular and flexible compared to previous detection techniques as we corroborate in Section 8.

Below we discuss each step of ESCORT design, which aims to overcome the design challenges in Section 4.2.

5.1 ESCORT Global Flow

Figure 3 demonstrates the global overview of ESCORT smart contract detection technique. Our framework has three stages: DNN classifier training, transfer learning, and deployment. We discuss each stage below.

Training. The top part of Figure 3 shows ESCORT’s training pipeline. To enable supervised learning for vulnerability detection, the defender first constructs the smart contracts bytecode dataset with corresponding labels, as detailed in Section 6. The defender specifies the system parameters, including the vulnerabilities of interests and the available hardware resources for ESCORT’s multi-output DNN design. Finally, the devised model is trained on the collected contract data with their corresponding labels, resulting in a converged DNN detector.

Transfer Learning. The middle part of Figure 3 shows the transfer learning phase of ESCORT. Given a trained detector and the bytecode of smart contracts with new vulnerabilities, ESCORT extends the DNN architecture devised in the original training phase

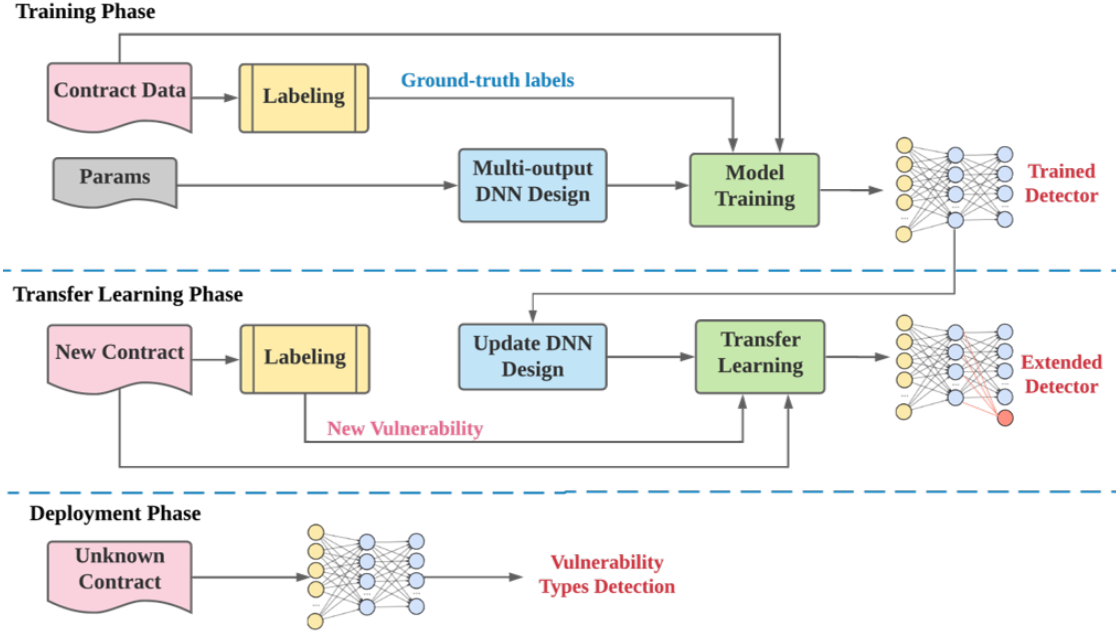


Figure 3: Overview of ESCORT framework for Ethereum smart contract vulnerability detection. ESCORT has three stages: training, transfer learning, and deployment as shown in the top, middle, and bottom part of the figure.

with new parallel branches. The layers in the expanded branch are then trained on the new vulnerability data with the associated labels to perform transfer learning. To the best of our knowledge, ESCORT is the first framework that supports transfer learning to accommodate new vulnerability types of the smart contracts.

Deployment. The bottom part of Figure 3 shows the workflow of ESCORT’s deployment stage. After the training/transfer learning phase completes, ESCORT returns a trained DNN classifier that can detect whether an unknown smart contract has any of the learned vulnerability types.

5.2 Neural Network Design

Prior ML-based vulnerability detection techniques only explored simple architectures, such as K-Nearest Neighbors, SVM, Decision Trees [87], or use a CNN [48]/RNN [80] to decide whether a smart contract is vulnerable or not. To overcome these constraints, ESCORT aims to design an **extensible** DNN detector that: (i) Provides the probability that the smart contract has certain vulnerabilities, instead of making a binary decision about contract security; (ii) Classifies multiple vulnerabilities with a single DNN. To this end, we propose an innovative **multi-output architecture** for **concurrent** detection of multiple vulnerability types. This neural network design step is shown by the ‘Multi-output DNN Design’ module in ESCORT’s global flow (Figure 3).

Figure 4 illustrates the generic architecture of our DNN detector. The stem and branch layers shown in the diagram are typical DNN layers such as the Dense layer, Dropout layer, and GRU layer. The multi-output model has two main components and we discuss each one below:

(i) **Feature Extractor.** The first component of ESCORT’s extensible DNN model is the common feature extractor (i.e., ‘stem’) shared by all the bottom branches. The feature extractor is a stack

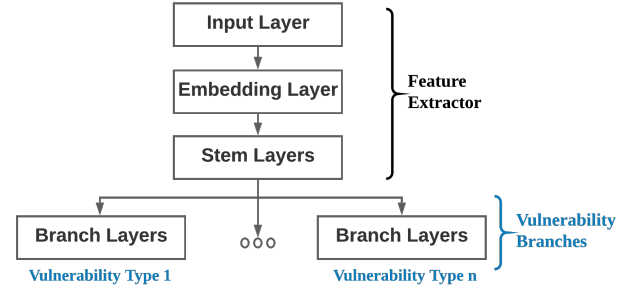


Figure 4: General architecture of ESCORT’s multi-output model for concurrent detection of multiple vulnerability types.

of layers that learn the fundamental features in the input data that are general and useful across different attributes. In the context of smart contracts, the feature extractor is trained to learn the semantic and syntactic information from the contracts’ bytecode. To this end, we incorporate several key layers in ESCORT’s feature extractor:

- **Embedding Layer.** The bytecode of smart contracts are long hexadecimal numbers, while DNNs typically work with fractional numbers to achieve high accuracy. The embedding can solve this discrepancy since it stores the word embedding in the numerical space and retrieves them using indices [90]. The embedding layer provides two key benefits: (i) Compressing the input via a linear mapping, thus reducing the feature dimension; (ii) Learning bytecode in the embedding space (fractional numbers). This facilitates representation exploration and gathers similar bytecode in the vicinity of each

other. ESCORT leverages the advantages of the embedding layer to capture the semantics in the input bytecode.

- **GRU/LSTM Layer.** The stem layers and branch layers in Figure 4 can include GRU/LSTM layers for processing sequential inputs. Gated Recurrent Units (GRU) and Long Short-Term Memory (LSTM) are two typical layers in recurrent neural networks that help to overcome the short-term memory constraint and vanishing gradient problem [52] using a ‘gating’ mechanism. More specifically, both types of layers have internal gates that regulate the information flow along the time sequences and decide which data shall be kept/forgotten. We mainly use GRU layers in ESCORT’s DNN design.

(ii) **Vulnerability Branches.** The second component of ESCORT’s multi-output DNN architecture is the **ensembling** of multiple vulnerability branches. Each branch is a stack of layers that are trained to learn the patterns/hidden representation of the corresponding vulnerability class. While there is no direct dependence between different branches, they share the same feature extractor, i.e., the input to each branch is the same. This is feasible since the branch input (which is also the feature extractor’s output) shall capture the semantics in the contract’s bytecode, which is common/general information useful for different vulnerabilities.

ESCORT’s ‘stem-branches’ architecture is similar to the ‘*mixture of experts*’ paradigm where the problem space is divided into homogeneous regions and individual expert models (learners) are trained to tackle each sector [54]. The main difference between ESCORT’s multi-output design and the mixture of experts model is that the latter one requires a trainable gating network to decide which expert shall be used for each input region, while our DNN model does not need such a gating mechanism since we aim to detect multiple vulnerability types of the input contract in parallel.

Note that the last layer of each vulnerability branch is a Dense layer with one neuron. The *sigmoid* evaluation of this neuron’s activation value gives the *probability* that the input contract has the specific vulnerability. As such, ESCORT engenders detection results with better **interpretability** by providing the confidence score for its diagnosis instead of the binary decision about vulnerability existence.

In summary, ESCORT’s multi-output architecture solves the feature extraction, efficiency and extensibility challenges (C2, C4, C5) identified in Section 4.2. In particular, ESCORT allows the defender to train a *single* DNN for detecting multiple vulnerability types instead of training an individual classifier for each attack, thus demonstrating superior efficiency compared to the prior work [48, 80, 87]. ESCORT design incurs minimal non-recurring engineering cost and is scalable as new vulnerabilities are identified.

