

A multi-class Function/Line Level Vulnerability Detection using Graph Neural Networks

Hamidreza M. Taheri, Alireza Shafieinejad*

Department of Electrical and Computer Engineering, Tarbiat Modares University, Jalal AleAhmad Nasr, P.O.Box 14115-111, Tehran, Iran

ARTICLE INFO

Keywords:

Vulnerability Detection
Line Level Vulnerability Detection
Multi-granularity Detection
Static Analysis
Deep Learning
Graph Neural Network

ABSTRACT

One of the challenging issues for software developers is detecting vulnerabilities at different development stages. Security researchers are always seeking new methods to detect vulnerabilities more precisely in a short amount of time. While there are many static and dynamic methods for detecting and discovering vulnerabilities, many of these approaches come with a high computational cost, leading to inefficiencies, particularly in large codebases. Recently, deep learning has gained prominence in extracting vulnerability features from code without the need for a cybersecurity expert.

In this paper, we propose a multi-class Vulnerability Detection scheme at both the Line Level (LLVD) and Function Level (FLVD) using graph neural networks based on node-level and graph-level prediction models, respectively. Moreover, by combining LLVD as a fine-grained method with FLVD as a coarse one, we propose a multi-granularity scheme called File/Line Level Vulnerability Detection (FLLVD) scheme. More specifically, it uses FLVD to detect the type of vulnerability while employing LLVD to identify its location in the source code. Our scheme's variants work with any abstraction graph extracted from incoming source code, such as Data Dependency Graph (DDG) and Program Dependency Graph (PDG).

We evaluate our schemes using both man-made and real-world datasets : SARD and BigVul. Particularly, LLVD and FLLVD achieve performance gains of 0.90 and 0.94, respectively, in terms of F_1 metrics for a subset of SARD with 20 vulnerability types. In contrast, for the combination of SARD and BigVul with 6 vulnerability types, LLVD and FLLVD have F_1 scores of approximately 0.76 and 0.82, respectively.

1. Introduction

Software vulnerabilities are the major source of security flaws and cyber attack. Thus, finding potential software vulnerabilities is an important step in creating defensive layer against cyber attacks. It is difficult and likely infeasible for developers to exactly determine the vulnerable parts of the source code in large-scale software systems. Thus, the main problem is to design an automated tool to detect efficiently and accurately vulnerabilities of a given source code.


Traditional solutions for automated detection are mainly rule-based obtained by human experts. Recently, a lot of research has been dedicated to using data-driven solutions for automated detection tools that take advantage of data mining and machine learning to predict the presence of software vulnerabilities. However, this approach is not straightforward and encounters many challenges in algorithm design and implementation such as selection of machine learning algorithm, proper representation of source code and fitting it to a constant


size vector for feeding to the input layer of learning model.

The most important issues in the state-of-the-art of vulnerability detection using machine learning are as follows:

- **Convolutional deep learning versus Graph Neural Network (GNN) model:** GNNs as part of deep learning technology, are a good choice for training an automatic vulnerability detection model since source code can be efficiently represented as graphs, either by focusing on control flow, data dependency among statements, or both. Recent work like [1–7] use GNN for vulnerability detection while earlier work [8–13] utilize conventional neural networks such as BLSTM (Bidirectional Long-Short Time Memory) for their building blocks.
- **Binary or multi-class classification:** Most of the prior researches focus on identifying a function as vulnerable or not-vulnerable regardless of vulnerability type. For instance, [2, 8, 9, 12, 14] are simple binary classification methods while the schemes in [3, 11, 13] are multi-class vulnerability detection

*Corresponding author

 h.moallem@modares.ac.ir (H.M. Taheri);
shafieinejad@modares.ac.ir (A. Shafieinejad)

 <https://www.modares.ac.ir/~shafieinejad> (A. Shafieinejad)

ORCID(s): 0000-0001-0000-0000 (A. Shafieinejad)

algorithms. The performance of multi-class algorithms is usually lower compared to binary classification due to factors such as unbalanced vulnerability datasets and the similarity between different vulnerabilities.

- **Code slicing:** refers to choosing more important parts of the code for embedding and passing to the neural network and thereby ignoring unnecessary details. For example VulDeePecker [8] uses *code attention* concept for choosing important parts of source code. It focuses only on parts of a given code that is related to code attentions. This refinement removes excessive parts of code and thereby restricting the input size of training model as well as decreasing noise during training.

However, the extraction of code attention is not straightforward. Currently, VulDeePecker only detects vulnerabilities caused by C/C++ library/API function calls, as it generates code attention specifically around an API system call. Identifying the rest of code attentions has been left as an open problem. Similarly, the schemes in [4–11, 15, 16] perform some type code slicing prior to word embedding or vectorization phase.

- **Detection level of vulnerability:** refers to the extent of vulnerability localization. The first option is a coarse-grain algorithm that detects vulnerable functions in the given code. Other interesting options include fine-grain algorithms that perform detection at the code segment, statement, or line level. The schemes in [1–3, 13–15, 17] are function level, i.e., only identifies vulnerable functions. While the schemes in [8, 9, 11, 16] detect a code segment with vulnerability, the work in [5, 7, 10, 12, 18] perform detection at the line level and the schemes in [4, 6, 19] aim to detect vulnerability at statement level. The detection accuracy in a fine-grained algorithm is likely reduced compared to a coarse-grained algorithm due to the likelihood of incorrect detection of line or statement.

The concept of line or statement level detection may be criticized by researchers who define vulnerability as a group of statements, neither a single statement nor a line, which their execution provides a situation that can be exploited by attackers. However, line level detection is interesting as it aims to help developers identify vulnerabilities through appropriate warnings from the vulnerability verification system.

- **Intra or Inter-function detection:** Inter-function detection refers to identifying vulnerabilities that originate from flaws in multiple functions. In other

words, this type of vulnerability cannot be detected by examining just a single function. Thus, an inter-function algorithm should extract features from different functions and aggregate them to classify the sample as either vulnerable or non-vulnerable.

Although most of the prior work, such as [2, 4–6, 8, 12–14], focuses on intra-function vulnerability detection, the schemes presented in [3, 7, 9–11, 16, 19] offer an inter-function approach for vulnerability detection.

- **Transforming to an intermediate language:** Intermediate language, e.g. Java bytecodes, is more closer to the hardware than a high level language like C++ and thus has more environment and runtime information which can help to detect software vulnerabilities more accurate or at the opposite increase noise in source data. Most of the work, such as [2, 8, 9, 14], does not utilize intermediate languages and instead analyzes the main source code directly. However, some approaches do employ intermediate languages as the basis for feature extraction during training. For example, ISVSF [4] utilizes a custom intermediate representation of the source code, as does [10].
- **Multi-granularity:** refers to simultaneously multi-level detection, e.g., detecting at both line and function level. Multi-granularity aims to detect the type of vulnerability at the function level while simultaneously identifying the location of the vulnerability through line or statement-level detection. Most of the prior works are often single-granularity. However, some schemes such as [5, 7, 10, 19] are multi-granularity.

Table 1 summarizes the state-of-the-art literature on vulnerability detection along with our scheme concerning the aforementioned features.

1.1. Related work

In this section, we review the related work in three subsections: deep learning (DL) approaches, graph neural network (GNN) approaches, and additional methods, including image processing, lexical analysis, and long language models (LLMs).

1.1.1. DL-based approaches with code slicing

Li et al. introduced VulDeePecker [8], a DL-based multi-class vulnerability detection approach using code gadgets to represent programs. This method extracts code slices from source codes based on code representation graphs like CFG. Finally code slices are converted to code gadgets and feed to BLSTM neural network.

Saccante et al. developed a prototype tool [20] using Long-Short Term Memory Recurrent Neural Networks (LSTM-RNN) to automate vulnerability detection in source code.

Fidalgo et al. proposed a deep learning method [21] for classifying PHP SQL injection (SQLi) vulnerabilities. It utilizes an intermediate language to represent code slices and applied NLP techniques.

Zou et al. introduced μ VulDeePecker [11], a deep-learning-based system for multi-class vulnerability detection. It focuses on a specific part of source code, namely code gadget, and then transfers it to a binary string to feed relevant information to a BLSTM networks. Each code gadget is mainly created around a line of source code containing vulnerable system call such as `strcpy`.

