

Vulnerability Detection with Fine-Grained Interpretations

Yi Li

New Jersey Inst. of Technology
New Jersey, USA
yl622@njit.edu

Shaohua Wang*

New Jersey Inst. of Technology
New Jersey, USA
davidsw@njit.edu

Tien N. Nguyen

University of Texas at Dallas
Texas, USA
tien.n.nguyen@utdallas.edu

ABSTRACT

Despite the successes of machine learning (ML) and deep learning (DL) based vulnerability detectors (VD), they are limited to providing only the decision on whether a given code is vulnerable or not, without details on *what part of the code is relevant to the detected vulnerability*. We present IVDETECT, an *interpretable vulnerability detector* with the philosophy of using Artificial Intelligence (AI) to detect vulnerabilities, while using Intelligence Assistant (IA) via providing VD interpretations in terms of vulnerable statements.

For vulnerability detection, we separately consider the vulnerable statements and their surrounding contexts via data and control dependencies. This allows our model better discriminate vulnerable statements than using the mixture of vulnerable code and contextual code as in existing approaches. In addition to the coarse-grained vulnerability detection result, we leverage *interpretable AI* to provide users with *fine-grained interpretations that include the sub-graph in the Program Dependency Graph (PDG) with the crucial statements* that are relevant to the detected vulnerability. Our empirical evaluation on vulnerability databases shows that IVDETECT outperforms the existing DL-based approaches by 43%–84% and 105%–255% in top-10 nDCG and MAP ranking scores. IVDETECT correctly points out the vulnerable statements relevant to the vulnerability via its interpretation in 67% of the cases with a top-5 ranked list. It improves over baseline interpretation models by 12.3%–400% and 9%–400% in accuracy.

CCS CONCEPTS

• Security and privacy → Software security engineering.

KEYWORDS

Vulnerability Detection; Deep Learning; Explainable AI; Interpretable AI

ACM Reference Format:

Yi Li, Shaohua Wang, and Tien N. Nguyen. 2021. Vulnerability Detection with Fine-Grained Interpretations. In *Proceedings of the 29th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE '21)*, August 23–28, 2021, Athens, Greece. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3468264.3468597>

*Corresponding Author

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
ESEC/FSE '21, August 23–28, 2021, Athens, Greece

© 2021 Association for Computing Machinery.
ACM ISBN 978-1-4503-8562-6/21/08...\$15.00
<https://doi.org/10.1145/3468264.3468597>

1 INTRODUCTION

Software vulnerabilities have caused substantial damage to society's software infrastructures. Automated vulnerability detection (VD) approaches can be broadly classified into two categories: program analysis (PA)-based [1, 2, 7–9, 33] and machine learning (ML)-based [23, 28, 30]. The PA-based VD techniques have often focused on solving the *specific types of vulnerabilities* such as BufferOverflow [3], SQL Injection [6], Cross-site Scripting [5], Authentication Bypass [4], etc. In addition to those types, the more general software vulnerabilities, e.g., in API usages of libraries/frameworks, have manifested in various forms. To detect them, machine learning (ML) and deep learning (DL) have been leveraged to implicitly learn the patterns of vulnerabilities from prior vulnerable code [15, 20, 37].

Despite several advantages, the ML/DL-based VD approaches are still limited to providing only coarse-grained detection results on whether an entire given method is vulnerable or not. In comparison with the PA-based approaches, they fall short in the ability to elaborate on the *fine-grained details* of the lines of code with specific statements that might be involved in the detected vulnerability. One could use fault localization (FL) techniques [16] to locate the vulnerable statements, however they require large, effective test suites. Due to such feedback at the coarse granularity from the existing ML/DL-based VD tools, developers would not know where and what to look for and to fix the vulnerability in their code. This hinders them in investigating the potential vulnerabilities.

To raise the level of ML/DL-based VD, we present IVDETECT, an *interpretable VD* with the philosophy of using *Artificial Intelligence* to detect coarse-grained vulnerability, while leveraging *Intelligence Assistant* via interpretable ML to provide fine-grained interpretations in term of vulnerable statements relevant to the vulnerability.

For coarse-grained *vulnerability detection*, our novelty is the *context-aware representation learning of the vulnerable code*. During training, the existing ML/DL-based VD approaches [20, 37] take the entire vulnerable code in a method as the input without distinguishing the vulnerable statements from the surrounding contextual code. Such distinction from vulnerable code and the contexts during training enable IVDETECT to better learn to discriminate the vulnerable code and benign ones. We represent source code via program dependence graph (PDG) and we treat the vulnerability detection problem as graph-based classification via Graph Convolution Network (GCN) [17] with feature-attention (FA), namely FA-GCN. The vulnerable statements, along with surrounding code, are encoded during the code representation learning.

For *fine-grained interpretation*, as the given method is deemed as vulnerable by IVDETECT, our novelty is to *leverage interpretable ML* [36] to provide the interpretation in term of the vulnerable statements as part of the PDG that are involved to the detected vulnerability. The rationale for choosing PDG sub-graph as an interpretation is

that a vulnerability often involves the data and control dependencies among the statements [26].

To derive the vulnerable statements as the interpretation, we leverage the interpretable ML model, GNNExplainer [36], that “explains” on why a model has arrived at its decision. Specifically, after vulnerability detection, to produce interpretation, IVDETECT takes as input the FA-GCN model along with its decision (vulnerable or not), and the input PDG G_M of the given method M . The goal is to find the interpretation subgraph, which is defined as a minimal sub-graph \mathcal{G} in the PDG of M that *minimizes the prediction scores between using the entire G_M and using \mathcal{G}* . To that end, we leverage GNNExplainer [36] in which the searching for \mathcal{G} is formulated as the learning of the edge-mask set EM . The idea is that if an edge belongs EM , i.e., *if it is removed from G_M , and the decision of the model is affected, then the edge is crucial and must be included in the interpretation for the detection result*. Thus, the minimal sub-graph \mathcal{G} in PDG contains the nodes and edges, i.e., the *crucial statements and program dependencies, that are most decisive/relevant to the detected vulnerability* when the decision is vulnerable.