5.3 Transfer Learning

Malicious parties have a strong incentive to discover and exploit new vulnerabilities of smart contracts due to the associated prodigious profits. As such, the contract inspection technique shall be extensible to learn new vulnerabilities as they are identified. We propose *transfer learning* as the solution to the challenge (C5) in Section 4.2. More specifically, we suggest to *expand* the pre-trained multi-output DNN model by adding new vulnerability branches for transfer learning. This process is demonstrated in the middle part

of Figure 3. The transfer learning capability of ESCORT ensures that our detection framework can be upgraded with the minimal cost to defend against emerging attacks on smart contracts.

Our transfer learning stage has two goals:

(G1) **Preserve Knowledge on Old Vulnerabilities.** On the one hand, the DNN detector shall retain knowledge about the previous vulnerability types that are used in the initial training phase. This property is important since ESCORT aims to provide a holistic and extensible solution to concurrent detection of multiple vulnerabilities. As such, maintaining high classification accuracy on the known attacks is essential.

(G2) **Learn New Vulnerabilities Quickly.** On the other hand, transfer learning aims to adapt the pre-trained model to achieve high accuracy on the new dataset in an efficient way. This is also required by the extensibility challenge (C4) in Section 4.2. To achieve fast adaptation, transfer learning shall yield minimal runtime overhead. This requirement is crucial for practical deployment, since training a new DNN model from scratch for the new vulnerabilities is prohibitively expensive and hard to maintain.

ESCORT’s transfer learning phase works as follows. When a new vulnerability is identified, the defender constructs a new training dataset accordingly and updates the converged DNN detector by adding a new vulnerability branch (i.e., the stack of layers). During transfer learning, the parameters of the common feature extractor and existing vulnerability branches are *fixed*. Only the parameters in the newly added branch are updated with the new vulnerability dataset. Freezing the feature extractor and the converged branches ensures that the updated DNN classifier preserves the detection accuracy on the old vulnerabilities (G1), and training a new branch enables the updated model to learn the new attack (G2).

Besides extensibility, ESCORT also enables lightweight and fast adaptation when *mode drift* occurs. Smart contracts running on Ethereum are dynamic and change over time, which might lead to a decrease of ESCORT’s performance. To alleviate the model drift concern, the contract developer can update the parameters of the vulnerability branches given a set of labeled new contracts while fixing the ones of the feature extractor.

6 DATASET CONSTRUCTION TOOLCHAIN

In this section, we present ContractScraper, the toolchain we built to construct the labeled smart contract dataset and to address the challenge (C1) defined in Section 4.2. Details about our dataset used in DNN training is given in Section 8.1. It is worth mentioning that a standard and public smart contract vulnerability dataset was not available before. We will open source ContractScraper toolchain and our dataset used in the evaluation to facilitate the development and comparison between emerging detection techniques.

6.1 Design Choices

To build a sufficiently large training set for our supervised DNN training, we make the design choice to work on the bytecode-level since the bytecode of smart contracts are publicly available. The accessibility of contract bytecode solves the challenge (C1) in Section 4.2. This also makes our approach agnostic to the programming languages of smart contracts, since smart contracts are written in different languages (e.g., Solidity [22], Viper [26], and Serpent [21])

get eventually compiled to the same bytecode. The EVM bytecodes are executed in a stack context, thus the control flow of a program at the bytecode level contains useful information for detection.

ContractScraper is designed to obtain the bytecode files of smart contracts from the Ethereum platform, label them using available bytecode-level detection tool(s), and store the result in a database. Note, that we concentrate on Ethereum smart contracts for exemplary purposes in this paper. Generally, ESCORT can also be used for vulnerability detection in smart contracts of other cryptocurrency platforms that use Ethereum-compatible EVM, such as Quorum [18], VeChain [25], Rootstock [19], and Tron [24], to name a few. This is possible since the bytecode of smart contracts on these platforms is compatible with Ethereum EVM.

6.2 ContractScraper Architecture

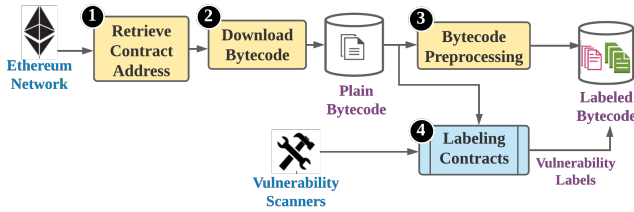


Figure 5: Generic workflow of ContractScraper for smart contract acquisition and labeling.

We show the generic workflow of ContractScraper toolchain in Figure 5. In step (1), the addresses of contracts are retrieved from the blockchain. Step (2) involves downloading the bytecode from the Ethereum network by address and extracting information into a database. In step (3), the bytecode is pre-processed for input efficiency. The last step (4) outputs the vulnerability types of the contracts using the bytecode-level detection tools. It is worth noticing that any vulnerability detection techniques that take bytecode as input can be used by ContractScraper for contract labeling, including existing methods such as Oyente [17], Mythril [14], Dedaub [2], and our proposed new method ESCORT. We show in Section 8 that ESCORT outperforms existing tools in terms of both detection effectiveness and efficiency.

For a concrete instantiation of the generic flow, we use the open-source tool Ethereum ETL [4, 5] to retrieve the addresses of contracts in step (1). The ETL tool connects to the Ethereum network and exports blockchain content into CSV files. For step (2), we also utilize existing tools: Multiple Web APIs exist to download the contracts’ bytecode by its address. We opt for two APIs for instantiating the bytecode downloading module of ContractScraper: Infura [1] and Dedaub’s Contract-Library [2]. For the instantiation of step (3), we develop an assistive Python module named *Contract Loader* to extract the information from smart contracts into a MySQL database [13]. In the last step (4), we use three tools for data labeling: Oyente [17], Mythril [14], and Dedaub [2]. We do not consider tools that scan smart contracts on the source code level since decompilation might result in information loss. The modular structure of ContractScraper makes it possible to easily extend the toolchain with other bytecode-level vulnerability detection tools including ESCORT.

6.3 Bytecode Acquisition

For dataset construction, ContractScraper first utilizes *Infura* API to download smart contracts. We were able to download 1.155.085 bytecodes from the first 5 million Ethereum blockchain blocks. This number corresponds to about 96% of the smart contracts in the targeted blocks. In some cases, the download resulted in an empty bytecode 0x. The possible reasons could be: (i) The Ethereum node is not fully synced with the network, thus the bytecode is not available; (ii) An empty contract is deployed; (iii) The smart contract is self-destructed. To tackle the situation where the smart contracts’ bytecode exists but cannot be downloaded via Infura API, ContractScraper uses another API provided by Dedaub tool to retrieve the missing bytecode files. Finally, at the end of the bytecode acquisition process, 1.156.611 smart contract bytecode files are available for ESCORT’s vulnerability analysis.

6.4 Bytecode Preprocessing

The downloaded bytecode consists of hexadecimal digits that represent particular operation sequences and parameters. In the pre-processing step, ContractScraper first transforms the collected raw bytecodes to sequences of operations divided by a unique separator and removes the input parameters from the bytecode to reduce the input size. Furthermore, it merges operations with the same functionality into one common operation. For instance, the similar commands *PUSH1* - *PUSH32* commands (represented by the bytes 0x60-0xf) are replaced with the *PUSH* operation (represented by 0x60). Note that some hexadecimal digits in the crawled bytecode do not correspond to any operations defined in the Ethereum Yellow Paper [88]. These bytes are considered as invalid operations and substituted with the value *XX*.

6.5 Labeling of Smart Contracts

A smart contract might have multiple vulnerabilities as introduced in Section 2.1. Each of the vulnerability detection tools used by ContractScraper is specialized for detecting a specific set of vulnerability types. Note that besides the bytecode-level detection tools themselves, ContractScraper also store the performance metrics (e.g., F1 score) of each tool on each vulnerability type that they can detect. The performance characterizations are obtained from the previous publications [2, 14, 17] with experts’ manual inspection to ensure the correctness. To determine if a given smart contract has a specific vulnerability type, ContractScraper selects the detection tool that features the highest F1 score on this vulnerability among all available ones and use it for contract labeling.

ContractScraper repeats the above process for each contract and each vulnerability type. We develop Python modules to perform the above task. In the end, 1.156.611 smart contracts are labeled. It is worth repeating that since there are no dependencies among the vulnerability scanning tools, the set of vulnerability types can be easily extended with other available tools.

7 IMPLEMENTATION

In this section, we instantiate the generic design of ESCORT described in Section 5 on eight vulnerability types and elaborate on implementation details.

7.1 Implementation of the DNN Model

We begin with the instantiation of the ESCORT’s DNN model. We utilize the *tf.keras* package [23] for model building, training, and inference. Furthermore, we devise a model serving API that takes the plain bytecode as input and returns the vulnerability detection labels as output. Our API provides end-to-end usage for the defender to update/deploy his own DNN and inspect unknown contract bytecode. We detail the implementation procedures as follows.