Li et al. introduced SySeVR [9], a framework for DL-based multi-class vulnerability detection in C/C++ programs. It focuses on extracting semantic and syntax features from the source code and combining them before feeding to neural network. They used 5-fold cross-validation to train eight standard models, including LR (Logistic Regression), MLP (Multi-Layer Perception), DBN (Deep Belief Network), CNN (Convolutional Neural Network), LSTM (Long Short-Term Memory), GRU (Gated Recurrent Unit), BLSTM, and BGRU (Bidirectional GRU), thereby identifying the model with the highest F_1 measure.

Li et al. introduced VulDeeLocator [10] a fine grained vulnerability detection system works at statement level. This work uses intermediate code for feature extraction. It can find the approximate location of a vulnerability across multiple files.

Zhang et al. proposed ISVSF [4], a framework for self-detecting statement level vulnerabilities in Java. Each sentence is defined as an intermediate representation of the sub-block obtained from control flow abstract syntax tree (CFAST). It considers the syntax characteristics of Java and adopts sentence-level method representation and pattern exploration.

1.1.2. GNN based approaches

Zhou et al. proposed Devign [2], a graph neural network model tailored for graph-level classification tasks for multi-class vulnerability detection. It leverages diverse code semantic representations and uses a Conv module to effectively extract valuable features from node representations. Devign detects vulnerability at function level as a result of using a graph level prediction model.

Haojie et al. introduced VULMG [3], a vulnerability detection approach framed as graph classification.

It utilizes VecG as a code property graph to extract lexical, grammatical, and semantic features from source code as well as adjacency matrix capturing structural, control, and dependency information.

DeepWukong [22], proposed by Cheng et al., is another GNN-based approach vulnerability detection in C/C++ programs. It embeds code fragments into compact representations preserving programming logic.

Renjith and Aji focused on Android Java vulnerabilities from 2015 to June 2021 [23]. Both vulnerable and fixed code snippets are transformed into graphical representations using intermediate graph formats and a GNN-based vulnerability detection mechanism is developed to analyze them.

BGNN4VD [24] extracts syntax and semantic information from source code using abstract syntax trees (AST), CFG and data flow graphs (DFG). It employs a Convolutional Neural Networks (CNNs) to classify vulnerabilities.

Hin et al. proposed LineVD [19], a statement-level vulnerability detection, utilizing GNNs and transformers to identify vulnerable statements. The feature extraction component utilizes CodeBERT, a transformer-based model, to generate a total of $n + 1$ feature embeddings: one for the overall function and n embeddings for each statement.

Yang et al. [15] proposed an algorithm to detect vulnerable sentences within a function using different code representation graphs including CDG, AST and DDG. The initial vectors, constructed using a pre-trained Word2Vec model, are fed into a vulnerability detection model based on a heterogeneous graph transformer.

Zou et al. introduced mVulPreter [5], a multi-granularity vulnerability detection model that combines function-level (coarse-grained) and slice-level (fine-grained) detection approaches. It tries to provide interpretable results.

SedSVD, proposed by Dong et al. [6], is a graph-based statement-level detection which leverages a Code Property Graph (CPG) to capture semantic and syntactic information. The important parts of the CPG, extracted by so-called central nodes, are used to construct subgraphs. They are embedded and trained using a Relational Graph Convolutional Network (RGCN) and a MLP.

Wu et al. introduced SlicedLocator [7], a dual-grained model detection that predicts both program-level and statement-level vulnerabilities. It incorporates a sliced dependence graph to preserve inter-procedural relations while filtering out vulnerability-irrelevant statements. It creates attention-based code

embedding networks and employs a new LSTM-GNN model for fusion of semantic and structural modeling.

Zhang et al. introduced SDV [17], a static method for detecting function-level vulnerabilities in C/C++ programs. It utilizes a code property graph with Jump Graph Attention Network layers and convolutional pooling for feature extraction.

Tang et al. introduced CSGVD [1], another function-level vulnerability detection. It integrates a PE-BL module for semantic feature extraction which consists of CodeBERT as embedding layer and a BiLSTM layer with GNNs to extract the structured information.

RGAN [14], proposed by Tang et al., addresses the issue of code imbalance in vulnerability detection, particularly in node-level graph prediction, where the proportion of non-vulnerable nodes is relatively high. RGAN utilizes a residual graph attention network with a mean bi-affine attention pooling mechanism to eliminate this issue.

Nguyen et al. introduced CodeJIT [16], a code-centric approach for just-in-time vulnerability detection that focuses on code changes, specifically commits in Git. It employs a graph-based representation, along with GNNs, to distinguish dangerous code changes from safe ones based on semantic encoding.

1.1.3. Other DL-based approaches

The schemes in [12, 25] employ DL-based image classification in vulnerability detection. The former, called VulCNN, addresses scalability and accuracy in detecting vulnerabilities within large-scale source code. It converts each function of source code into an image while preserving program intricacies. The latter, called VulGAI, extracts degree, closeness and second order metrics from the source code to build an image. Both schemes use a CNN model to detect vulnerabilities. performance.

Russell et al. proposed a vulnerability detection using deep representation learning that directly interprets lexed source code[26]. The scheme uses a custom C/C++ lexer to create a generic representation of function and explores both CNNs and RNNs for feature extraction from the embedded source representation.

Wartschinski et al. introduced VUDENC [27], a DL-based vulnerability detection tool that automatically learns features of vulnerable code from a large Python codebase. It applies a Word2Vec model to identify semantically similar code tokens and to provide a vector representation. Then, these vectors are fed into a LSTM network to classify vulnerable code token sequences.

Fu et al. introduced VulExplainer [13], a transformer based vulnerability classification for addressing diversity challenges. First, based on CWE abstract types, it redistributes labels into sub-distributions to reach a set of groups with more balanced label distribution. Then, it trains TextCNN teachers on simplified distributions and a Transformer student model for generalization across different CWE-IDs.

Steenhoek et al. [28] conducted an empirical study replicating state-of-the-art deep learning models on vulnerability detection datasets. They analyzed model capabilities, training data impacts, and interpretability challenges across different vulnerability types and dataset compositions.

Li et al. [18] proposed a DL-based approach for line level vulnerability detection. It uses Clang to tokenize each line of function to get the tokens and their types. Further, position information for each token is added using positional encoding. Then, these vectors are fed into a hybrid neural network consisting of memory networks and multi-head attention mechanisms.

Zhang et al. worked on ChatGPT prompts [29] to develop a vulnerability detection engine based on Large language model (LLM). It integrates structural and sequential auxiliary information, optimizes prompt configurations for vulnerability detection tasks, and leverages ChatGPT's ability of memorizing multi-round dialogue.

1.2. Contributions

In this paper, we focus on a multi-class and fine grained vulnerability detection algorithm. To solve this problem, we utilized a node level prediction model of graph neural networks in addition to graph prediction model. Further, we take a combination of node and graph prediction model into account to get higher precision. Below are our contributions.

- We propose a multi-class Line-Level Vulnerability Detection (LLVD) system using a graph neural network, which not only detects the type of vulnerability but also identifies its location within the source code.
- We use a combination form of both node and graph level prediction to obtain a better accuracy for vulnerability detection without redesigning and retraining of learning model.
- Our scheme works with any desired abstraction level like CFG and DDG. However, our evaluations show that DDG obtains the best results in terms of precision, recall and F_1 -score.

Table 1

Summary of our existing DL-based work on vulnerability detection (F: Function, C: Code segment, L: Line and S: Statement)

Feature	ours	[8]	[2]	[9]	[10]	[11]	[3]	[4]	[15]	[19]	[5]	[6]	[7]	[12]	[13]	[1]	[18]	[17]	[16]	[14]
Graph neural networks	✓	×	✓	×	×	×	✓	✓	✓	✓	✓	✓	✓	×	×	✓	✓	✓	✓	✓
Detection level	L	C	F	C	L	C	F	C	F	S	L	S	L	L	F	F	L	F	C	F
Multi-class	✓	×	×	×	×	✓	✓	×	×	×	×	×	×	×	✓	×	×	×	×	×
Inter-function	×	×	×	✓	✓	✓	✓	×	×	✓	×	×	✓	×	×	×	×	×	✓	×
Code slicing	×	✓	×	✓	✓	✓	×	✓	✓	×	✓	✓	✓	×	×	×	×	×	✓	×
Intermediate language	×	×	×	×	✓	×	×	✓	×	×	×	×	×	×	×	×	✓	×	×	×
Multi-granularity	✓	×	×	×	✓	×	×	×	×	✓	✓	×	✓	×	×	×	×	×	×	×

- Our algorithm encounters a stage so-called *update prediction* which provides a fine grain detection at line of source code.

Our work is different from the schemes in [8–13] which we use a graph neural network for learning model.