Using our results, a practitioner would 1) examine the ranked list of potentially vulnerable methods, and 2) use the interpretation to further investigate what statements in the code that caused the model to predict that vulnerability.

We conducted several experiments to evaluate IVDETECT in both vulnerability detection at the method level and interpretation in term of vulnerable statements. We use 3 large C/C++ vulnerability datasets: Fan [13], Reveal [11] and FFMpeg+Qemu [37]. For the method-level VD, our results show that IVDETECT outperforms the existing ML/DL-based approaches [11, 19, 20, 27, 37] by 43%–84% and 105%–255% at the top 10 list for two ranking scores nDCG and MAP, respectively. For the statement-level interpretation, IVDETECT correctly points out the vulnerable statements relevant to the vulnerability in 67% of the cases with a top-5 ranked list. It improves over the baseline ATT [36] and GRAD [36] interpretation models by 12.3%–400% and 9%–400% in accuracy, respectively.

The contributions of this paper include:

A. Interpretable VD with Fine-grained Interpretations

a. Vulnerability Detection with Fine-grained Interpretations: IVDETECT is the first approach to leverage interpretable ML to enhance VD with *fine-grained* details on PDG sub-graphs, statements, and dependencies relevant to the detected vulnerability.

b. Context-aware Representation Learning of vulnerable code: The novelty of our representation learning of vulnerable code is *our consideration of the contextual code surrounding the vulnerable statements and fixes* to better train the VD model.

B. Empirical Evaluation. Our results show IVDETECT’s high accuracy in both detection and interpretation (See data/results at [10]).

2 MOTIVATION

2.1 Motivating Example

Figure 1 shows the method `ec_device_ioctl_xcmd` in Linux 4.6, which constructs the I/O control command for the ChromeOS devices. This is listed as a vulnerable code within Common Vulnerabilities and Exposures (CVE-2016-6156) in the National Vulnerability Database.

The commit log of the corresponding fix stated that

```

1  static long ec_device_ioctl_xcmd(struct cros_ec_dev *ec, void __user *arg)
2  {
3      long ret;
4      struct cros_ec_command u_cmd;
5      struct cros_ec_command *s_cmd;
6      if (copy_from_user(&u_cmd, arg, sizeof(u_cmd)))
7          return -EFAULT;
8      if ((u_cmd.outsize > EC_MAX_MSG_BYTES) || (u_cmd.insize > EC_MAX_MSG_BYTES))
9          return -EINVAL;
10     s_cmd = kmalloc(sizeof(*s_cmd) + max(u_cmd.outsize, u_cmd.insize), GFP_KERNEL);
11     if (!s_cmd)
12         return -ENOMEM;
13     if (copy_from_user(s_cmd, arg, sizeof(*s_cmd) + u_cmd.outsize) {
14         ret = -EFAULT;
15         goto exit;
16     }
17     + if (u_cmd.outsize != s_cmd->outsize ||
18     +   u_cmd.insize != s_cmd->insize) {
19     +   ret = -EINVAL;
20     +   goto exit;
21     + }
22     s_cmd->command += ec->cmd_offset;
23     ret = cros_ec_cmd_xfer(ec->ec_dev, s_cmd);
24     /* Only copy data to userland if data was received. */
25     if (ret < 0)
26         goto exit;
27     - if (copy_to_user(arg, s_cmd, sizeof(*s_cmd) + u_cmd.insize)
28     + if (copy_to_user(arg, s_cmd, sizeof(*s_cmd) + s_cmd->insize)
29         ret = -EFAULT;
30 exit:
31     kfree(s_cmd);
32     return ret;
33 }

```

Figure 1: CVE-2016-6156 Vulnerability in Linux 4.6

“At line 6 and line 13, the driver fetches user space data by pointer `arg` via `copy_from_user()`. The first fetched value (stored in `u_cmd`) (line 6) is used to get the `in_size` and `out_size` elements and allocation a buffer (`s_cmd`) at line 10 so as to copy the whole message to driver later at line 13, which means the copy size of the whole message (`s_cmd`) is based on the old value (`u_cmd.outsize`) from the first fetch. Besides, the whole message copied at the second fetch also contains the elements of `in_size` and `out_size`, which are the new values. The new values from the second fetch might be changed by another user thread under race condition, which will result in a double-fetch bug when the inconsistent values are used.”

Thus, to fix this bug, a developer added the code at lines 17–21 to make sure that `u_cmd.outsize` and `u_cmd.insize` have not changed due to race condition between the two fetching calls. Moreover, memory access might be also beyond the array boundary, causing a buffer overflow within the method call `cros_ec_cmd_xfer(...)`, when the command is transferred to the ChromeOS device at line 23.

Another issue is at line 27 with `copy_to_user`. The method call `cros_ec_cmd_xfer(...)` can set `s_cmd->insize` to a lower value. Thus, the new smaller value must be used to avoid copying too much data to the user: `u_cmd.insize` at line 27 is changed into `s_cmd->insize`.