7.1.1 Dataset Chunking. To enable our supervised learning, a designated dataset with the preprocessed bytecode and associated labels is required. The labeled bytecode dataset for the model training process is constructed as explained in Section 6. Furthermore, we develop a Python module to read data from the database and generate CSV files as the inputs to ESCORT’s DNN model. Since the size of our bytecode dataset is large, loading the full dataset into memory at one time is likely to incur out of memory (OOM) error. To alleviate the memory constraint (Challenge (C1) in Section 4.2), we perform *data chunking* as follows: (i) The bytecode samples and the corresponding labels in the CSV file are merged into one *DataFrame*; (ii) The dataset is then shuffled and stored as a temporary file; and (iii) We create chunks from the shuffled dataset and store the segments as chunk files based on the pre-defined chunk size (default is 1024). The chunk files are passed directly to the DNN model. We provide more details about our data chunking process in Appendix Section A.2.

Note that the collected contract data might have the class imbalance issue [87] (Challenge (C3) in Section 4.2). In our work, we construct a balanced training set to ensure that each batch of data fed into ESCORT’s DNN model has a comparable number of vulnerable and safe contracts. Details about our dataset balancing is given in Section 8.1.

7.1.2 Model Building and Training. To realize ESCORT’s DNN-based vulnerability scanning approach, we instantiate a Multi-Output-Layer (MOL) DNN model based on our generic architecture explained in Section 5.2.

Figure 6 visualizes the actual model implemented in our experiments. We first build a MOL-DNN model with six output branches for main model training and then extend it with two new branches for transfer learning. Recall that our objective is to enable concurrent detection of multiple vulnerability types. To this end, we construct a recurrent neural network by stacking embedding layer, GRU layers, Dropout layers, and Dense layers to process sequential bytecode inputs. The feature extractor is concatenated with multiple branches to output the probability of having each specific vulnerability. After building the MOL-DNN model, we specify the model hyper-parameters as shown in Table 2, and launch the training process as depicted in Figure 3. Since the bytecode dataset is chunked to avoid OOM errors, ESCORT’s DNN model iterates over all chunks multiple times during model training.

We define ‘#Global_Epochs’ to be the total number of times that the full training set is iterated by our DNN model and ‘#Local_Epochs’ as the number of times that each chunk is used to update the model. Before the chunked data are passed to the model’s input layer, the bytecode sequence needs to be vectorized. This is realized by a *tokenizer*, which transforms the hexadecimal data into

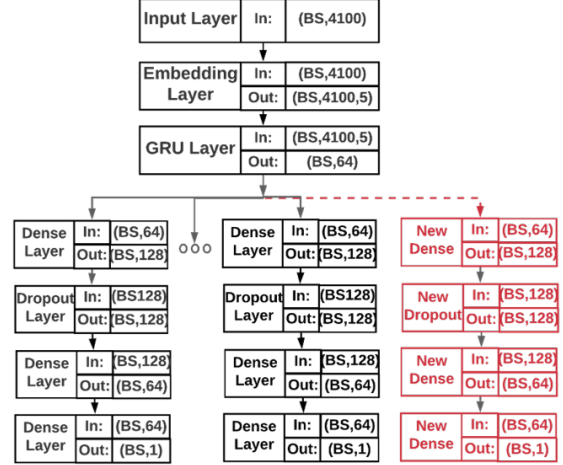


Figure 6: Multi-output DNN architecture of ESCORT for concurrent detection of multiple vulnerability types. Here, *BS* is the batch size, the rightmost red branch denotes adding a new branch for transfer learning on new vulnerabilities.

Variable	Setting
Layer Type	Embedding, GRU, Dense, Dropout
#Hidden Unites	GRU:64, Dense:[128,64,1]
Optimizer	Adam
Loss Function	Binary Cross-Entropy
Learning Rate	0.001
Dropout Values	0.2
#Local_Epochs	1
#Global_Epochs	1
Batch Size	32
MAX Seq. Length	4100

Table 2: Model Hyper-parameters.

numeric vectors. After tokenization, a hyper-parameter *MAX_SEQUENCE_LENGTH* is applied to the input vectors. Sequences shorter than this length are zero-padded, while sequences longer than this length are truncated. We empirically study the distribution of the bytecode length and show the results in Figure 7. The hyperparameter *MAX_SEQUENCE_LENGTH* is set to 4100 to ensure more than 98.5% of the contracts are not truncated¹.

The tokenized data are then passed to the DNN model for training. We assess the classification accuracy of the model after each Local Epoch (L_Ep., train on one chunk) and each Global Epoch (G_Ep., train on all chunks). The evaluation results are stored in a metric history object to keep track of the progress. In the testing phase, the remaining (unseen) data chunks are passed to the model to compute the detection metrics. At the end of the model training stage, we save the converged MOL-DNN model, the used tokenizer, and the evaluation metrics files to wrap the model as an API service.

¹Both excessive padding and truncation may worsen performance. For instance, truncated contracts may end up being mislabeled as benign if a vulnerability resides in the truncated part of the contract.

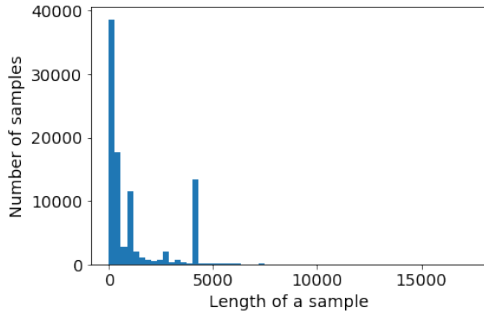


Figure 7: Distribution of bytecodes length in the dataset.

7.2 Model Serving API

Our trained DNN model can detect pre-defined vulnerability types in smart contracts. We wrap the model within an API to ensure that we can serve predictions to end-users on the fly. We utilize Flask [7] and devise a Python module to provide a REST API endpoint for running model inference on bytecode files. Our API also performs automated bytecode transformation (input preprocessing) to remove the manual efforts for the defender. The learned ESCORT model and associated configurations are passed to our API. The Python module creates two different API endpoints. The first endpoint shows the configuration passed to the module and the second one triggers model inference. Listing 1 shows an example of the plain bytecode of the smart contracts, which is passed to the second endpoint for vulnerability detection.

Listing 1: Sample request body when calling ESCORT’s prediction endpoint.

```
1 { "smart_contract": "606060405236150100000..." }
```

In the second endpoint, the bytecode of the input smart contract is vectorized using the same tokenizer as the one used in ESCORT’s model training step. The processed sequence is then fed as the input to the trained MOL-DNN model for vulnerability detection. In addition to the detected vulnerability types, the prediction time of ESCORT is tracked and shown to the user. An example response from ESCORT’s API endpoint is shown in Listing 2.

Listing 2: ESCORT’s response to the sample request when calling our prediction endpoint. The response includes the analysis results of multiple vulnerability types.

```
1 { "prediction": {
2   "ASSERT_VIOLATION": 0.0001,
3   "ACCESSIBLE_SELFDESTRUCT": 0.9998,
4   "DoS (UNBOUNDED_OP)": 0.9996,
5   "MULTIPLE_SENDS": 0.0012,
6   "TAINTED_SELFDESTRUCT": 0.9998,
7   "CALLSTACK": 0.9995,
8   "MONEY_CONCURRENCY": 0.0013,
9   "REENTRANCY": 0.0009},
10  "prediction_time_in_second": "0.02" }
```

8 EVALUATION

We assess ESCORT on the large-scale smart contract dataset built as described in Section 6. In this section, we explain our experimental setup and the evaluation metrics to characterize the performance of ESCORT’s DNN model.

8.1 Dataset

To build our dataset, we collected ~ 1.2 million smart contracts from the first 5 million Ethereum blockchain blocks using ContractScraper, as explained in Section 6.5 and label them accordingly. Note that vulnerability scanners can normally detect multiple vulnerability types, thus there might also be overlap in coverage. For instance, a reentrancy bug is detected by all three tools, Oyente, Mythril, and Dedaub.

To select only one label for redundantly labeled vulnerability types, we first compute respective F1 scores of different tools [12] based on the true positives, false positives, and false negatives reported in the papers. The tool with the highest F1 score on this particular vulnerability class is then used to determine if the collected smart contracts are vulnerable to this vulnerability class. For instance, Oyente [17] is used to detect the reentrancy bug since it yields the highest score.

For the purpose of evaluation, we include eight vulnerability types in our dataset. However, as discussed in Section 5.3, we stress that our approach is not limited only to eight classes and can be easily extended with new attacks.

In the following, we define eight vulnerability types included in our evaluation and refer to vulnerability categories presented in Section 2.1 for more detailed vulnerability description.