Unlike most schemes, such as Devign [2], which function as binary classifiers, our scheme addresses multi-class vulnerability detection.

Further, our algorithm tends to detect the vulnerability at line level of source code while most of the schemes like [13] find it at function or code segment level. More specifically, few schemes including [3, 11, 13] which address multi-class vulnerability detection, are coarse grained algorithms where the vulnerability is identified at function or code segment level. Moreover, like [5, 7, 10, 19], our scheme has multi-granularity in which vulnerability is detected at both function and line level.

The rest of this paper is organized as follows: Section 2 introduces the system model and problem formulation. We present three algorithms which respectively detect vulnerability at function level (FLVD), line-level (LLVD) and combination of function and line level (FLLVD). The evaluation results are presented in Section 3. Finally, Section 4 concludes the paper.

2. Proposed algorithms for vulnerability detection

2.1. Problem Formulation

We need to design a scheme that identifies the type of vulnerability among a set of possible vulnerability types, as well as its location within a given source code. Let N_{class} denote the number of different vulnerability types. Further, assume that the target program source code contains a set of functions or procedures. The objective is to label each line of any function with a number from set $0, 1, 2, \dots, N_{class}$, where type-0 identifies non-vulnerable code while type- i ($1 \leq i \leq N_{class}$) corresponds to a Common Weakness Enumeration Identifier (CWE-ID).

Note that both type and location of vulnerability must be predicted accurately. That is we will encounter a false alarm by detecting either a wrong type or an incorrect location. The problem is quite general, but under specific assumptions, it becomes a well-known problem. For example, in the case of $N_{class} = 2$, it reduces to a binary classification problem that distinguishes between vulnerable lines and non-vulnerable ones. Moreover, if we assign the same label to each line of a function in the given source code, it transforms into a file-level vulnerability detection problem.

2.2. System Model

The overview of the proposed framework is shown in Figure 1. The primary objective of design is to ensure a high degree of flexibility. For instance, the framework supports both fine-grained and coarse-grained vulnerability detection by operating at the line and function levels, respectively. Further, our scheme is independent of programming language, supporting various source codes, including Java and C++. It converts these source codes to graph abstraction levels, such as Control Flow Graphs (CFG) or Data Dependency Graphs (DDG), preparing them for input into the neural network model.

According to Figure 1, our scheme composed of the following steps occurred in sequence:

- Preprocessing
- Embedding
- Training
- Test and Update prediction

The first step deals with primary input as a latest patch of software package and its previous vulnerable version to create a clean dataset as labeled source code. It contains a set of source files categorized into vulnerable and non-vulnerable files. It is supposed that each vulnerable file at least contains a line of code that is labeled as a specific vulnerability. Basically, these files are gathered from the latest patch of a software package

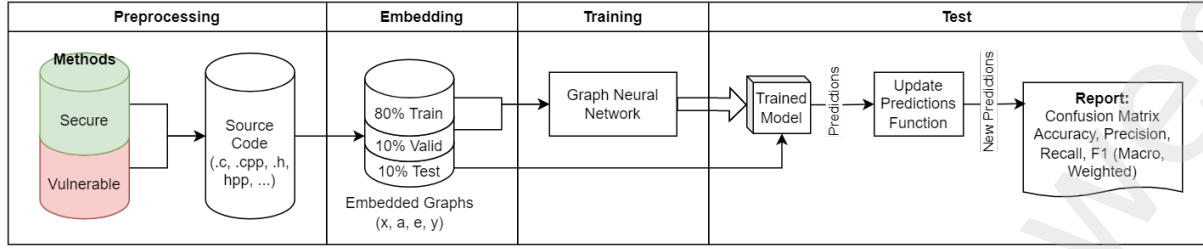


Figure 1: System Model Big Picture.

which eliminate a set of specific vulnerabilities. Thus, for each vulnerable file, there is a non-vulnerable version which leads to a balanced dataset.

For efficiency, in each vulnerable file, we select only a set of functions that contain at least one vulnerable line, removing the rest of the functions. In summary, we will have a set of function pairs (f_s, f_v) , where f_s is the secure version and f_v is the vulnerable version. Now, we extract a desired abstraction graph model for each function pair like CFG or DDG depending on our configuration with a suitable software tool. The output of this stage is graph pairs (G_s, G_v) , where the former is the secure version and the latter is the vulnerable version. The label of each node in G_s is set to 0 while some nodes of G_v is marked with 1 denoting that these nodes are corresponding to a vulnerable line.

The second stage performs graph embedding in which each G_s and G_v generated in previous step are transformed to a proper format which is suitable for feeding to neural network model. Basically, this step encounters word embedding where each node or edge containing string labels, i.e., part of source code, will be converted to a vector using a word embedding technique. The embedded has a specific data structure representing nodes, edges, adjacency matrix and nodes vulnerability labels.

The third step involves training where in each embedded graph of training set is fed into a customized neural network. More specifically, *GatedGraphConv* [30] model is used for training CFG graphs while *ECCConv* [31] is used for the graphs with edge-label like DDG. The result of this step is a trained model.

The final step performs testing to evaluate the trained model's performance. First, each graph of test set is fed into model to predict node or graph label. Then update prediction is done in which predicted node labels are further processed. This function alters the predicted output in order to improve the evaluation metrics. The basic idea of this phase is that unified the prediction labels to a set of nodes relating to a single line of source code. We will show that this step

increase the precision of the model even for real-world applications.

2.3. Line Level Vulnerability Detection

The first scheme (LLVD) is a fine-grain method which predicts vulnerabilities at node level using GNN.

2.3.1. Preprocessing

Since this step is highly tied with the incoming dataset, we will describe it in more detail in section 3.

2.3.2. Embedding

As mentioned earlier, we deal with a set of pair (f_s, f_v) denoting the safe and vulnerable version of a given function, respectively. Without loss of generality, we use the following abstraction graph models for converting a desired f_s or f_v into a graph model:

- CFG (Control Flow Graph)
- DDG (Data Dependency Graph)
- PDG (Program Dependency Graph)

The first item focuses on the flow of control through a program while the second concentrates on the flow of data and dependencies between operations. Further, PDG maintain some control dependency information in addition to data dependencies. The output of CFG is denoted by a pair (\mathbf{x}, \mathbf{a}) while the output of DDG and PDG is represented by a triple $(\mathbf{x}, \mathbf{a}, \mathbf{e})$ where :

- \mathbf{x} denotes the set of nodes,
- \mathbf{a} denotes the adjacency matrix which illustrates the transition between nodes and
- \mathbf{e} denotes the set of edges which describes data dependency between nodes.

Each statement in source code converts to a single or multiple nodes, based on its complexity. A simple assignment statement is mapped to a single node while complex instruction like if-else or function call is transformed to multiple nodes. Each node has a *code* field

which corresponds to a whole or part of a statement. Further, as each edge label in DDG corresponds to data dependency, it contains a variable or an expression.

To feed x and e into the learning model, the source code related to each node or edge, represented as an ASCII string, must be converted into an array of floats. This process, known as word/code embedding, is mainly implemented by means of either Word2Vec or BERT.

We examined both methods, but we disregarded Word2Vec due to its poor results. The second solution, which is based on the pre-trained model CodeBERT for the C++ language [32], yields better results. The output length of this model is 768. Let n and m denotes the number of nodes and edges, respectively. Thus, x and e are embedding to array of $n \times 768$ and $m \times 768$ float numbers. Moreover, a is embedding into array of n^2 elements. In CFG each element is 0 or 1 while in DDG it can 0 or a positive integer larger than 1 denoting the number of edges between two distinct nodes. In both cases, a will be a sparse matrix which is efficiently handled by suitable structure to optimize memory usage.

2.3.3. Training

We design two distinct neural networks for LLVD shown in Figure 2 and Figure 3, respectively. The former is used for graphs without edge like CFG while the latter is used for graphs with edge label like DDG. In the first model, the input x after normalization along with a are feeding to a 4-layer gated graph convolution networks. The input part a is not updated and is feeding to each hidden layer without any change while the input part x is gradually update in each hidden layer. After the hidden layers, the output is forwarded into a Dropout layer with rate 0.5 to avoid over-fitting. After that, the output is feeding into a Dense layer which reduces the output to $n \times N_{class}$ where N_{class} denotes the number of vulnerability classes.

In the second model, the input x as well as input e is passed through a normalization layer and a Dense layer with scale factor 768 to 128. Then, both of them along with a are feeding to an ECCConv layer. After that, the output is forwarded into a Dropout layer with rate 0.5 to avoid over-fitting. Finally, the output is feeding into a Dense layer which reduces the output to $n \times N_{class}$.