This vulnerable code could potentially cause the damages such as denial of service, buffer overflow, program crash, etc. Deep learning (DL) advances enable several approaches [20, 37] to *implicitly learn* from the history the patterns of vulnerable code, and to detect *more general vulnerabilities*. However, they are still limited in comparison with program analysis-based approaches in the ability to provide any detail on the *fine-grained* level of the vulnerable statements, and on why the model has decided on the vulnerability. For example, the PA-based approaches, e.g., a race detection technique could potentially detect the involvement of the two fetching statements at line 6 and line 13. The method in Figure 1 might be deemed as vulnerable by a DL-based model. But without any fine-grained details, a developer would not know where and what to investigate

```

1 static long ec_device_ioctl_xcmd(struct cros_ec_dev *ec, void __user *arg)
2 {
3     long ret;
4     struct cros_ec_command u_cmd;
5     struct cros_ec_command *s_cmd;
6     if (copy_from_user(&u_cmd, arg, sizeof(u_cmd)))
7         return -EFAULT;
8     if ((u_cmd.outsize > EC_MAX_MSG_BYTES) || (u_cmd.insize > EC_MAX_MSG_BYTES))
9         return -EINVAL;
10    s_cmd = kmalloc(sizeof(*s_cmd) + max(u_cmd.outsize, u_cmd.insize),
11                  GFP_KERNEL);
12    if (!s_cmd)
13        return -ENOMEM;
14    if (copy_from_user(s_cmd, arg, sizeof(*s_cmd) + u_cmd.outsize) ||
15        goto exit;
16    }
17    if (u_cmd.outsize != s_cmd->outsize ||
18        u_cmd.insize != s_cmd->insize) {
19        ret = -EINVAL;
20        goto exit;
21    }
22    s_cmd->command += ec->cmd_offset;
23    ret = cros_ec_cmd_xfer(ec->ec_dev, s_cmd);
24    /* Only copy data to userland if data was received. */
25    if (ret < 0)
26        goto exit;
27    if (copy_to_user(arg, s_cmd, sizeof(*s_cmd) + u_cmd.insize))
28        if (copy_to_user(arg, s_cmd, sizeof(*s_cmd) + s_cmd->insize))
29            ret = -EFAULT;
30    exit:
31    kfree(s_cmd);
32    return ret;
33 }

```

Figure 2: Interpretation Sub-Graph for Figure 1

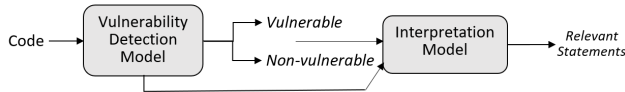


Figure 3: Overview of IVDetect

next. This would make the output of a DL model less constructive in VD. Moreover, a fault localization technique [16], which locates buggy statements, would need a large, effective test suite.

Regarding detection, the existing DL-based approaches [20, 37] do not fully exploit all the available information on the vulnerable code during training. For example, during training, we know that lines 23 and 27 are vulnerable/buggy, and other *relevant statements via data/control dependencies provide contextual information for the vulnerable ones*. However, the existing approaches [20, 37] do not consider the vulnerable statements and do not use the contextual code to help a model discriminate the vulnerable and non-vulnerable ones. The entire method would be fed to a DL model.

IVDetect Approach. We introduce IVDetect, an DL-based, *interpretable vulnerability detection* approach that goes beyond the decision of vulnerability by providing the *fine-grained* interpretation in term of the vulnerable statements. Specifically, as the method is deemed as vulnerable by IVDetect, it will provide a *list of important statements as part of the program dependence graph (PDG) that are relevant to the detected vulnerability*. For example, it provides the partial sub-graph of the PDG including the statements at the lines 13–15, 22–23, and 25–27 in Figure 2 for the vulnerable code at line 23 and line 27. We use the PDG sub-graph including important statements for fine-grained VD since they will give a developer the hints on the program dependencies relevant to the vulnerability for further investigation. Moreover, if our model determines the code as non-vulnerable, it can also produce the key sub-graph of the PDG with key statements that are deemed to be safe.

2.2 Key Ideas and Architecture Overview

IVDetect has two main modules (Figure 3): graph-based vulnerability detection model, and graph-based interpretation model. The input is the source code of all methods in a project. The output is the ranked list of methods with the detection result/score and the interpretation (PDG sub-graph). Let us explain our key ideas.

2.2.1 Graph-based Vulnerability Detection Model (Section 3). As seen in Section 2.1, a vulnerability is usually exhibited as multiple statements are exploited, thus, it is natural to capture the vulnerable code as a sub-graph in the PDG with the data and control flows. To do so, we model the vulnerability detection via the Graph Convolutional Network (GCN) [17] as follows. The PDG of a method M is represented as a graph $G_N = (V, E)$ in which V is a set of nodes representing the statements, and E is a set of edges representing the data/control dependencies. A feature description x_V is for every node v , which represents a property of a node, e.g., variable name, etc. Features are summarized in a $N \times D$ feature matrix X_M (N : number of nodes and D is the number of input features). Let f be a label function on the statements and methods $f: V \rightarrow \{1, \dots, C\}$ that maps a node in V and an entire method to one of the C classes. In IVDetect, $C=2$ for vulnerable (\mathcal{V}) and non-vulnerable ($\mathcal{N}\mathcal{V}$).

For training on (non-)vulnerable code in the training set, GCN performs similar operations as CNN where it learns the features with a small filter/window sliding over PDG sub-structure. Differing from image data with CNN, the neighbors of a node in GCN are unordered and variable in size. To predict if a method M is vulnerable, its PDG G_M with the associated feature set $X_M = \{x_j | v_j \in G_M\}$ are built. GCN learns a conditional distribution $P(Y|G_M, X_M)$, where Y is a random variable representing the labels $\{1, \dots, C\}$. That distribution indicates the probability of the graph G_M belonging to each of the classes $\{1, \dots, C\}$, i.e., M is vulnerable or not (Section 3).

2.2.2 Distinction between Vulnerable Statements and Surrounding Contexts. During training, for each vulnerable statement s in a method in the training dataset, we distinguish s and the surrounding contextual statements for s . A context consists of the statements with data and/or control dependencies with s . This is expected to help our model recognize better the vulnerable code appearing in specific surrounding contexts, and have better discriminating the vulnerable code from the benign one. For example, the existing approaches feed the entire PDG of the method in Figure 2 into a model. IVDetect distinguishes and learns the vector representation for the vulnerable statement at line 27 while considering as contexts the statements with data/control dependencies with line 27: the data-dependency context (lines 31, 22, 13, 10, and 6), and the control-dependency context (lines 29, 25, 23, and 13).