- **Callstack Depth [cl. 1]:** This vulnerability class belongs to the category **Programming Error** (cf. Section 2.1) and exploits the stack size limit issues of the EVM.
- **Reentrancy [cl. 2]:** Reentrancy bugs are caused by **External Calls** and allow an attacker to drain funds, as explained in Section 2.1.
- **Multiple Sends [cl. 3]:** This class hinges on the exploitation of a smart contract’s **Execution Costs** to induce DoS.
- **Accessible selfdestruct [cl. 4]:** This **Programming Error** can be exploited to terminate a contract such that the remaining funds are sent to a predefined address.
- **DoS (Unbounded Operation) [cl. 5]:** An attacker can exploit the limited **Execution Costs** of a smart contract when the execution time is dependent on input from an external caller.
- **Tainted selfdestruct [cl. 6]:** This vulnerability is an extension of vulnerability class 4 (cl. 4). The attacker here can set to which address the remaining balance of the smart contract is sent.
- **Money concurrency [cl. 7]:** This vulnerability is also known as Transaction Ordering Dependence (TOD) and belongs to the **Influence by Miners** category.
- **Assert violation [cl. 8]:** This **Programming Error** leads to a constant error state of the smart contract, which can be exploited by an attacker.

For each vulnerability class, we select at least 15,000 samples with this specific vulnerability from our raw dataset (~ 1.2 million contracts) and concatenate them to construct the vulnerable contract set with the equally-sized distribution. We empirically set the minimal sample number to 15,000 since this is the smallest size of the well-represented vulnerability types in our dataset. The above dataset was then completed with 15,000 completely clean smart contracts where no vulnerabilities were detected by the tools

used in this work (Section 6.5). Figure 8 shows the vulnerability class distribution of our balanced, labeled dataset. One can see from Figure 8 that our dataset construction method solves the challenge (C3) formulated in Section 4.2.

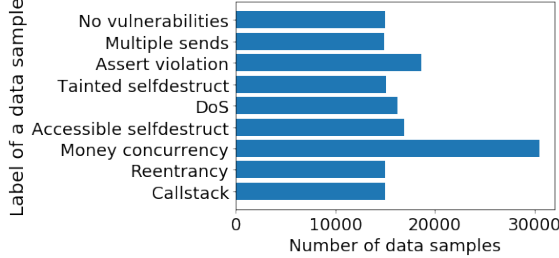


Figure 8: Vulnerability class distribution of our dataset.

It is worth mentioning that the actual size of our labeled dataset is 93,497 samples instead of $15000 \times (8 + 1)$. This is because ESCORT formulates vulnerability detection as *multi-label classification*, meaning that a contract might have multiple labels and repeatedly appear in the selection of several vulnerability types described above. We take 80% and 20% of the labeled dataset to construct the training and test dataset for ESCORT. Note that we set aside 10% of the training data as the validation data to prevent model over-fitting.

8.2 Evaluation Metrics

ESCORT provides concurrent detection of multiple vulnerability types. We evaluate the performance of ESCORT’s multi-output DNN model with F1 score, precision, and recall. In addition, we calculate two of the most common performance metrics, hamming loss, and Jaccard similarity [45]. We detail each metric as follows.

8.2.1 Base Values. The results of true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN) are the base values to compute other metrics. The true values represent the number of correctly predicted results, which can be either true positive or true negative. The false values indicate that the DL model gives the wrong outputs [8].

8.2.2 Precision and Recall. The precision metric describes the ratio of truly positive values to all positive predictions.

This indicates the reliability of the classifier’s positive prediction [71]. The recall (or sensitivity) metric shows the proportion of actual positives that are correctly classified. The formulas to compute these two metrics are given below:

$$\text{Precision} = \frac{TP}{TP + FP}, \text{ Recall} = \frac{TP}{TP + FN}. \quad (1)$$

8.2.3 F1 Score. The F1 score metric is commonly used in information retrieval and it quantifies the overall decision accuracy using precision and recall. The F1 score is defined as the *harmonic mean* of the precision and recall:

$$\text{F1_score} = \frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}. \quad (2)$$

The best and the worst value of the F1 score is 1 and 0, respectively. The F1 score can be calculated for each class label or globally [20].

In our evaluation, we use the weighted F1 score where the per-class F1 scores are weighted by the number of samples from that class [12].

8.2.4 Jaccard Similarity. In multi-label classification, Jaccard similarity (Jaccard index) is defined as the size of the intersection divided by the size of the union of two label sets. It is used to compare the set of predicted labels for a sample to the corresponding set of true labels. It ranges from 0 to 1 where 1 is the perfect score. The Jaccard similarity does not consider the correct classification of negatives [71] and can be computed as follows:

$$\text{Jaccard_Similarity} = \frac{TP}{TP + FP + FN}. \quad (3)$$

8.2.5 Hamming Loss. The Hamming loss gives the percentage of wrong labels to the total number of labels. A lower hamming loss indicates the better performance of a model. For an ideal classifier, the hamming loss is 0. In multi-label classification, the hamming loss is defined as the *hamming distance* [69] between the ground-truth label y and the prediction value \hat{y} :

$$\mathcal{L}_{\text{Hamming}} = \frac{\# \text{Mismatch}(y, \hat{y})}{\text{Length}(y)}. \quad (4)$$

where $\text{Length}(y)$ is the total number of vulnerability types.

8.2.6 DNN Loss Function. The loss function is a crucial part of DNN training since the training process aims to minimize the loss for obtaining a high task accuracy. As such, the loss value quantifies how well a classifier performs on the given dataset. In our experiments, we use *Binary Cross-Entropy* (BCE) loss to train ESCORT’s multi-output DNN model. Given the expected value y and the prediction \hat{y} , the BCE loss is computed as:

$$\mathcal{L}_{\text{BCE}}(y, \hat{y}) = -(y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})). \quad (5)$$

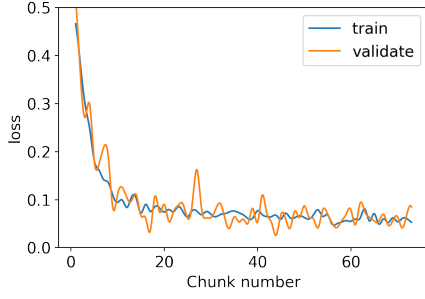
8.3 Experimental Setup

All of our experiments are conducted on a machine with Arch Linux OS having AMD Ryzen 3 3200G and NVIDIA GeForce RTX 3090 GPU with 32 and 24 GB of RAM, respectively. The software versions are as follows: Tensorflow 2.3.1, CUDA 11.1, NVIDIA driver 455.38, cuDNN 8.0.5.39-1, and kernel 5.4.77-1.

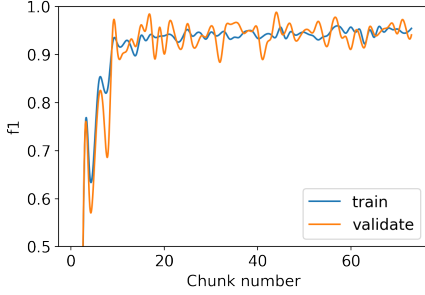
8.4 Evaluation Results

8.4.1 Classifier Learning. We train ESCORT’s model on the training set with hyper-parameters listed in Table 2 and assess it on the test set (consisting of 18,700 contracts). To corroborate the detection effectiveness of ESCORT, we plot the learning curves of our multi-output DNN in Figure 9. The training and the validation curves demonstrate the *time-evolving* performance of ESCORT for supervised learning and the *generalization* capability of the model, respectively. The learning curves show that our DNN model can achieve an average F1 score higher than 95% on both the training and validation set.

8.4.2 Sensitivity to Training Configurations. ESCORT’s training pipeline is described in Section 7.1. We use the MOL-DNN architecture in Figure 6 with six vulnerability branches in ESCORT’s training phase. This model has a total number of 115,846 learnable



(a) ESCORT’s BCE loss on training and validation set.



(b) ESCORT’s F1 score on training and validation set.

Figure 9: The model learning process of ESCORT.

parameters. We investigate whether the detection performance of ESCORT can be improved by increasing the number of global and local epochs. To characterize ESCORT’s performance, we compute the metrics defined in Section 8.2 and summarize the results in Table 3. It can be seen that ESCORT’s accuracy has a slight increase with a longer training time. Since the accuracy improvement is minor w.r.t. the large increase of training time, we use the first configuration (G_Ep=1, L_Ep=1) as our default setting for the rest of our evaluations.

8.4.3 Class-wise Detection Performance. We describe ESCORT’s overall performance on the vulnerability types introduced in Section 8.1. Here, we provide a fine-grained insight into ESCORT’s capability on each vulnerability class. Table 4 shows the class-specific metrics obtained by our MOL-DNN model. ESCORT achieves an average of 97% precision and 95% F1 score across all six vulnerability types. In the case of the Accessible selfdestruct vulnerability (cl. 4), ESCORT’s recall score is relatively lower than others. The low recall value indicates that ESCORT yields more false negatives on this vulnerability class according to Equation (1).