2.3.4. Update Predictions

As mentioned above, during CFG or DDG extraction, a single statement can be mapped into multiple nodes. Since LLVD detects vulnerabilities at the node level, it can be very fine-grained; that is, it may mark a

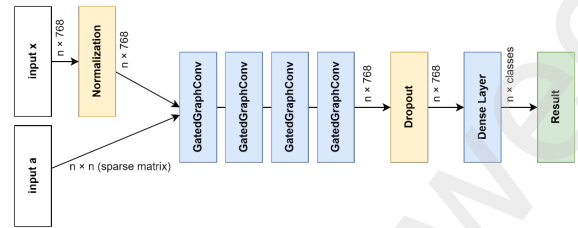


Figure 2: The node level prediction model for graphs without edge labels like CFG.

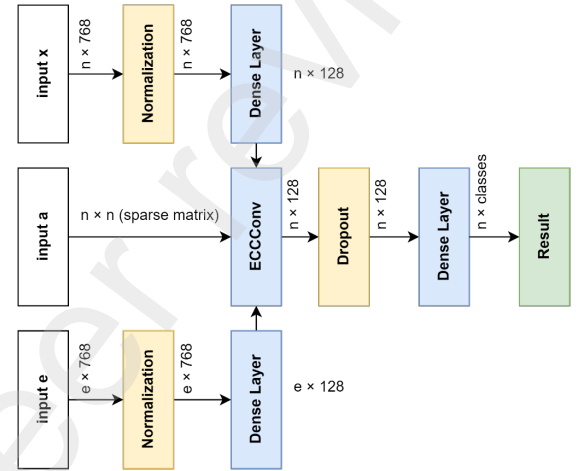


Figure 3: The node level prediction model for graphs with edge labels like DDG.

part of a statement as vulnerable while leaving the rest as non-vulnerable. To address this issue, we consider a new stage as *Update-prediction* which aggregates predicted label for multiple nodes corresponding to a single statement of source code. It ensures that the nodes extracted from a single statement share the same prediction.

For better understanding, consider the example in Figure 4, where a single vulnerability was detected in a node extracted from line 4 of the source code. Since two more nodes (6 and 7) are associated with line 4 of source code, *UpdatePrediction* will replace their labels with vulnerability class 2 which is currently assigned to node 5.

Algorithm 1 outlines the procedure for updating node-level predictions. It begins by initializing an index variable, α , at the first node. Then, a while loop iterates through the nodes, aggregating predictions for contiguous nodes associated with a specific line of code. Once these nodes are identified, the algorithm checks for the maximum predicted value among them. If a non-zero value is found, it will be assigned to all nodes in this

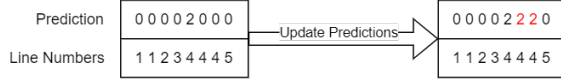


Figure 4: A simple example for update-prediction.

segment. Then, α is updated to the beginning of next segment, and the loop iterates until all nodes have been processed.

Algorithm 1 UpdatePrediction

Input : $P_1, \dots, P_n, L_1, \dots, L_n$

Output : P_1, \dots, P_n

```

1:  $\alpha \leftarrow 1$ 
2: while ( $\alpha < n$ ) do
3:    $\beta \leftarrow \alpha$ 
4:   while ( $\beta < n$  and  $L_{\beta+1} = L_\alpha$ ) do
5:      $\beta++$ 
6:   end while
7:    $\gamma \leftarrow \max_{i=\alpha}^{\beta} P_i$ 
8:   if ( $\gamma > 0$ ) then
9:      $(P_\alpha, \dots, P_\beta) \leftarrow (\gamma, \dots, \gamma)$ 
10:  end if
11:   $\alpha \leftarrow \beta + 1$ 
12: end while
13: return ( $P_1, \dots, P_n$ )

```

Algorithm 1 describes the procedure for updating node-level predictions in the FLLVD model. Similar to the LLVD algorithm, this method also begins by initializing the index α at the first node. The algorithm iterates over the nodes to identify contiguous segments corresponding to the same line. However, instead of summing the prediction scores, this algorithm assigns each node the maximum prediction score within the segment. This approach ensures that all nodes in a segment reflect the highest level of predicted vulnerability, thereby aligning with the function-level predictions of the FLLVD model.

2.4. Function Level Vulnerability Detection

Function Level Vulnerability Detection (FLVD) is our second scheme based on graph level prediction. It is a coarse-grain method and useful for enhancing other methods to create a multi-granularity vulnerability detection. It predicts the vulnerability class for the entire source code of a function without any localization. The training model for FLVD is the same as that for LLVD, except that the final Dense layer reduces the output to N_{class} bits in a one-hot encoding format. Note that FLVD can utilize both CFG and DDG abstractions during the embedding and training phases.

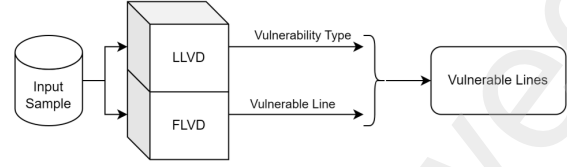


Figure 5: FLLVD as the combination of node-level and graph-level prediction model.

2.5. Function/Line Level Vulnerability Detection

The third scheme is Function/Line Level Vulnerability Detection (FLLVD). It is a multi-granularity method obtained through the combination of FLVD and LLVD models. This scheme prioritizes FLVD in order that if FLVD does not detect any vulnerabilities in a function, all nodes will be labeled as non-vulnerable, regardless of the LLVD predictions.

Otherwise, if FLVD detects a specific vulnerability y_g , we will take the LLVD prediction into account. More specifically, the nodes that LLVD indicates as vulnerable are marked with y_g , regardless of the type of vulnerability predicted by LLVD. Briefly, LLVD is used for the localization of vulnerabilities while the type of vulnerability is identified by FLVD. The details of FLLVD is depicted in Figure 5 as well as explained in Algorithm 2.

Algorithm 2 FLLVD-Prediction

Input : G, L_1, \dots, L_n

Output : P_1, \dots, P_n

```

1:  $y_G \leftarrow \text{FLVD-Prediction}(G)$ 
2:  $(P_1, \dots, P_n) \leftarrow (0, \dots, 0)$ 
3: if ( $y_G = 0$ ) then
4:   return ( $P_1, \dots, P_n$ )
5: end if
6:  $(Q_1, \dots, Q_n) \leftarrow \text{LLVD-Predict}(G)$ 
7:  $(q_1, \dots, q_n) \leftarrow \text{UpdatePrediction}(Q_1, \dots, Q_n, L_1, \dots, L_n)$ 
8: for  $i \leftarrow 1$  to  $n$  do
9:   if ( $q_i > 0$ ) then
10:     $P_i \leftarrow y_G$ 
11:   end if
12: end for
13: return ( $P_1, \dots, P_n$ )

```

3. Evaluations

3.1. Data preparation

3.1.1. Data gathering

We examined different variants of our scheme with both man-made and real world datasets : SARD

and BigVul. SARD (Software Assurance Reference Dataset) is a collection of test programs with documented vulnerability which vary from small synthetic programs to large applications in C, C++, Java, PHP, and Csharp, and cover over 150 classes of weaknesses. We selected a relatively balanced set of SARD datasets which summarized in Table 2. We will explain more details of process in subsection 3.1.2. For BigVul, the available samples are very limited in both vulnerability types and the number of samples. Therefore, we combined BigVul samples with SARD to achieve a sufficient number of samples while also ensuring a balanced set of different vulnerabilities

3.1.2. Graph generation

The preparation of dataset along with graph generation is shown in Figure 6. In SARD, both vulnerable and non-vulnerable code are included in a single file, with comments detailing each vulnerable statement. Therefore, we first separate the vulnerable code from the non-vulnerable code by creating two distinct files, effectively doubling the number of samples in the dataset. Each sample is then processed with the Joern tool [33] to extract the desired graph. The output of Joern can be further processed and converted into various abstraction graphs. By default, the CFG graph is extracted for each function and is used to ensure that the given code is a patched version of the existing code. It is done by applying the DotDiff algorithm to extracted graphs of both safe and vulnerable code (f_s and f_v) to verify that f_v is a patched version of f_s by comparing respective nodes and edges. It detects changes by comparing DOT files of each function in the vulnerable and non-vulnerable versions while ignoring non-affecting characters such as spaces and tabs during comparison.