2.2.3 Graph-based Interpretation Model for Vulnerability Detection (Section 4). After prediction, IVDetect performs fine-grained interpretation. It uses both the PDG G_M of the method M and the GCN model as the input to obtain the interpretation. To that end, we leverage the interpretable ML technique *GNExplainer* [36]. Its goal is to take the GCN and a specific input graph G_M , and produce the *crucial sub-graph structures and features* in G_M that affect the decision of the model. GNExplainer’s idea is that *if removing or altering a node/feature does affect the prediction outcome, the node/feature is considered as essential and thus must be included*

in the *crucial set* (let us call it the *interpretation set*). GNNExplainer searches for a sub-graph \mathcal{G}_M in G_M that minimizes the difference in the prediction scores between using the whole graph G_M and using the minimal graph \mathcal{G}_M (Section 4). Because without that subgraph \mathcal{G}_M in the input PDG G_M , GCN model would not decide G_M as vulnerable, \mathcal{G}_M is considered as crucial **PDG sub-graph** consisting of **crucial statements** and data/control dependencies relevant to the detected vulnerability (if the outcome is \mathcal{V}). If the outcome is non-vulnerability, \mathcal{G}_M can be considered as the safe statements in PDG for the model to decide the input method M as benign code.

3 GRAPH-BASED VULNERABILITY DETECTION MODEL

3.1 Representation Learning

Let us present how we build the vector representations for code features. For a *statement*, we extract the following types of **features**:

1. Sequence of Sub-tokens of a Statement. At the lexical level, we capture the content of a statement in term of the sequence of sub-tokens. We choose the sub-token granularity because the sub-tokens are more likely to be repeated than the entire lexical tokens in source code [31]. We tokenize each statement and keep only the variables, method and class names. The names are broken into sub-tokens using CamelCase or Hungarian convention. We remove the sub-tokens with one character to avoid the influence of noises. For example, in Figure 4, the tokens of S_{27} are collected and broken down into the sequence: copy, to, user, arg, etc. Then, we use GloVe [25], to build the vectors for tokens, together with Gate Recurrent Unit (GRU) [12] to build the feature vector for the sequence of sub-tokens for S_{27} . GloVe is known to capture well semantic similarity among tokens. GRU is chosen to summarize the sequence of vectors into one feature vector for the next step.

2. Code Structure of a Statement. We capture code structure via the AST sub-tree. In Figure 4, the AST sub-tree for S_{27} is extracted and fed to Tree-LSTM [32] to capture the structure into a vector F_2 .

3. Variables and Types. For each node (i.e., a statement), we collect the names of the variables and their static types at their locations, break them into the sub-tokens. For example, we collect the variable `s_cmd` and its static type `cross_ec_command`. We use the same vector building techniques as for the sub-token sequences as in feature 1, including GloVe and GRU, to apply on the sequences of sub-tokens built from the variables' names (e.g., `s_cmd`) and those from the variables' types (e.g., `cross_ec_command`).

4. Surrounding Contexts. During training, for a statement s , we also encode the statements surrounding s , which we refer to as *context*. Data- and Control-dependency contexts contain the statements having such dependencies with the current statement. For example, the data-dependency context for S_{27} includes the statements at the lines 31, 22, 13, 10, and 6. If the control dependencies are considered, the statements with control dependencies with S_{27} at the lines 29, 25, 23, and 13 are included. The vectors for the statements in the context are calculated via GloVe and GRU as described earlier. The number of dependencies could be different, then the lengths of the GRU model inputs could be different. Therefore, we apply zero padding with a masking layer, which allows the model

to skip the zeros at the end of the sequence of sub-tokens. Those zeros will not be included in the training.

5. Attention-based Bidirectional GRU. After having all vectors for the features F_1, F_2, \dots , we use a bi-directional GRU and an attention layer to learn the weight vector W_i for each feature F_i , based on the hidden states from that model. Then, we compute the weighted vector for each feature by multiplying the original vector for the feature by the weight: $F'_i = W_i.F_i$.

Finally, we need to consider the impacts from the *dependent statements to the current statement in the PDG*. The rationale is that those neighboring statements in the PDG must have the influence on the current statement if one of them is vulnerable. For example, the neighboring statements for S_{27} in the PDG include the statements at lines 6, 22, 25, and 29. Thus, we combine and summarize them into the final feature vector $F_{S_{27}}$ for the statement S_{27} as follows:

$$F_{S_{27}} = \sum_i W_i \text{Concat}(h(F'_i, j)) \quad (1)$$

W_i is the trainable weight for combination; *Concat* is the concatenate layer to link all values into one vector; h is the hidden layer to summarize vector into a value; $i = S_6, S_{22}, S_{25}, S_{27}, S_{29}$; j is feature index. $F_{S_{27}}$ is used in the next step with GCN model for detection.

3.2 Vulnerability Detection with FA-GCN

Figure 5 presents how we use Feature-Attention GCN model (FA-GCN) [29] for detection. The rationale is that FA-GCN can deal well with the graphs with sparse features (not all the statements share the same properties), and potentially noisy features in a PDG. First, we parse the method M into PDG. Similar to CNN using the filter on an image, FA-GCN performs sliding a small window along all the nodes (statements) of the PDG. For example, in Figure 5, the window marked with (A) for the node S_{27} consists of itself and the neighboring statements/nodes $S_6, S_{22}, S_{25},$ and S_{29} . Another window (marked with (B)) is for the node S_{23} , including itself and the neighboring nodes: S_{22} and S_{25} . For each window, FA-GCN generates the feature representation matrix for the statement at the center. For example, for the window centered at S_{27} , it generates the feature vector $F_{S_{27}}$ for S_{27} , using the process explained in Figure 4. From the representation vectors for all statements, FA-GCN uses a join layer to link all these vectors into the Feature Matrix \mathcal{F}_m for method M . A row in \mathcal{F}_m corresponds to a window in PDG.