8.5 Extensibility/Generalizability Performance

To corroborate the extensibility of ESCORT, we first train a MOL-DNN model on six vulnerability types (cl. 1-6), and then perform transfer learning on the remaining two vulnerability types (cl. 7-8). The details of the transfer learning procedures can be found in Section 5.3. The two newly added branches (for cl. 7-8) have in total 33,282 trainable parameters. For transfer learning, we used the same hyper-parameters listed in Table 2 and trained the model on vulnerability types 7 and 8 simultaneously. This process takes only 2 hours and 43 minutes, which is 47.39% of the time used by model

Metrics (Avg.)		G_Ep.:1 L_Ep.:1	G_Ep.:1 L_Ep.:3	G_Ep.:3 L_Ep.:1	G_Ep.:3 L_Ep.:3
Training	Loss	0.06	0.05	0.06	0.05
	Precision	0.99	1.00	0.99	0.97
	Recall	0.91	0.90	0.90	0.95
	F1 score	0.95	0.95	0.94	0.96
	Hamming Loss	0.02	0.02	0.02	0.01
	Jaccard Similarity	0.90	0.90	0.89	0.92
Test	Loss	0.06	0.05	0.05	0.05
	Precision	0.98	0.99	0.99	0.97
	Recall	0.90	0.91	0.91	0.95
	F1 score	0.95	0.95	0.95	0.96
	Hamming Loss	0.02	0.02	0.02	0.01
	Jaccard Similarity	0.90	0.90	0.90	0.92
Time	Training (H:M)	5:44	17:11	17:40	49:41
	Prediction (Sec.)	0.02	0.02	0.02	0.02

Table 3: The averaged performance of our multi-label classification model across initial vulnerability types (cl. 1-6) with different G_Ep. and L_Ep. configurations.

Metrics	Initial vulnerability types						New Vuln. Classes	
	cl. 1	cl. 2	cl. 3	cl. 4	cl. 5	cl. 6	cl. 7	cl. 8
Loss	0.08	0.05	0.07	0.09	0.01	0.09	0.14	0.08
Precision	0.99	0.97	0.99	0.99	1.00	0.98	0.95	0.98
Recall	0.88	0.96	0.91	0.85	0.99	0.88	0.90	0.90
F1 score	0.93	0.96	0.95	0.91	0.99	0.93	0.92	0.93
FPR	0.00	0.01	0.00	0.00	0.00	0.00	0.02	0.00
FNR	0.12	0.06	0.10	0.15	0.01	0.12	0.10	0.10

Table 4: Class-specific metrics (for unseen/test data) retrieved by our multi-output model for initial training phase (cl. 1-6) and after transfer learning phase (cl. 7-8).

training with six vulnerabilities. The detection result on the new vulnerability types are shown in the last two columns of Table 4. It can be seen that ESCORT achieves 92% and 93% F1 score on the two new vulnerabilities cl. 7 and cl. 8, respectively. The empirical results corroborate that ESCORT is extensible to new attacks by enabling lightweight and effective transfer learning.

We also demonstrate the superior efficiency and effectiveness of ESCORT’s innovative multi-output architecture compared to the existing ML-based detection techniques. In this comparison experiment, we define a regular RNN consisting of the three top layers and two branches shown in Figure 6 as the baseline model. This new model is trained *from scratch* on the two new vulnerability types. For ESCORT, we use the pre-trained feature extractor obtained from the main training stage and append two new branches of layers to it. As we explain in Section 5.3, only the parameters in the new vulnerability branches are trained during transfer learning.

Note that the size of the transfer learning dataset is considerably smaller than the data size in the main training phase since the new smart contracts are collected in a shorted time period. The

limited training data raises concerns about model *under-fitting* if the amount of trainable parameters is too large. We show the empirical comparison results in Table 5 when adapting the ML model to the new vulnerability types (cl. 7 and cl. 8). It can be seen that ESCORT reduces the training time by $\sim 43\%$ while yielding a similar F1 score compared to the baseline. The advantage of ESCORT is derived from our multi-output RNN design which decomposed the contract vulnerability problem into two sub-tasks: contract representation learning (performed by feature extractor), and vulnerability pattern identification (performed by each branch). As such, when extending the ML-based detector to new vulnerability types, ESCORT can focus on distinguishing the new vulnerability pattern without ‘*re-learning*’ the contract representation.

Transfer Learning	FPR	FNR	F1 score	Training Time (H:M)
ESCORT	0.02	0.10	0.93	2:43
Baseline	0.01	0.10	0.94	4:46

Table 5: Performance comparison of transfer learning between ESCORT and baseline (training-from-scratch).

9 RELATED WORK

Various vulnerability inspection tools have been developed for smart contracts to ensure the security of the cryptocurrency system. We categorize the existing detection techniques based on their working mechanism and discuss each type below.

9.1 Static Vulnerability Detection Methods

Static detection techniques analyze the smart contract in a static environment by examining its source code or bytecode.

9.1.1 Information Flow Analysis-based. Slither [39] uses *taint analysis* [83] to detect vulnerabilities in Solidity source code. It can find nearly all vulnerabilities related to the user inputs or critical data flows while the inspection time might be prohibitively long. Dedaub’s Contract-Library by Dedaub [3] provides multiple different features via an online API. It collects bytecode of the smart contracts and performs vulnerability classifications using tool MadMax [43] that performs flow and loop analysis to detect gas-focused vulnerabilities [31].

9.1.2 Symbolic Execution-based. Oyente [59] detects vulnerabilities in the source code or bytecode of Solidity contracts using *symbolic execution*. Symbolic execution represents the program’s behavior as built formula and uses symbolic inputs to decide if a certain path can be reached [53]. As such, its performance depends on the number of explored paths and the program’s complexity [34, 77]. Oyente constructs the control flow graph of the contract and uses it to create inputs for symbolic execution. Manticore [61] analyzes the contract by repeatedly executing symbolic transactions against the bytecode, tracking the discovered states, and verifying code invariants [11]. Securify [84] first obtains the contract’s semantic information by performing symbolic analysis of the dependency graph, then checks the predefined compliance and violation patterns for vulnerability detection. teEther [55] searches for certain critical paths in the control flow graph of the smart contract and uses symbolic execution for vulnerability identification.

9.1.3 Logic Rules-based. Vandal [32] is a logic-driven static program analysis framework. It converts the low-level EVM bytecode to semantic logic relations and describes the security analysis problems with logic rules. The datalog engine executes the specifications for input relations and outputs the vulnerabilities. eThor [76] is a static analysis technique built on top of reachability analysis achieved by Horn clause resolution. NeuCheck [58] adopts a syntax tree in a syntactical analyzer to transform source code of smart contracts to an intermediate representation (IR). Vulnerabilities are identified by searching for detection patterns in the syntax tree. SmartCheck [81] converts the Solidity source code to XML-based IR and verifies it against detection patterns defined in XPath language [27].

9.1.4 Composite Methods-based. Mythril [63] combines multiple vulnerability detection approaches, including symbolic execution, taint analysis, and Satisfiability Modulo Theories (SMT). SMT solving converts the contract to SMT constraints to reveal program flaws. Zeus [51] uses symbolic model checking, abstract interpretation, and constrained horn clauses to verify contracts’ security. Osiris [82] combines symbolic execution and taint analysis to precisely identify integer bugs in smart contracts.

9.2 Dynamic Vulnerability Detection Methods

Dynamic testing techniques execute the program and observe its behaviors to determine the vulnerability’s existence.

9.2.1 Fuzzing-based. MythX [15] combines synthetic execution and *code fuzzing*. It provides a cloud-based API for developers to inspect smart contracts. Fuzzing [41] is a testing method that attempts to expose the vulnerabilities by executing the program with invalid, unexpected, or random inputs. The brute-force nature determines that fuzzing incurs large runtime overhead and might have poor code coverage due to its dependency on the inputs [34]. ReGuard [57] is another fuzzing tool specialized in the Reentrancy bug. It creates an IR for the smart contract. A fuzzing engine is used to generate random byte inputs and analyze the execution traces for reentrancy bugs detection. ContractFuzzer [50] generates fuzzing inputs based on the ABI specifications of smart contracts. Test oracles are defined to monitor and analyze the contract’s runtime behaviors for vulnerability detection. Echidna [44] is a fuzzer that generates random tests to detect violations in assertions and custom properties. ILF [46] uses symbolic execution to generate contract inputs and employs imitation learning to design a neural network-based fuzzer from symbolic execution. sFuzz [65] is an adaptive fuzzer for smart contracts that combines the AFL fuzzer and multi-objective strategy to explore hard-to-cover branches. Harvey [89] is a greybox fuzzer that predicts new inputs to cover new paths and fuzzes the transaction sequence in a demand-driven manner.