Although CFG is generated by default, other abstraction formats, such as DDG and PDG can be created upon request. All extracted graphs are saved in DOT format. Each non-trivial line of code can be mapped to one or more nodes in the code graph. Since comments are omitted in the extracted graph, we label the nodes corresponding to a vulnerable line as vulnerable nodes and leave the others as non-vulnerable. Against SARD, in BigVul, patched functions are located in separate files. By comparing a desired function in the patched file with the one in previous file version, we can identify the vulnerable function.

After identifying patched functions, we can mark a specific node in the extracted graph as either vulnerable or non-vulnerable, as well as labeling the whole graph accordingly. Now, labeled graphs are ready for the next

stage, i.e., embedding. The output of SARD preparation and graph generation is shown in Table 2. The second column denotes the number of samples for each vulnerability class. The number of functions extracted from each class samples is shown in third column. Subsequent columns G_v , G_s , |CFG| and |DDG| denote the number of vulnerable and non-vulnerable graphs, respectively.

Figure 7 shows extracted dot files for a simple C++ program. The source code, shown in Figure 7a, defines a main function consists of 9 lines. For simplicity and efficient use of space we removed some unnecessary details from DOT files. The DOT file for CFG, DDG and PDG are shown in Figure 7b, Figure 7c, Figure 7d, respectively. Moreover, the corresponding graphs are depicted in Figure 7e, Figure 7f, Figure 7g, respectively. DOT file consists of two main parts: nodes and edges definitions. Each node corresponds to a statement in source file. The related line of source code is identified by the phrase L_i at the end of node definition. We see that trivial lines such as "", "else" and "" does not mapped to any node. Against, complex source codes, such as if-condition in line 5, are mapped to more than one nodes. Further, most of non-trivial code lines (such as 4, 6 and 8) are mapped to a single node. Nodes in DDG and PDG are almost the same as CFG with a slight difference that some declarations such as argument passing are considered as nodes.

The edges in a CFG strictly represent control dependencies between nodes, while in a DDG, they reflect data dependencies among nodes. In contrast, edges in a PDG may represent both control and data dependencies among nodes. The control and data dependencies for edges in the PDG DOT file are discriminated by the keywords CDG and DDG, respectively. Further, in graphs represented in Figure 7, control and data dependencies are distinguished by solid and dashed lines, respectively. Each data dependency edge is labeled by a variable or expressions which caused data dependency among nodes. There can be multiple edges between the same pair of nodes due to various variables that create data dependencies

3.2. Results for SARD

we use the following metrics for evaluations:

$$Recall = \frac{TP}{TP+FN} \quad (1)$$

$$Precision = \frac{TP}{TP+FP} \quad (2)$$

$$F_1 = \frac{2Precision \times Recall}{Recall + Precision} = \frac{2TP}{2TP+FP+FN} \quad (3)$$

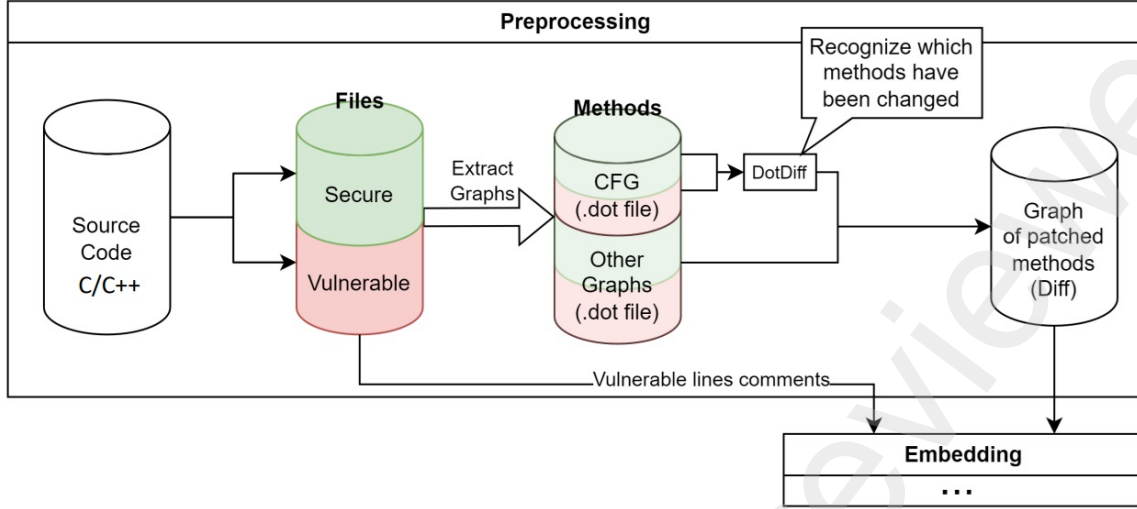


Figure 6: Preparation of dataset.

Table 2

SARD dataset used in our evaluation.

CWE	Samples	Functions	$ G_v $	$ G_s $
23	4100	11850	3700	8150
36	4100	11850	3700	8150
78	8200	23700	7400	16300
121	8177	24503	7308	17195
122	10116	30452	9088	21364
124	3456	10500	3084	7416
126	2388	7518	2136	5382
127	3456	10500	3084	7416
134	5160	24630	5100	19530
190	6660	27140	5891	21249
191	4854	20211	4364	15847
194	1968	5688	1776	3912
195	1968	5688	1776	3912
369	1548	5940	1332	4608
401	2812	10917	1657	9260
415	1724	6581	1482	5099
590	4355	12663	3040	9623
690	1640	4740	1480	3260
762	6388	24552	5488	19064
789	1720	6600	1480	5120
Average	3784	14311	3718	10593

Where TP, FP and FN denote the number of true positive, false positive and false negative, respectively. For multi-class evaluation, we use macro average F_1 which is computed by taking the arithmetic mean of all the per-class F_1 scores as :

$$\overline{F_1} = \frac{1}{N_{class}} \sum_{i=1}^{N_{class}} F_1(c) \quad (4)$$

Where $F_1(c)$ denotes F_1 score for class c .

The result of binary classification for SARD is shown in Figure 8. For each vulnerability type, we trained a specific model with a dataset containing that vulnerability along with its non-vulnerable samples. We examined both CFG and DDG with LLVD and FLLVD. The results shows that :

- In most cases, FLLVD has better performance in terms of F_1 for both code representation graph. More specifically, the improvement is about 1.6% for DDG and 6.2% for CFG. Since FLLVD only updates the type of prediction made by LLVD, it demonstrates that graph-level prediction is more accurate than node-level prediction in detecting vulnerability types.
- The performance of DDG is better than CFG regardless the detection algorithm. More specifically, the gain of DDG over CFG is about 15% for LLVD as well as 10.3% for FLLVD. It demonstrates that data dependency provides more information about vulnerabilities than control dependency.

However, some exceptions are observed. For example, FLLVD_{DDG} outperforms other methods except for CWE-127. Further, the performance for DDG is approximately higher than CFG in terms of F_1 metric. FLLVD_{CFG} exhibits the same performance as FLLVD_{DDG} for CWE 23 and 36. These vulnerabilities correspond to relative and absolute 'Path Traversal. It seems that the code embedding for file-system calls (e.g., open) is sufficient to catch this vulnerability, regardless of its dependency on previous or subsequent statements.

```

1. int main(int argc,
    char *argv[])
2. {
3.     char input[128];
4.     scanf("%s", input);
5.     if (strlen(input) > 10)
6.         printf("%s", input);
7.     else
8.         puts("Nothing");
9. }

```

```

#nodes
"1" (METHOD,main) L1
"2" (scanf,scanf("%s", input)) L4
"3" (strlen,strlen(input)) L5
"4" (strlen(input) > 10) L5
"5" (printf,printf("%s", input)) L6
"6" (puts,puts("Nothing")) L8
"7" (METHOD_RETURN,int) L1

```

```

#edges
"1" -> "2"
"2" -> "3"
"3" -> "4"
"4" -> "5"
"4" -> "6"
"5" -> "7"
"6" -> "7"

```

```

#nodes
"1" (METHOD,main) L1
"2" (METHOD_RETURN,int) L1
"3" (PARAM,int argc) L1
"4" (PARAM,char *argv) L1
"5" (scanf,scanf("%s", input)) L4
"6" (strlen(input) > 10) L5
"7" (printf,printf("%s", input)) L6
"8" (strlen,strlen(input)) L5
"9" (puts,puts("Nothing")) L8