Next, FA-GCN performs the convolution operation by first calculating the symmetric normalized Laplacian matrix \tilde{A} [17], and then calculating the convolution to generate the representation matrix M_m for the method m . After that, we use the traditional steps as in a CNN model: using a spatial pyramid pooling layer (to normalize the method representation matrix into a uniform size, and reduce its total size), and connecting its output to a fully connected layer to transform the matrix into a vector V_m to represent m . With V_m , we perform classification by using two hidden layers (controlling the length of vectors and output) and a softmax function to produce a prediction score for m . We use those scores as *vulnerability scores to rank the methods* in a project. The decision for m as \mathcal{V} or $\mathcal{N}\mathcal{V}$ is done via a trainable threshold on the prediction score [18, 20].

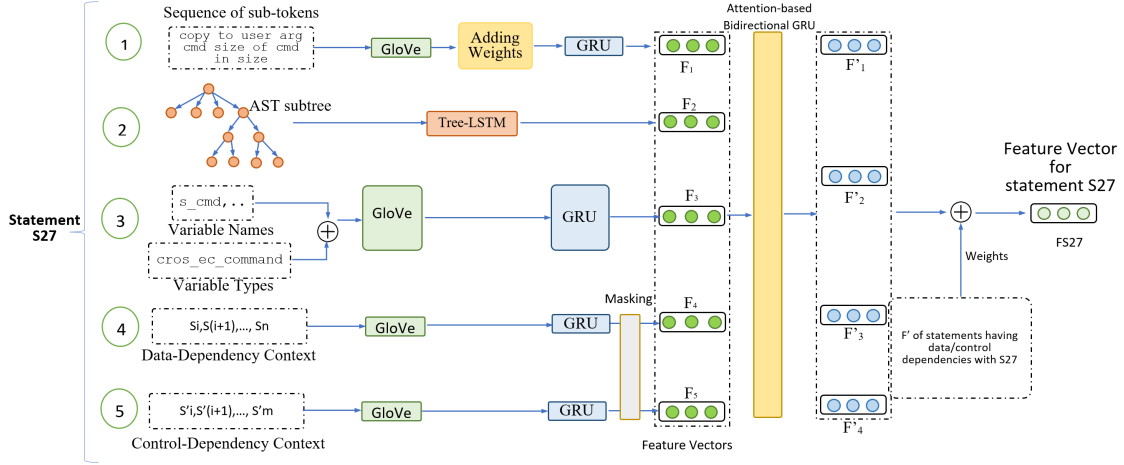


Figure 4: Code Representation Learning for Statement S27 in Graph-based Vulnerable Code Detection

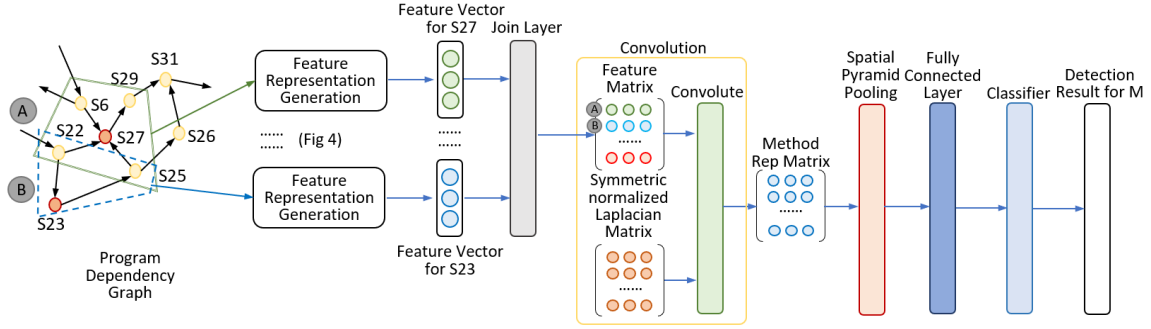


Figure 5: Vulnerability Detection with FA-GCN

4 GRAPH-BASED INTERPRETATION MODEL

Let us explain how we use GNNExplainer [36] to build our graph-based interpretation. The input includes the trained FA-GCN model, the PDG (G_M) of the method M , and the detection result \mathcal{V} or $\mathcal{N}\mathcal{V}$, and prediction score. Figure 6 illustrates our process for the case of \mathcal{V} (Vulnerable) (the case of $\mathcal{N}\mathcal{V}$ is done similarly).

To derive the interpretations, the key goal is to find a sub-graph \mathcal{G}_M in the PDG G_M of the method M that minimizes the difference in the prediction scores between using the entire graph G_M and using the minimal graph \mathcal{G}_M . To do so, we use GNNExplainer with the *masking technique* [36], which treats the searching for the minimal graph \mathcal{G}_M as a learning problem of the *edge-mask* set EM of the edges. The idea is that learning EM helps IVDetect derive the interpretation sub-graph \mathcal{G}_M by masking-out the edges in EM from G_M (“masked-out” is denoted by \odot):

$$\mathcal{G}_M = G_M \odot EM \quad (2)$$

Figure 6 illustrates GNNExplainer’s principle. As an edge-mask set is applied, GNNExplainer checks if the FA-GCN model produces the same result (in this case the result is \mathcal{V}). If yes, the edge in the edge-mask is not important and is not included in \mathcal{G}_M . Otherwise, the edge is important and included in \mathcal{G}_M . Because the numbers

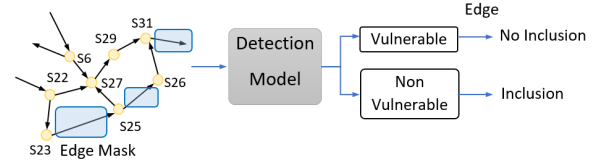


Figure 6: Masking to Derive Interpretation Sub-Graphs

of possible sub-graphs and the edge-mask sets are untractable, GNNExplainer uses a learning approach for the edge-mask EM .