9.2.2 Validation-based. ContractLarva [70] is a runtime verification tool for smart contracts where a violation of defined properties can lead to various handling strategies, such as a system stop. These properties can include undesired event traces of control or data flow. Maian [66, 67] combines symbolic analysis and concrete validation to inspect the smart contract’s bytecode. In concrete validation, the contract is executed on a fork of Ethereum for tracing and validation. By passing symbolic inputs to the contract, the execution

trace is analyzed to identify the vulnerabilities. Sereum [73] uses runtime monitoring and verification to protect existing smart contracts against reentrancy attacks without modifications or semantic knowledge of the contracts. It detects inconsistent states in the contract via dynamic taint tracking and data flow monitoring during contract execution.

9.3 Machine Learning for Vulnerability Detection

Several works have attempted to perform automated contract scanning using machine learning techniques. We discuss their working mechanisms and limitations below.

ContractWard. ContractWard detects smart contracts vulnerability in the *opcode-level* by extracting *bigram* features from the simplified opcode and training individual binary ML classifiers for each vulnerability class [87]. The paper targets six vulnerabilities and experiments with Random Forests, K-Nearest Neighbors, SVM, AdaBoost, and XGBoost classifiers.

Compare to our work, ContractWard has three main limitations: (i) *Requires source code of smart contracts.* This approach analyzes smart contracts with opcodes. To do so, it decompiles source codes and converts them to opcodes. It is worth mentioning that decompilation might result in information loss. (ii) *Not extensible to new exploitation attacks.* This paper uses ‘One vs. Rest’ algorithms and designs separate ML models to detect each vulnerability class. This means that supporting a new vulnerability class requires training a new ML model from scratch, which is costly. (iii) *Not scalable to long contracts.* The *bigram* language model has a short window size in the Markov chain model. This determines that ContractWard is not scalable to long contracts and cannot capture long-term dependency in the code. However, using an *n*-gram model with a larger window size increases the feature size, thus complicating model training due to the high dimensionality of data.

LSTM-based. The paper [80] proposes a sequence learning approach to detect weakness in the opcode of smart contracts. Particularly, this paper uses one-hot encoding and an embedding matrix to represent the contract’s opcode. The obtained code vectors are used as input to train an LSTM model for determining whether the given smart contract is safe or vulnerable (i.e., binary classification).

The LSTM-based scheme yields limited detection performance since: (i) The reported F1 score of 86% is relatively low. We hypothesize that this might be because different vulnerability types have diverse behaviors, thus making it hard to distinguish the group of multiple vulnerabilities from the safe contracts. (ii) The LSTM model only provides a binary decision about contract security without distinguishing vulnerability types.

AWD-LSTM based. A sequence-based multi-class classification scheme is presented in [42]. This paper adapts ‘Average Stochastic Gradient Descent Weighted Dropped LSTM’ (AWD-LSTM) [60] for vulnerability detection. The proposed model consists of two parts: a pre-trained encoder for language tasks [47], and an LSTM-based classifier for vulnerability classification. This method works on the opcode-level and can detect three vulnerability types.

Compared to ESCORT, the AWD-LSTM based detection method has the following constraints: (i) *Non-uniform effectiveness.* The multi-class detection performance of [42] is not uniformly effective across different vulnerabilities. In particular, [42] yields an F1

score of 95% on safe contracts and 30% on Prodigal contracts [68]. ESCORT features a much smaller performance divergence across different classes as can be seen from Table 4. (ii) *Not extensible.* The AWD-LSTM based model in [42] is a fixed design to detect pre-specified vulnerability types. The extension to incorporate new attacks is not considered in [42].

CNN-based. The paper [48] transforms the contract bytecode into fix-sized RGB color images and trains a convolution neural network for vulnerability detection.

Similarly to ESCORT, CNN-based classifier uses multi-label classification, which has a low confidence score when determining the exact vulnerability types.

Compare to our work, the CNN-based detection scheme has the following limitations: (i) The multi-label classification performance is not satisfying due to its low confidence level. We hypothesize that this is because image representation of the bytecode and the CNN architecture ignore the sequential information existing in the contract. (ii) The extensibility/generalization ability of the CNN-based detection method is neither discussed nor evaluated.

GNN-based. The paper [92] proposes a graph neural network (GNN)-based approach. In particular, this work builds a ‘contract graph’ from the contract’s source code where nodes and edges represent critical function calls/variables and temporal execution trace, respectively. This graph is normalized to highlight important nodes and passed to a temporal message propagation (TMP) network for vulnerability detection.

While supporting multi-class detection, the GNN-based method has the following drawbacks: (i) *Restricted applicability.* The paper [92] requires to build a graph from the contract’s source code. However, the source code is typically hard to obtain from the public blockchain. ESCORT operates on the contract bytecode which is publicly available from the EVM, thus can be deployed in more scenarios. (ii) *Limited effectiveness.* The GNN-based technique yields an average F1 score of 77% across all three vulnerabilities. We hypothesize that the graph normalization process in [92] does not preserve the malicious nodes responsible for vulnerability exploitation, leading to the low F1 score.

10 CONCLUSION

To ensure the safety of the Ethereum cryptocurrency system, we present ESCORT, the first deep learning-based automated framework that supports concurrent detection of multiple vulnerability classes and lightweight transfer learning. We identify two key components of vulnerability detection: feature extraction of general smart contracts and each particular vulnerability class, and disentangle these two subtasks. ESCORT’s multi-output RNN design is highly modular, scalable, efficient, and extensible as opposed to the previous works. Empirical results show that ESCORT achieves an average of 95% detection accuracy in terms of F1 score across various vulnerability classes and can be quickly adapted on the new vulnerability data. Given an unknown smart contract, ESCORT can provide parallel detection of eight vulnerabilities in 0.02 second. As a separate contribution, we devise ContractScraper, a toolchain for dataset construction and labeling based on smart contracts’ bytecode downloaded from Ethereum blockchain and existing vulnerability detection tools. We will open source ContractScraper toolchain and our dataset to promote research in this area.

A APPENDIX

In this section, we provide a detailed introduction about cryptocurrency systems and illustrate how we perform dataset chunking for ESCORT’s training.

A.1 Ethereum Platform

Blockchain is proposed as a distributed ledger that records transactions between two parties in a verifiable and permanent way [91]. Ethereum is an open-sourced cryptocurrency platform based on blockchain and provides a Turing-complete Ethereum Virtual Machine (EVM) that enables developers to deploy decentralized applications. A cryptocurrency platform has the following key characteristics:

Decentralized Nature. In contrast to conventional currencies, virtual money is not administered by a central authority but by a distributed peer-to-peer network. A network of nodes, the so-called ‘miners’, are responsible to perform money transactions, data storage, and updates [6]. Note that on a blockchain, all code, data, and transactions are shared and available for inspection on every single node. All actions performed inside of this network need to be confirmed by the majority of all participating nodes [75].

Mathematical Algorithm as a Basis of Cryptocurrency Value. Money, in the Ethereum context called Ether, can be initially earned by solving a complex mathematical problem that can be accepted by the other nodes. The process is called *mining*. Once a transaction is triggered, a state value for the new block is calculated and verified by the participating nodes in the network [6]. The amount of money available in the network is also ensured to be limited. As such, the Ether and the associated owners can be tracked at all times [38].

Resilience to Data Manipulations from Outside. The information of the mined money is stored in a public data structure, i.e., the blockchain. More specifically, modifications and transactions are stored in blocks that are appended to the chain in a chronological order after checking its validity [59]. Afterward, the network rejects any attempts to alter the blockchain entries. Therefore, the data is immutable and irreversible [38].

Pseudonymous Nature. In general, registration is not required to use cryptocurrency. Users that perform transactions inside the network are identified by a public key and a private key. All transactions are associated with the addresses instead of explicit users [6]. Therefore, it is hard to determine the identity of a user although all transactions are stored publicly on the blockchain.

A.2 Data Chunking Process

As mentioned in the challenge (C1) of Section 4.2, bytecode files contain long sequences, thus consuming large memory during DNN training. To solve this challenge, the collected bytecode data needs to be partitioned. Note that the pre-labeling process discussed in Section 6.5 is based on the smart contract as a whole. It is not guaranteed that vulnerabilities will occur in the chunked segments of a smart contract since direct partition of the contract removes the inherent dependency traces of potential attacks. As such, we decide to chunk the dataset constructed in Section 6 while keeping the bytecode of each contract complete to relax the memory usage of ESCORT’s DNN model.

We describe how we create data chunks from the labeled bytecode data in Section 7.1. Figure 10 shows how DNN training is performed on the data chunks. Two hyper-parameters are specified for model training purpose: *Local_Epoch* and *Global_Epoch*. As shown in the figure, *Local_Epoch* defines how often each chunk is used for training before switching to the next chunk, *Global_Epoch* defines how often the union of all the data chunks should be iterated over during the entire training process.