```

```

#edges
"3" -> "2" "argc"
"4" -> "2" "argv"
"6" -> "2" "strlen(input)"
"6" -> "2" "strlen(input) > 10"
"7" -> "2" "input"
"7" -> "2" "printf("%s", input)"
"1" -> "3"
"1" -> "4"
"1" -> "5"
"8" -> "6" "input"
"1" -> "6"
"5" -> "8" "input"
"1" -> "8"
"1" -> "7"
"8" -> "7" "input"
"1" -> "9"

```

```

#nodes
"1" (METHOD,main) L1
"2" (METHOD_RETURN,int) L1
"3" (PARAM,int argc) L1
"4" (PARAM,char *argv) L1
"5" (scanf,scanf("%s", input)) L4
"6" (strlen(input) > 10) L5
"7" (printf,printf("%s", input)) L6
"8" (strlen,strlen(input)) L5
"9" (puts,puts("Nothing")) L8

```

```

#edges
"3" -> "2" "DDG: argc"
"4" -> "2" "DDG: argv"
"6" -> "2" "DDG: strlen(input)"
"6" -> "2" "DDG: strlen(input) > 10"
"7" -> "2" "DDG: input"
"7" -> "2" "DDG: printf("%s", input)"
"1" -> "3" "DDG: "
"1" -> "4" "DDG: "
"1" -> "5" "DDG: "
"8" -> "6" "DDG: input"
"1" -> "6" "DDG: "
"5" -> "8" "DDG: input"
"1" -> "8" "DDG: "
"1" -> "7" "DDG: "
"8" -> "7" "DDG: input"
"1" -> "9" "DDG: "
"6" -> "7" "CDG: "
"6" -> "9" "CDG: "

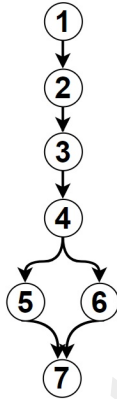
```

(a) C++ Source code

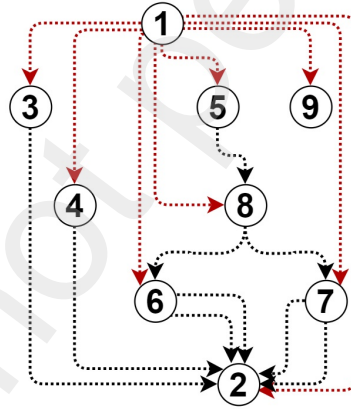
(b) CFG DOT file

(c) DDG DOT file

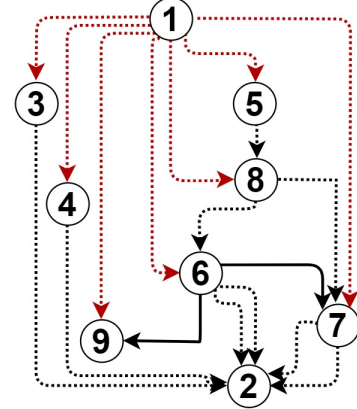
(d) PDG DOT file



(e) CFG



(f) DDG



(g) PDG

Figure 7: DOT files and dependency graph for different code graphs.

We examined our multi-class vulnerability detection algorithms using DDG and PDG on a subset of SARD, as shown in Table 2. For each vulnerability class, 1,000 graphs were randomly chosen, resulting in a balanced dataset across all vulnerability classes. However, the imbalance between vulnerable and non-vulnerable instances still remains. The results are shown in Figure 9. In contrast to binary classification, this experiment analyzes all vulnerabilities simultaneously. In addition to the 20 vulnerability types,

a separate class for non-vulnerable cases (referred to as 'Benign') is also included in the detection process. The overall F_1 (macro average) is 0.833, 0.902, 0.787 and 0.938 for $LLVD_{DDG}$, $FLLVD_{DDG}$, $LLVD_{PDG}$ and $FLLVD_{PDG}$, respectively. The important issues that we observe are :

- Except a few cases, FLLVD has better performance in terms of F_1 for both code representation graph. Further, the improvement for PDG is higher than DDG. More specifically, the gap between $LLVD_{DDG}$

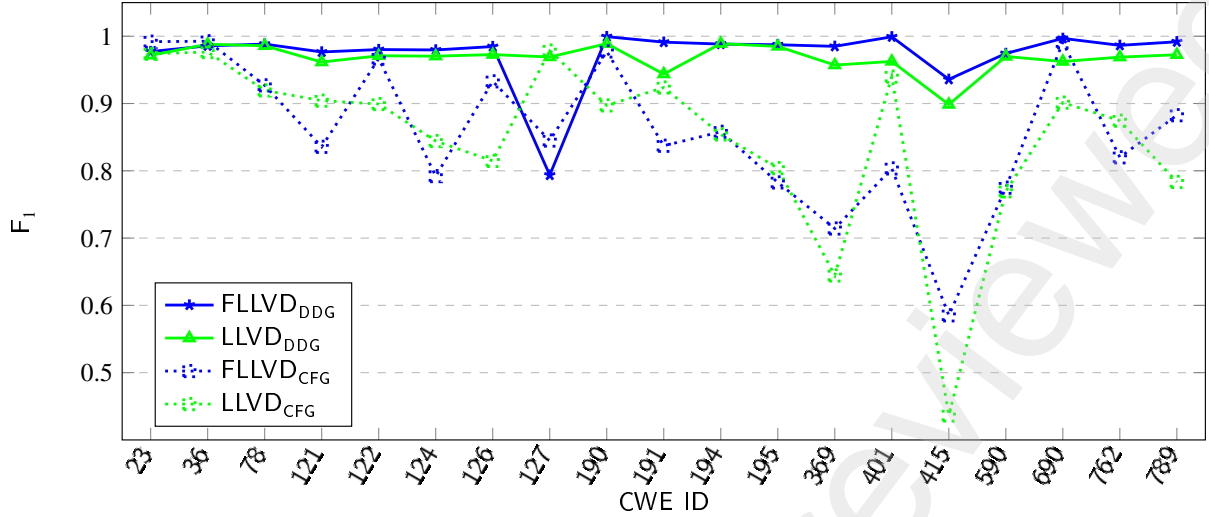


Figure 8: Results of our binary classification in terms of F_1 metric

and FLLVD_{DDG} is about 0.07 while the gap between LLVD_{PDG} and FLLVD_{PDG} reaches to 0.15. Similar to binary classification, we find that graph-level prediction operates more precisely than node-level prediction in detecting vulnerability types.

- The performance of PDG is better than DDG for FLLVD and vice versa for LLVD. More specifically, the gain of PDG over DDG is about 3.6% for FLLVD while it becomes -1.5% for LLVD. It demonstrates that control dependency in PDG improves FLLVD performance comparing to DDG while it provides some noise for LLVD.
- The detection rate of FLLVD for non-vulnerable code is approximately 100%, while for LLVD, it is about 99%. It is mainly originated from the fact that non-vulnerable portion of the dataset is larger than the vulnerable portion. Table 2 shows that vulnerable functions in SARD account for approximately 26% of the total functions. This imbalance ratio worsens at the node level, with the proportion of vulnerable nodes constituting about 2% of the total number of nodes

In summary, we can argue that FLLVD_{PDG} outperforms other variants except for CWE 369, which corresponds to 'divide by zero' where FLLVD_{DDG} demonstrates better performance by about 2%.

3.3. Results for SARD+BigVul

As mentioned above, BigVul samples are very limited in both types and number. More specifically the available vulnerability types are CWE-ID 78, 134, 191,

369 and 415. Therefore, we combined BigVul samples with SARD to achieve a sufficient number of samples while also ensuring a balanced set of different vulnerabilities. In this way, we arrive at a dataset containing approximately 600 graphs for each vulnerability type. It is worth noting that the graphs are randomly selected from both vulnerable and non-vulnerable samples.

The results for both LLVD and FLLVD, are represented at Table 3. We can see that :

- Similar to SARD result, Except CWE 369, FLLVD has better performance in terms of F_1 for both DDG and PDG. More specifically, the gap between LLVD_{DDG} and FLLVD_{DDG} is about 0.06 while the gap between LLVD_{PDG} and FLLVD_{PDG} reaches to 0.1.
- In contrast to SARD results, the performance of FLLVD_{DDG} is better than FLLVD_{PDG} in terms of overall F_1 . It is mainly originated from the ppor performance of FLLVD_{PDG} in detection of CWE 415 corresponding to "Double free". It seems that control dependency in PDG provides some noise for FLLVD in detecting CWE 415.