Let us formally explain how GNNExplainer [36] works. It formulates the problem by maximizing the mutual information (MI) between the minimal graph \mathcal{G}_M and the input G_M :

$$\max_{\mathcal{G}_M} MI(Y, \mathcal{G}_M) = H(Y) - H(Y|G = \mathcal{G}_M) \quad (3)$$

Y is the outcome decision by the FA-GCN model. Thus, the entropy term $H(Y)$ is constant for the trained FA-GCN model. Maximizing the MI value for all \mathcal{G}_M is equivalent to minimizing conditional entropy $H(Y|G = \mathcal{G}_M)$, which by definition of conditional entropy can be expressed as

$$- \mathbb{E}_{Y|\mathcal{G}_M} [\log P_{FA-GCN}(Y|G = \mathcal{G}_M)] \quad (4)$$

Table 1: Three Datasets

Dataset	Fan	Reveal	Devign
Vulnerabilities	10,547	1,664	10,067
Non-vulnerabilities	168752	16505	12,294
Ratio (Vul:Non-vul)	1:16	1:9.9	1:1.2

The meaning of this conditional entropy formula is a measure of how much uncertainty remains about the outcome Y when we know $G = \mathcal{G}_M$. GNNExplainer also limits the size of \mathcal{G}_M by K_M , i.e., taking K_M edges that give the highest mutual information with the prediction outcome Y . Direct optimization of the formula 4 is not tractable, thus, GNNExplainer treats \mathcal{G}_M as a random graph variable \mathcal{G} . The objective in Equation 4 becomes:

$$\min_{\mathcal{G}} \mathbb{E}_{\mathcal{G}_M \sim \mathcal{G}} H(Y|G = \mathcal{G}_M) \quad (5)$$

$$\min_{\mathcal{G}} H(Y|G = \mathbb{E}_{\mathcal{G}}[\mathcal{G}_M]) \quad (6)$$

From Equation 5, we obtain Equation 6 with Jensen’s inequality. The conditional entropy in Equation 6 can be optimized by replacing $\mathbb{E}_{\mathcal{G}}[\mathcal{G}_M]$ to be optimized by masking with EM on the input graph \mathcal{G}_M . Now, we can reduce the problem to learning the mask EM . Details on training can be found in [36]. The resulting sub-graph \mathcal{G}_M is directly used as an interpretation. We can similarly produce the interpretations for the cases of non-vulnerability result.

5 EMPIRICAL EVALUATION

5.1 Research Questions

To evaluate IVDETECT, we seek to answer the following questions:

RQ1. Comparison on the Method-Level Vulnerability Detection (VD). How well does IVDETECT perform in comparison with the state-of-the-art method-level Deep Learning VD approaches?

RQ2. Comparison with other Interpretation Models for Fine-grained VD Interpretation. How well does IVDETECT perform in comparison with the state-of-the-art interpretation models for fine-grained VD interpretation to point out vulnerable statements?

RQ3. Vulnerable Code Patterns and Fixing Patterns. Is IVDETECT useful in detecting vulnerable code patterns and fixes?

RQ4. Sensitivity Analysis for Internal Features. How do internal features affect the overall performance of IVDETECT?

RQ5. Sensitivity Analysis on Training Data. How do different data splitting schemes affect IVDETECT’s performance?

RQ6. Time Complexity. What is time complexity of IVDETECT?

5.2 Datasets

We have conducted our study on three vulnerability datasets including Fan *et al.*’s [13], Reveal [11] and FFMPEG+Qemu [37] (Table 1). Fan *et al.* [13] dataset covers the CWEs from 2002 to 2019 with 21 features for each vulnerability. At the method level, the dataset contains +10K vulnerable methods and fixed code. The Reveal dataset [11] contains +18K methods with 9.16% of the vulnerable ones. The FFMPEG+Qemu dataset has been used in Devign study [37] with +22K data, and 45.0% of the entries are vulnerable.

5.3 Experimental Methodology

RQ1. Comparison on Method-Level DL-based VD Approaches.

Baselines. We compare IVDETECT with the state-of-the-art DL-based vulnerability detection approaches: 1) **VulDeePecker** [20]: a DL-based approach using Bidirectional LSTM on the statements and their data/control dependencies. 2) **Devign** [37]: an DL-based approach that uses GGCN model with Gated Graph Recurrent Layers on the AST, CFG, DFG, and code sequences for graph classification. 3) **SySeVR** [19]: in addition to statements and program dependencies, this approach also uses program slicing and leverages several DL models (LR, MLP, DBN, CNN, LSTM, etc.). 4) **Russell et al.** [27]: This DL approach encodes source code as matrices of code tokens and leverages convolution model with random forest (RF) via ensemble classifier. 5) **Reveal** [11]: This approach uses GGNN, MLP, and with Triplet Loss on graph representations of source code.

Procedure. A dataset contains a number of vulnerable and non-vulnerable methods. We first randomly split all of its vulnerable methods into 80%, 10%, and 10% to be used for training, tuning, and testing, respectively. For training, we add to that 80% part the same number of non-vulnerable methods as the vulnerable ones to obtain the balanced training data. For tuning and testing, we also add the non-vulnerable methods but we use the real ratio between vulnerable and non-vulnerable methods in the original dataset to build tuning/testing data. We use AutoML [21] on all models to automatically tune hyper-parameters on the tuning dataset.

We also performed the evaluation across the datasets. We first trained our model on the combination of two datasets Reveal and FFMPEG+Qemu, which has a balanced number of vulnerable methods and non-vulnerable ones. We then tested the model on Fan dataset, which has a more realistic ratio of vulnerable and non-vulnerable methods. To ensure the model suitable for cross-data evaluation, we also used 20% of Fan dataset for tuning the parameters and performed prediction on the remaining 80%.

Evaluation Metrics: We use the following evaluation metrics.

Mean Average Precision $MAP = \frac{\sum_{q=1}^Q AvgP(q)}{Q}$, with *Average Precision* $AvgP = \sum_{k=1}^n P(k)rel(k)$, where n is the total number of results k is the current rank in the list, $rel(k)$ is an indicator function equaling to 1 if the item at rank k is actually vulnerable, and to zero otherwise. Q is the total number of classification types. It is 1 because we only have two types including vulnerable and non-vulnerable classes, however, we rank all the methods based on their scores (1 indicates vulnerable, and 0 otherwise).