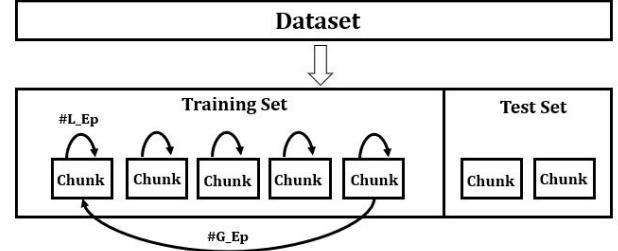


Figure 10: Illustration of our dataset chunking method.

REFERENCES

- [1] Accessed 2020. Ethereum API | IPFS API Gateway | ETH Nodes as a Service. <https://infura.io/>
- [2] Accessed 2020. Ethereum Contract Library by Dedaub. <https://contract-library.com/>
- [3] Accessed 2020. Ethereum Contract Library by Dedaub. <https://contract-library.com/>
- [4] Accessed 2020. Ethereum ETL. <https://github.com/blockchain-etl/ethereum-etl>
- [5] Accessed 2020. Ethereum in BigQuery: How we built this dataset. <https://cloud.google.com/blog/>
- [6] Accessed 2020. Ethereum Whitepaper. <https://ethereum.org>
- [7] Accessed 2020. Flask API. <https://www.flaskapi.org/>
- [8] Accessed 2020. Google Machine Learning Glossary. <https://developers.google.com/machine-learning/glossary>
- [9] Accessed 2020. Hyperledger Project. <https://www.hyperledger.org/>
- [10] Accessed 2020. Keras: Multiple outputs and multiple losses. <https://www.pyimagesearch.com/2018/06/04/keras-multiple-outputs-and-multiple-losses/>
- [11] Accessed 2020. Manticore. <https://github.com/trailofbits/manticore>
- [12] Accessed 2020. Multi-Class Metrics Made Simple, Part II: the F1-score. <https://towardsdatascience.com/multi-class-metrics-made-simple-part-ii-the-f1-score-ebe8b2c2ca1>
- [13] Accessed 2020. MySQL. <https://www.mysql.com/>
- [14] Accessed 2020. Mythril. <https://github.com/ConsenSys/mythril>
- [15] Accessed 2020. MythX: Smart contract security service for Ethereum. <https://mythx.io/>
- [16] Accessed 2020. Overview · Smart Contract Weakness Classification and Test Cases. <http://swcregistry.io/>
- [17] Accessed 2020. Oyente: An analysis tool for smart contracts. <https://github.com/melonproject/oyente>
- [18] Accessed 2020. Quorum. <https://consensys.net/quorum/>
- [19] Accessed 2020. Rootstock (RSK) platform for smart contracts. <https://en.bitcoinwiki.org/wiki/Rootstock>
- [20] Accessed 2020. scikit-learn: F-measure. https://scikit-learn.org/stable/modules/generated/sklearn.metrics.f1_score.html
- [21] Accessed 2020. Serpent assembly programming language. <https://github.com/ethereum/serpent>
- [22] Accessed 2020. The Solidity Contract-Oriented Programming Language. <https://github.com/ethereum/solidity>
- [23] Accessed 2020. TensorFlow. <https://www.tensorflow.org/>
- [24] Accessed 2020. Tron: Decentralize the web. <https://tron.network/>
- [25] Accessed 2020. Vechain. <https://www.vechain.com/>
- [26] Accessed 2020. Vyper documentation. <https://vyper.readthedocs.io/en/stable/>
- [27] Accessed 2020. XML Path Language (XPath). <https://www.w3.org/TR/xpath20/>
- [28] Ethem Alpaydin. 2020. *Introduction to machine learning*. MIT press.
- [29] Leonardo Alt and Christian Reitwießner. 2018. SMT-Based Verification of Solidity Smart Contracts. In *International Symposium on Leveraging Applications of Formal Methods*.
- [30] Nicola Atzei, Massimo Bartoletti, and Tiziana Cimoli. 2017. A Survey of Attacks on Ethereum Smart Contracts (SoK). In *International Conference on Principles of Security and Trust*.