3.4. Comparison to previous schemes

Table 4 presents comparative results for the state of the art in vulnerability detection. It highlights the respective datasets, classification variables, detection levels, and average F1 scores. The third column indicates the datasets on which the respective schemes are trained or tested, including various combinations of existing datasets such as SARD, NVD, and BigVul, or a

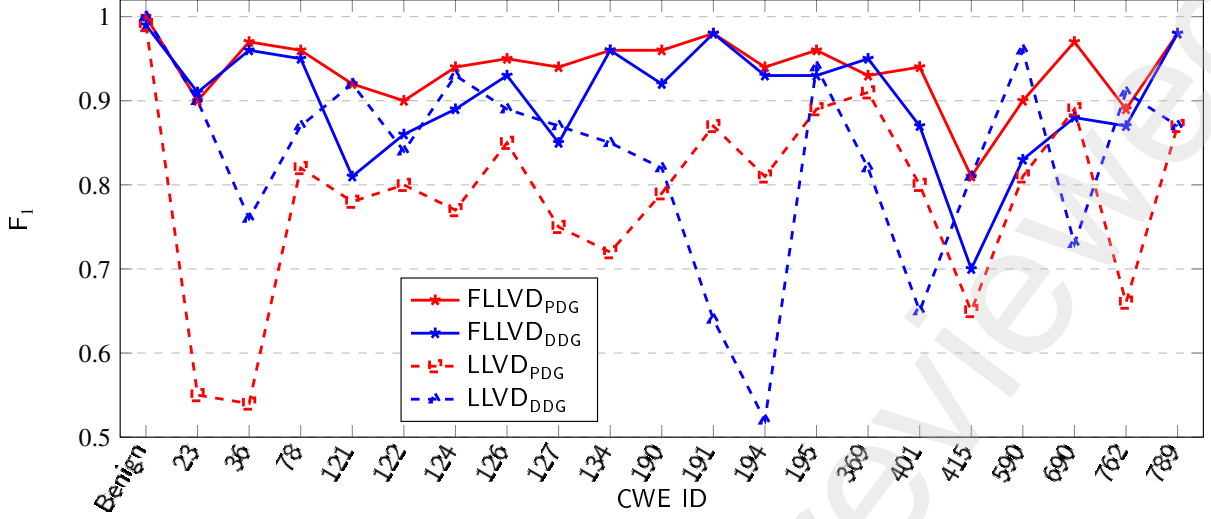


Figure 9: Performance of our multi-class detection in terms of F_1 metric

Table 3
Results for SARD + BigVul dataset in terms of F_1 .

CWE ID	LLVD _{DDG}	FLLVD _{DDG}	LLVD _{PDG}	FLLVD _{PDG}
0	0.99	0.99	0.99	0.99
78	0.83	0.85	0.79	0.87
134	0.69	0.96	0.68	0.88
191	0.73	0.82	0.27	0.78
369	0.78	0.62	0.84	0.76
415	0.56	0.73	0.36	0.26
$\overline{F_1}$	0.76	0.82	0.66	0.76

customized dataset, demonstrating the breadth and contextual focus of the datasets. The forth column indicates the number of distinct classes used in each scheme. A higher value suggests a more complex classification task that the model is designed to handle.

We can see that most of the work involves binary classification, just distinguishing between vulnerable and non-vulnerable code, regardless of the type of vulnerability. This limitation may restrict their applicability in scenarios with diverse codebases or various vulnerability types that must be accounted for. VulDeePecker+, μ VolDeePecker, and VULEXPLAINER are the only schemes that provide multiclass vulnerability detection. However, these schemes are coarse-grained, detecting vulnerabilities at the function or code segment level and leaving ambiguity about the exact location of the vulnerabilities.

Here, **FLLVD_{PDG}** demonstrates a superior performance with an F_1 score of 0.938, indicating a strong

balance between precision and recall. This performance is very close to the performance of VulDeePecker+ and μ VolDeePecker as multi-class detection schemes by regarding that they are more coarse grained than our scheme.

The analysis distinctly emphasizes the trade-off between the number of classes, detection level and detection performance. Our proposed methods demonstrate that a higher number of classes can still achieve high detection accuracy relative to others in the same category. Overall, our schemes provide evidence of state-of-the-art effectiveness in vulnerability detection when compared to various alternatives across different detection levels and datasets, particularly in the context of software vulnerability detection, where precise and contextual detection is crucial.

4. Conclusions

We proposed a novel multi-class vulnerability detection scheme leveraging graph neural networks to address the pressing need for effective vulnerability identification in software systems. By integrating Line Level Vulnerability Detection (LLVD) and Function Level Vulnerability Detection (FLVD), we established a comprehensive multi-granularity approach known as the File/Line Level Vulnerability Detection (FLLVD) scheme. This method not only enhances detection accuracy but also enables us to find both the type and location of each vulnerability within the given source code. The evaluation demonstrates promising results,

Table 4

Comparison of our schemes with prior work in terms of N_{class} , detection level and F_1 metric.

Scheme	Cite	Dataset	N_{class}	Detection level	$\overline{F_1}$
FLLVD_{PDG}	This work	SARD	21	Line	0.938
FLLVD_{DDG}	This work	SARD+BigVul	6	Line	0.828
SlicedLocator	[7]	SARD+NVD+CVEFixes[34]	2	Line	0.687
MAVLD	[18]	SARD+Draper	2	Line	0.99
VULGAI	[12]	SARD+NVD	2	Line	0.845
LineVD	[19]	BigVul	2	Statement	0.360
ISVSF	[4]	SARD+NVD	2	Statement	0.965
SedSVD	[6]	SARD+NVD	2	Statement	0.951
mVulPreter	[5]	NVD	2	Code segment	0.581
VulDeePecker	[8]	SARD	2	Code segment	0.934
VulDeePecker	[8]	NVD	2	Code segment	0.905
VulDeePecker+	[8]	SARD+NVD	41	Code segment	0.936
μ VulDeePecker	[11]	SARD+NVD	41	Code segment	0.946
CODEJIT	[20]	Custom	2	Code segment	0.84
Devign	[2]	NVD	2	Function	0.557
VULMG	[3]	SARD	2	Function	0.944
VDRG	[15]	SARD+Deving+Reveal	2	Function	0.812
VULEXPLAINER	[13]	BigVul	45	Function	0.639
CS-GVD	[1]	CodeXGLUE	2	Function	0.603
SDV	[17]	CHRRDB	2	Function	0.396
RGAN	[14]	Reveal	2	Function	0.473

with LLVD and FLLVD achieving substantial improvements in their F_1 metrics.

Our findings underscore the significance of employing a multi-granularity perspective in vulnerability detection, effectively bridging the gap between fine-grained and coarse-grained analyses. The proposed schemes can serve as foundational tools for further advancements in the field, paving the way for more reliable and secure software development practices. Future work should focus on expanding the dataset variations, refining the graph representations, and exploring the adaptation of our scheme to emerging vulnerability types, thereby enhancing its applicability in a continuously evolving threat landscape.

A. My Appendix

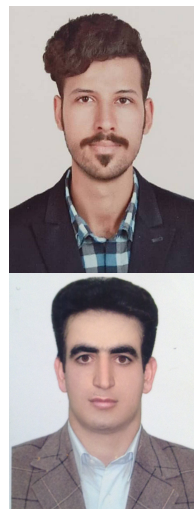
CRedit authorship contribution statement

Hamidreza M. Taheri: Algorithm Implementation, Software, Data curation, Writing - Original draft preparation. **Alireza Shafieinejad:** Conceptualization of this study, Methodology, Verification.

References

- [1] W. Tang, M. Tang, M. Ban, Z. Zhao, and M. Feng, "Csgvd: A deep learning approach combining sequence and graph embedding for source code vulnerability detection," *Journal of Systems and Software*, vol. 199, p. 111623, 2023.
- [2] Y. Zhou, S. Liu, J. Siow, X. Du, and Y. Liu, "Devign: Effective vulnerability identification by learning comprehensive program semantics via graph neural networks," *Advances in neural information processing systems*, vol. 32, 2019.
- [3] Z. Haojie, L. Yujun, L. Yiwei, and Z. Nanxin, "Vulmg: A static detection solution for source code vulnerabilities based on code property graph and graph attention network," in *2021 18th International Computer Conference on Wavelet Active Media Technology and Information Processing (ICCWAMTIP)*, pp. 250–255, IEEE, 2021.
- [4] H. Zhang, Y. Bi, H. Guo, W. Sun, and J. Li, "ISVSF: Intelligent vulnerability detection against java via sentence-level pattern exploring," *IEEE Systems Journal*, vol. 16, no. 1, pp. 1032–1043, 2021.
- [5] D. Zou, Y. Hu, W. Li, Y. Wu, H. Zhao, and H. Jin, "mvul-preter: A multi-granularity vulnerability detection system with interpretations," *IEEE Transactions on Dependable and Secure Computing*, 2022.
- [6] Y. Dong, Y. Tang, X. Cheng, Y. Yang, and S. Wang, "SedSVD: Statement-level software vulnerability detection based on Relational Graph Convolutional Network with subgraph embedding," *Information and Software Technology*, vol. 158, p. 107168, 2023.
- [7] B. Wu, F. Zou, P. Yi, Y. Wu, and L. Zhang, "Slicedlocator: Code vulnerability locator based on sliced dependence graph," *Computers & Security*, vol. 134, p. 103469, 2023.