Normalized DCG at rank k , $nDCG_k = \frac{DCG_k}{IDCG_k}$, with *Discounted Cumulative Gain* at rank k , $DCG_k = \sum_{i=1}^k \frac{r_i}{\log_2(i+1)}$; and *Ideal DCG* at rank k $IDCG_k = \sum_{i=1}^{|R_k|} \frac{2^i - 1}{\log_2(i+1)}$; where r_i is the score of the result at position i , and R_k the rank of the actual vulnerable methods (ordered by their scores) in the resulting list up to the position k .

First Ranking (FR) is the rank of the first correctly predicted vulnerable method. **Average ranking (AR)** is the average rank of the correctly predicted vulnerable methods in the top-ranked list.

Accuracy under curve (AUC) is defined as $AUC = P(d(m_1) > d(m_2))$ in which P is the probability, d is the detection model (can be regarded as a binary classifier), m_1 is a randomly chosen positive instance, and m_2 is a randomly chosen negative instance.

Precision (P) is the fraction of relevant instances among the retrieved ones. It is calculated as $Precision = \frac{TP}{TP+FP}$ while TP is

the number of true positives and the FP is the number of false positives.

Recall (R) is the fraction of relevant instances that were retrieved. It is calculated as $Recall = \frac{TP}{TP+FN}$ while TP is the number of true positives and the FN is the number of false negatives.

F score (F) is the harmonic mean of precision and recall. It is calculated as $Fscore = 2 \frac{Precision * Recall}{Precision + Recall}$.

RQ2. Comparison with other Interpretation Models for Fine-grained Interpretation.

Baselines. We compare IVD_{DETECT} with the following interpretation models. 1) **ATT** [36]: This approach is a graph attention network that uses the attention mechanism to evaluate the weights (importance levels) of the edges in the input graph. 2) **GRAD** [36]: This approach is a gradient-based method that computes the gradient of the GNN’s loss function *w.r.t.* the adjacency matrix.

Procedure. Our goal here is to evaluate how well IVD_{DETECT} produces the fine-grained interpretations pointing to vulnerable statements. Thus, to train/test the interpretation model, we need to use the Fan dataset because it contains the vulnerable statements and respective fixes. The other two datasets contain only the vulnerabilities at the method level and no fixes. Therefore, in this RQ2, for the vulnerability prediction part, we used the GCN-FA model that was trained on Reveal and FFMpeg+Qemu and predicted on the Fan dataset. For the methods that are vulnerable, but predicted as non-vulnerable, we considered those cases as incorrect because the resulting interpretations do not make sense for incorrect detection. For the methods that are actually non-vulnerable (regardless of the predictions), we could not use them because the non-vulnerable methods do not have the fixed statements as the ground truth for interpretations. Thus, we use the set of methods that are vulnerable and correctly detected as vulnerable for the evaluation of the interpretation model. Let us use D to denote this set.

For the interpretation, we randomly split D into 80%, 10%, and 10% for training, tuning, and testing. For training, we used *the fixed statements* as the labels for interpretation because those fixed ones were the vulnerable ones. For testing, we compared the relevant statements from the interpretation model against the actual fixed statements. Each method in the testing set and the trained GCN-FA model are the input of the interpretation model in this RQ2.

Evaluation Metrics. Given an interpretation sub-graph \mathcal{G}_M generated from the graph-based interpretation model, we evaluate the accuracy of the interpretation for a model as follows. For a method, if \mathcal{G}_M has an overlap with any statement in the code changes that fix the vulnerability, \mathcal{G}_M is considered as a correct interpretation, i.e., relevant to the VD. We then calculate Accuracy as the ratio between the number of correct interpretations over the total number of interpretations. Because code changes could include addition, deletion, and modification, we further define such overlap as follows.

If one of the statements S in the vulnerable version was *deleted* or *modified* for fixing, and if $\mathcal{G}_M \ni S$, then we consider the interpretation sub-graph \mathcal{G}_M is correct, otherwise incorrect. If one of the statements S' was *added* to the vulnerable version for fixing, we check on the fixed version whether \mathcal{G}_M contains any statement with data or control dependencies with S' , we consider it as correct, otherwise, incorrect. For example, in Fig. 2, \mathcal{G}_M contains the statement S23 with data and control dependencies with one of the

Table 2: RQ1. Top-10 Vulnerability Detection Ranked Results on FFMpeg+Qemu Dataset. 0: incorrect, 1: correct

Top-10 result	1	2	3	4	5	6	7	8	9	10	Total
VulDeePecker	0	0	0	0	0	0	1	0	1	1	3
SySeVR	0	0	0	0	0	1	1	1	0	1	4
Russell <i>et al.</i>	0	0	0	0	1	0	1	0	1	1	4
DeVign	0	0	0	0	1	0	1	1	1	0	4
Reveal	0	0	0	1	0	1	0	1	1	1	5
IVDETECT	1	0	1	1	1	0	1	1	0	0	6

added lines from 17–21. Thus, \mathcal{G}_M is correct. The rationale is that if the interpretation sub-graph \mathcal{G}_M contains some statement relevant to the added statement to fix the vulnerability, that interpretation is useful in pointing out the code relevant to the vulnerability.

We also use Mean First Ranking (MFR), i.e., the mean of the rankings for the first statement that needs to be fixed in the interpretation statements, and Mean Average Ranking (MAR), i.e., the mean of the rankings for all statements to be fixed in the interpretation statements. If a statement to be fixed has not been selected as interpretation, we do not consider it when calculating MFR/MAR.

RQ3. Vulnerable Code Patterns and Fixing Patterns.

Procedure. We use a mining algorithm on the set of interpretation sub-graphs to mine patterns of vulnerable code. We also mine fixing patterns for those vulnerabilities. See details in Section 6.3.

Evaluation Metrics. We counted the identified patterns.

RQ4. Sensitivity Analysis for Features.