- [31] Lexi Brent, Neville Grech, Sifis Lagouvardos, Bernhard Scholz, and Yannis Smaragdakis. 2020. Ethainter: A Smart Contract Security Analyzer for Composite Vulnerabilities. In *PLDI*.
- [32] Lexi Brent, Anton Jurisevic, Michael Kong, Eric Liu, Francois Gauthier, Vincent Gramoli, Ralph Holz, and Bernhard Scholz. 2018. Vandal: A scalable security analysis framework for smart contracts. *arXiv preprint arXiv:1809.03981* (2018).
- [33] Peter F Brown, Vincent J Della Pietra, Peter V Desouza, Jennifer C Lai, and Robert L Mercer. 1992. Class-Based n-gram Models of Natural Language. *Computational linguistics* (1992).
- [34] Jun Cai, Shangfei Yang, Jinqian Men, and Jun He. 2014. Automatic software vulnerability detection based on guided deep fuzzing. In *International Conference on Software Engineering and Service Science*.
- [35] Michael del Castillo. Accessed 2020. The DAO Attacked: Code Issue Leads to \$60 Million Ether Theft. <https://www.coindesk.com/dao-attacked-code-issue-leads-60-million-ether-theft>
- [36] Christopher D Clack, Vikram A Bakshi, and Lee Braine. 2016. Smart contract templates: foundations, design landscape and research directions. *arXiv preprint arXiv:1608.00771* (2016).
- [37] Thomas Cook, Alex Latham, and Jae Hyung Lee. 2017. DappGuard : Active Monitoring and Defense for Solidity Smart Contracts. (2017).
- [38] Kevin Delmolino, Mitchell Arnett, Ahmed Kosba, Andrew Miller, and Elaine Shi. 2016. Step by step towards creating a safe smart contract: Lessons and insights from a cryptocurrency lab. In *International Conference on Financial Cryptography and Data Security*.
- [39] Josselin Feist, Gustavo Grieco, and Alex Groce. 2019. Slither: A Static Analysis Framework For Smart Contracts. In *2019 IEEE/ACM 2nd International Workshop on Emerging Trends in Software Engineering for Blockchain (WETSEB)*.
- [40] Ying Fu, Meng Ren, Fuchen Ma, Heyuan Shi, Xin Yang, Yu Jiang, Huizhong Li, and Xiang Shi. 2019. EVMFuzzer: detect EVM vulnerabilities via fuzz testing. In *ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*.
- [41] Patrice Godefroid, Michael Y Levin, and David Molnar. 2012. SAGE: whitebox fuzzing for security testing. *Queue* (2012).
- [42] Ajay K Gogineni, S Swayamjyoti, Devadatta Sahoo, Kisor K Sahu, et al. 2020. Multi-Class classification of vulnerabilities in Smart Contracts using AWD-LSTM, with pre-trained encoder inspired from natural language processing. *arXiv preprint arXiv:2004.00362* (2020).
- [43] Neville Grech, Michael Kong, Anton Jurisevic, Lexi Brent, Bernhard Scholz, and Yannis Smaragdakis. 2018. Madmax: Surviving out-of-gas conditions in Ethereum smart contracts. *ACM on Programming Languages* (2018).
- [44] Gustavo Grieco, Will Song, Artur Cygan, Josselin Feist, and Alex Groce. 2020. Echidna: Effective, Usable, and Fast Fuzzing for Smart Contracts. In *SIGSOFT International Symposium on Software Testing and Analysis*.
- [45] Gavin Hacking. 2017. *Mastering Machine Learning with scikit-learn*. Packt Publishing Ltd.
- [46] Jingxuan He, Mislav Balunović, Nodar Ambroladze, Petar Tsankov, and Martin Vechev. 2019. Learning to Fuzz from Symbolic Execution with Application to Smart Contracts. In *ACM SIGSAC Conference on Computer and Communications Security*.
- [47] Jeremy Howard and Sebastian Ruder. 2018. Universal language model fine-tuning for text classification. *arXiv preprint arXiv:1801.06146* (2018).
- [48] TonTon Hsien-De Huang. 2018. Hunting the Ethereum smart contract: Color-inspired inspection of potential attacks. *arXiv preprint arXiv:1807.01868* (2018).
- [49] Christoph Jentzsch. 2016. Decentralized autonomous organization to automate governance. *White paper, November* (2016).
- [50] Bo Jiang, Ye Liu, and WK Chan. 2018. ContractFuzzer: Fuzzing smart contracts for vulnerability detection. In *2018 33rd IEEE/ACM International Conference on Automated Software Engineering (ASE)*.
- [51] Sukrit Kalra, Seep Goel, Mohan Dhawan, and Subodh Sharma. 2018. ZEUS: Analyzing Safety of Smart Contracts.. In *NDSS*.
- [52] Jaeyoung Kim, Mostafa El-Khamy, and Jungwon Lee. 2017. Residual LSTM: Design of a deep recurrent architecture for distant speech recognition. *arXiv preprint arXiv:1701.03360* (2017).
- [53] James C King. 1976. Symbolic execution and program testing. *Commun. ACM* (1976).
- [54] Shun Kiyono, Jun Suzuki, and Kentaro Inui. 2019. Mixture of expert/imitator networks: Scalable semi-supervised learning framework. In *AAAI Conference on Artificial Intelligence*.
- [55] Johannes Krupp and Christian Rossow. 2018. teEther: Gnawing at Ethereum to Automatically Exploit Smart Contracts. In *27th USENIX Security Symposium (USENIX Security 18)*.
- [56] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. 2015. Deep learning. *Nature* (2015).
- [57] Chao Liu, Han Liu, Zhao Cao, Zhong Chen, Bangdao Chen, and Bill Roscoe. 2018. Regard: finding reentrancy bugs in smart contracts. In *International Conference on Software Engineering: Companion (ICSE-Companion)*.
- [58] Ning Lu, Bin Wang, Yongxin Zhang, Wenbo Shi, and Christian Esposito. 2019. NeuCheck: A more practical Ethereum smart contract security analysis tool. *Software: Practice and Experience* (2019).
- [59] Loi Luu, Duc-Hiep Chu, Hrishi Olickel, Prateek Saxena, and Aquinas Hobor. 2016. Making smart contracts smarter. In *ACM SIGSAC conference on computer and communications security*.
- [60] Stephen Merity, Nitish Shirish Keskar, and Richard Socher. 2017. Regularizing and optimizing LSTM language models. *arXiv preprint arXiv:1708.02182* (2017).
- [61] Mark Mossberg, Felipe Manzano, Eric Hennenfent, Alex Groce, Gustavo Grieco, Josselin Feist, Trent Brunson, and Artem Dinaburg. 2019. Manticore: A User-Friendly Symbolic Execution Framework for Binaries and Smart Contracts. *arXiv preprint arXiv:1907.03890* (2019).
- [62] Bernhard Mueller. 2018. Smashing ethereum smart contracts for fun and real profit. In *9th Annual HITB Security Conference (HITBSecConf)*. 54.
- [63] Bernhard Mueller. Accessed 2020. MythX Tech: Behind the Scenes of Smart Contract Security Analysis. <https://blog.mythx.io/features/mythx-tech-behind-the-scenes-of-smart-contract-analysis/>
- [64] Satoshi Nakamoto. Accessed 2020. Bitcoin: A peer-to-peer electronic cash system. <http://bitcoin.org/bitcoin.pdf>
- [65] Tai D Nguyen, Long H Pham, Jun Sun, Yun Lin, and Quang Tran Minh. 2020. sFuzz: An Efficient Adaptive Fuzzer for Solidity Smart Contracts. *arXiv preprint arXiv:2004.08563* (2020).
- [66] Ivica Nikolic. Accessed 2020. MAIAN. <https://github.com/ivicanikolicsg/MAIAN>
- [67] Ivica Nikolić, Aashish Kolluri, Ilya Sergey, Prateek Saxena, and Aquinas Hobor. 2018. Finding the greedy, prodigal, and suicidal contracts at scale. *Annual Computer Security Applications Conference* (2018).
- [68] Ivica Nikolić, Aashish Kolluri, Ilya Sergey, Prateek Saxena, and Aquinas Hobor. 2018. Finding the greedy, prodigal, and suicidal contracts at scale. In *Annual Computer Security Applications Conference*.
- [69] Mohammad Norouzi, David J Fleet, and Russ R Salakhutdinov. 2012. Hamming distance metric learning. In *Advances in neural information processing systems*.
- [70] Gordon Pace and Joshua Ellul. 2018. Runtime Verification of Ethereum Smart Contracts. In *2018 14th European Dependable Computing Conference (EDCC)*. <https://doi.org/10.1109/EDCC.2018.00036>
- [71] David MW Powers. 2020. Evaluation: from precision, recall and F-measure to ROC, informedness, markedness and correlation. *arXiv preprint arXiv:2010.16061* (2020).
- [72] Rob Price. Accessed 2020. Digital currency Ethereum is cratering because of a \$50 million hack. uk.businessinsider.com/dao-hacked-ethereum-crashing-in-value-tens-of-millions-allegedly-stolen-2016-6
- [73] Michael Rodler, Wenting Li, Ghassan O Karam, and Lucas Davi. 2018. Sereum: Protecting Existing Smart Contracts Against Re-Entrancy Attacks. *arXiv preprint arXiv:1812.05934* (2018).
- [74] Rebecca Russell, Louis Kim, Lei Hamilton, Tomo Lazovich, Jacob Harer, Onur Ozdemir, Paul Ellingwood, and Marc McConley. 2018. Automated vulnerability detection in source code using deep representation learning. In *International Conference on Machine Learning and Applications (ICMLA)*.
- [75] Alexander Savelyev. 2017. Contract law 2.0: ‘Smart’ contracts as the beginning of the end of classic contract law. *Information & Communications Technology Law* (2017).
- [76] Clara Schneidewind, Ilya Grishchenko, Markus Scherer, and Matteo Maffei. 2020. *eThor: Practical and Provably Sound Static Analysis of Ethereum Smart Contracts*.
- [77] Edward J Schwartz, Thanassis Avgerinos, and David Brumley. 2010. All you ever wanted to know about dynamic taint analysis and forward symbolic execution (but might have been afraid to ask). In *2010 IEEE symposium on Security and privacy*.
- [78] David Siegel. Accessed 2020. Understanding The DAO Attack. <https://www.coindesk.com/understanding-dao-hack-journalists>
- [79] Matt Suiche. Accessed 2020. The \$280M Ethereum’s Parity bug. <https://blog.comae.io/the-280m-ethereums-bug-f28e5de43513>
- [80] Wesley Joon-Wie Tann, Xing Jie Han, Sourav Sen Gupta, and Yew-Soon Ong. 2018. Towards safer smart contracts: A sequence learning approach to detecting security threats. *arXiv preprint arXiv:1811.06632* (2018).
- [81] Sergei Tikhomirov, Ekaterina Voskresenskaya, Ivan Ivanitskiy, Ramil Takhaviev, Evgeny Marchenko, and Yaroslav Alexandrov. 2018. Smartcheck: Static analysis of Ethereum smart contracts. In *International Workshop on Emerging Trends in Software Engineering for Blockchain*.
- [82] Christof Ferreira Torres, Julian Schütte, and Radu State. 2018. Osiris: Hunting for integer bugs in Ethereum smart contracts. In *Annual Computer Security Applications Conference*.
- [83] Omer Tripp, Marco Pistoia, Stephen J Fink, Manu Sridharan, and Omri Weisman. 2009. TAJ: effective taint analysis of web applications. *ACM Sigplan Notices* (2009).
- [84] Petar Tsankov, Andrei Dan, Dana Drachler-Cohen, Arthur Gervais, Florian Buenzli, and Martin Vechev. 2018. Securify: Practical security analysis of smart contracts. In *ACM SIGSAC Conference on Computer and Communications Security*.
- [85] Petr Vytovtov and Kirill Chuvilin. [n.d.]. Prediction of Common Weakness Probability in C/C++ Source Code Using Recurrent Neural Networks.
- [86] Hanna M Wallach. 2006. Topic modeling: beyond bag-of-words. In *International Conference on Machine Learning*.

- [87] Wei Wang, Jingjing Song, Guangquan Xu, Yidong Li, Hao Wang, and Chunhua Su. 2020. Contractward: Automated vulnerability detection models for Ethereum smart contracts. *IEEE Transactions on Network Science and Engineering* (2020).
- [88] Gavin Wood et al. 2014. Ethereum: A secure decentralised generalised transaction ledger. *Ethereum project yellow paper* (2014).
- [89] Valentin Wüstholz and Maria Christakis. 2020. Harvey: A greybox fuzzer for smart contracts. In *ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*.
- [90] Hamed Zamani and W Bruce Croft. 2017. Relevance-based word embedding. In *International ACM SIGIR Conference on Research and Development in Information Retrieval*.
- [91] Zibin Zheng, Shaoan Xie, Hongning Dai, Xiangping Chen, and Huaimin Wang. 2017. An overview of blockchain technology: Architecture, consensus, and future trends. In *2017 IEEE international congress on big data (BigData congress)*.
- [92] Yuan Zhuang, Zhenguang Liu, Peng Qian, Qi Liu, Xiang Wang, and Qinming He. 2020. Smart Contract Vulnerability Detection Using Graph Neural Networks. (2020). <https://doi.org/10.24963/ijcai.2020/450>