- [8] Z. Li, D. Zou, S. Xu, X. Ou, H. Jin, S. Wang, Z. Deng, and Y. Zhong, "Vuldeepecker: A deep learning-based system for vulnerability detection," *arXiv preprint arXiv:1801.01681*, 2018.
- [9] Z. Li, D. Zou, S. Xu, H. Jin, Y. Zhu, and Z. Chen, "Sysevr: A framework for using deep learning to detect software vulnerabilities," *IEEE Transactions on Dependable and Secure Computing*, vol. 19, no. 4, pp. 2244–2258, 2021.
- [10] Z. Li, D. Zou, S. Xu, Z. Chen, Y. Zhu, and H. Jin, "Vuldeeloctor: a deep learning-based fine-grained vulnerability detector," *IEEE Transactions on Dependable and Secure Computing*, vol. 19, no. 4, pp. 2821–2837, 2021.
- [11] D. Zou, S. Wang, S. Xu, Z. Li, and H. Jin, "μVulDeePecker: A Deep Learning-Based System for Multiclass Vulnerability Detection," *IEEE Transactions on Dependable and Secure Computing*, vol. 18, no. 5, pp. 2224–2236, 2021.
- [12] C. Zhang and Y. Xin, "VulGAI: vulnerability detection based on graphs and images," *Computers & Security*, vol. 135, p. 103501, 2023.
- [13] M. Fu, V. Nguyen, C. K. Tantithamthavorn, T. Le, and D. Phung, "Vulexplainer: A transformer-based hierarchical distillation for explaining vulnerability types," *IEEE Transactions on Software Engineering*, 2023.
- [14] M. Tang, W. Tang, Q. Gui, J. Hu, and M. Zhao, "A vulnerability detection algorithm based on residual graph attention networks for source code imbalance (rgan)," *Expert Systems with Applications*, vol. 238, p. 122216, 2024.
- [15] H. Yang, H. Yang, L. Zhang, and X. Cheng, "Source Code Vulnerability Detection Using Vulnerability Dependency Representation Graph," in *2022 IEEE International Conference on Trust, Security and Privacy in Computing and Communications (TrustCom)*, pp. 457–464, IEEE, 2022.
- [16] S. Nguyen, T.-T. Nguyen, T. T. Vu, T.-D. Do, K.-T. Ngo, and H. D. Vo, "Code-centric learning-based just-in-time vulnerability detection," *Journal of Systems and Software*, p. 112014, 2024.
- [17] C. Zhang and Y. Xin, "Static vulnerability detection based on class separation," *Journal of Systems and Software*, vol. 206, p. 111832, 2023.
- [18] M. Q. Li, B. C. Fung, and A. Diwan, "A Novel Deep Multi-head Attentive Vulnerable Line Detector," *Procedia Computer Science*, vol. 222, pp. 35–44, 2023.
- [19] D. Hin, A. Kan, H. Chen, and M. A. Babar, "Linevd: Statement-level vulnerability detection using graph neural networks," in *Proceedings of the 19th international conference on mining software repositories*, pp. 596–607, 2022.
- [20] N. Saccante, J. Dehlinger, L. Deng, S. Chakraborty, and Y. Xiong, "Project achilles: A prototype tool for static method-level vulnerability detection of Java source code using a recurrent neural network," in *2019 34th IEEE/ACM International Conference on Automated Software Engineering Workshop (ASEW)*, pp. 114–121, IEEE, 2019.
- [21] A. Fidalgo, I. Medeiros, P. Antunes, and N. Neves, "Towards a deep learning model for vulnerability detection on web application variants," in *2020 IEEE International Conference on Software Testing, Verification and Validation Workshops (ICSTW)*, pp. 465–476, IEEE, 2020.
- [22] X. Cheng, H. Wang, J. Hua, G. Xu, and Y. Sui, "Deepwukong: Statically detecting software vulnerabilities using deep graph neural network," *ACM Transactions on Software Engineering and Methodology (TOSEM)*, vol. 30, no. 3, pp. 1–33, 2021.
- [23] G. Renjith and S. Aji, "Vulnerability Analysis and Detection Using Graph Neural Networks for Android Operating System," in *Information Systems Security: 17th International Conference, ICISS 2021, Patna, India, December 16–20, 2021, Proceedings 17*, pp. 57–72, Springer, 2021.
- [24] S. Cao, X. Sun, L. Bo, Y. Wei, and B. Li, "Bgnn4vd: Constructing bidirectional graph neural-network for vulnerability detection," *Information and Software Technology*, vol. 136, p. 106576, 2021.
- [25] Y. Wu, D. Zou, S. Dou, W. Yang, D. Xu, and H. Jin, "VulCNN: An Image-inspired Scalable Vulnerability Detection System," 2022.
- [26] R. Russell, L. Kim, L. Hamilton, T. Lazovich, J. Harer, O. Ozdemir, P. Ellingwood, and M. McConley, "Automated vulnerability detection in source code using deep representation learning," in *2018 17th IEEE international conference on machine learning and applications (ICMLA)*, pp. 757–762, IEEE, 2018.
- [27] L. Wartschinski, Y. Noller, T. Vogel, T. Kehrner, and L. Grunske, "VUDENC: vulnerability detection with deep learning on a natural codebase for python," *Information and Software Technology*, vol. 144, p. 106809, 2022.
- [28] B. Steenhoek, M. M. Rahman, R. Jiles, and W. Le, "An empirical study of deep learning models for vulnerability detection," in *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*, pp. 2237–2248, IEEE, 2023.
- [29] C. Zhang, H. Liu, J. Zeng, K. Yang, Y. Li, and H. Li, "Prompt-enhanced software vulnerability detection using chatgpt," in *Proceedings of the 2024 IEEE/ACM 46th International Conference on Software Engineering: Companion Proceedings*, pp. 276–277, 2024.
- [30] Y. Li, D. Tarlow, M. Brockschmidt, and R. Zemel, "Gated graph sequence neural networks," *arXiv preprint arXiv:1511.05493*, 2015.
- [31] M. Simonovsky and N. Komodakis, "Dynamic edge-conditioned filters in convolutional neural networks on graphs," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3693–3702, 2017.
- [32] https://huggingface.co/neulab/codebert_cpp, "Hugging face: Codebert c++ pretrained model," 2024. Accessed: 2024-9-20.
- [33] <https://joern.io>, "Joern, the bug hunter's workbench," 2024. Accessed: 2024-9-20.
- [34] G. Bhandari, A. Naseer, and L. Moonen, "Cvefixes: automated collection of vulnerabilities and their fixes from open-source software," in *Proceedings of the 17th International Conference on Predictive Models and Data Analytics in Software Engineering*, pp. 30–39, 2021.



Hamidreza M. Taheri was born in Tehran, Iran in 1999. He received the B.S. degree in computer engineering from Isfahan University of Technology, in 2021. He was also received M.Sc. degree in computer engineering from Tarbiat Modares University, Tehran, Iran, in 2024. His research area includes vulnerability detection and graph neural networks. He also works as a security engineer.

Alireza Shafieinejad is an Assistant Professor in the Department of Electrical and Computer Engineering, Tarbiat Modares University, Tehran, Iran. He received the B.S. degree in computer engineering from the Sharif University of Technology, Tehran, in 1999. He was also received M.Sc. and Ph.D. degrees in electrical engineering from Isfahan University of Technology, Isfahan, Iran,

respectively in 2002 and 2013. Since 2015, he has been an Assistant Professor at Tarbiat Modares University, Tehran, Iran. His research interests are in the area of Wireless Network Coding, Network Security, Design and Analysis of Cryptography Algorithms. At Tarbiat Modares University, he leads research in network security and penetration testing in SE-Lab (Security Evaluation Laboratory).