Procedure. We first built a base model with only the feature that represents the code as the sequence of tokens. We then built other variants of our model by gradually adding one more feature in Section 3.1 to the base model including the sequence of sub-tokens, AST subtree, variable names, data dependencies, and control dependencies. We measured accuracy for each variant. We used the Fan dataset and the same experiment setting as in RQ1.

Evaluation Metrics. We use the same metrics as in RQ1.

RQ5. Sensitivity Analysis for Training Data. We used different ratios in data splitting for training, tuning, and testing: (80%, 10%, 10%), (70%, 15%, 15%), (60%, 20%, 20%), and (50%, 25%, 25%). We used the same Fan dataset and setting as in RQ1.

Evaluation Metrics. We use the same metrics as in RQ1.

RQ6. Time Complexity Analysis. We measure the actual training and predicting time.

6 EXPERIMENTAL RESULTS

6.1 RQ1. Comparison on Method-Level VD

In Table 2, among the top 10 prediction results, IVD_{DETECT} has the most correct predictions (6 vulnerable methods). The vulnerable methods correctly detected by IVD_{DETECT} are also pushed higher in the top-10 ranked list with 4 correct results out of 5 top results. All other baselines have only 0–1 correct detection in the top-5 list. Importantly, the first rank for IVD_{DETECT} (i.e., the rank of the first correctly detected vulnerable methods) is 1st, while those of the baselines are 4th, 5th, 5th, 6th, and 7th (the bold values in Table 2). Moreover, IVD_{DETECT} can detect 14, 35, and 64 vulnerabilities among top-20, top-50, and top-100 prediction results.

Tables 3, 4, and 5 show the comparison among the approaches on three datasets. IVD_{DETECT} consistently performs better in all the metrics (Table 3). For nDCG@{1,3}, all the baselines get zeros

Table 3: RQ1. Method-Level VD on FFMpeg+Qemu Dataset

	VulDee- Pecker	SySeVR	Russell <i>et al.</i>	Devign	Reveal	IVDETECT
nDCG@1	0	0	0	0	0	1
nDCG@3	0	0	0	0	0	0.63
nDCG@5	0	0	0.43	0.45	0.5	0.65
nDCG@10	0.37	0.44	0.45	0.46	0.5	0.68
nDCG@15	0.45	0.48	0.49	0.52	0.55	0.75
nDCG@20	0.48	0.51	0.54	0.56	0.6	0.82
MAP@1	0	0	0	0	0	1
MAP@3	0	0	0	0	0	0.83
MAP@5	0	0	0.20	0.20	0.25	0.80
MAP@10	0.22	0.31	0.30	0.32	0.38	0.78
MAP@15	0.29	0.33	0.34	0.37	0.41	0.72
MAP@20	0.32	0.35	0.37	0.42	0.45	0.69
FR@1	n/a	n/a	n/a	n/a	n/a	1
FR@3	n/a	n/a	n/a	n/a	n/a	1
FR@5	7	6	5	5	4	1
FR@10	7	6	5	5	4	1
FR@15	7	6	5	5	4	1
FR@20	7	6	5	5	4	1
AR@1	n/a	n/a	n/a	n/a	n/a	1
AR@3	n/a	n/a	n/a	n/a	n/a	2
AR@5	n/a	n/a	5	5	4	3.3
AR@10	8.7	7.8	7.8	7.4	7.4	4.7
AR@15	11.2	10	9.5	10	9.1	7.6
AR@20	13.3	12.1	12.6	12.1	12.4	10.3
AUC	0.68	0.72	0.79	0.77	0.79	0.84

because they did not have correct detections in top-3 results. IVDETECT can improve nDCG@10 from 43%–84% and nDCG@20 from 37%–71% as compared to the baselines. Higher nDCG indicates that IVDETECT achieves the ranking closer to the perfect ranking and the correct vulnerable methods appear higher in the top list.

For MAP scores, IVDETECT relatively improves over the baselines from 105%–255% for top-10 and from 53%–116% for top 20. With higher MAP, IVDETECT has higher precision on average for all the top-ranked positions in the top list. That is, the top-ranked result is highly precise in detecting the vulnerable methods.

IVDETECT also achieves better first ranking (FR) and average ranking (AR). While its best FR is 1 and that of next best performer is 4. For AR@10, a correct vulnerable method is on average ranked by IVDETECT 2.7–4.0 positions higher in the ranked list than by the baselines. Our tool also has relatively higher AUC from 6%–24%.

The comparative results on Fan and Reveal datasets are similar (Tables 4 and 5). In Fan dataset, IVDETECT can improve the nDCG and MAP scores over the baselines by 26%–43%, 50%–170% for top-10, and 21%–475%, 40%–250% for top-20. IVDETECT’s FRs and ARs are better from 2–6 positions and 0.7–2.7 positions for top 10, and 2–13 positions and 1.6–9.1 positions for top 20. In Reveal dataset, the improvements in nDCG, MAP, FR, and AR are 33%–73%, 42%–209%, 2–7 positions, and 0–4 positions for top 10, and 19%–111%, 28%–236%, 2–12 positions, and 1.2–6.2 positions for top 20.

The results on three datasets are different due to the ratio between the vulnerable and non-vulnerable methods. That ratio is 1:16 and 1:9.9 in Fan and Reveal datasets. That number is 1:1.2 in FFMpeg+Qemu dataset, thus, there are more vulnerable methods, and the results are consistently higher across all the models.

Table 6 shows the results of Precision and Recall of our IVDETECT and the baselines. Specifically, IVDETECT has higher precision than all the baselines on three datasets. IVDETECT can improve the Precision by 2.6%–105%. For the Recall, IVDETECT is marginally lower than Reveal on Fan and FFMpeg+Qemu datasets (i.e., 1.4%

Table 4: RQ1. Method-Level VD on Fan Dataset

	VulDee- Pecker	SySeVR	Russell <i>et al.</i>	Devign	Reveal	IVDETECT
nDCG@1	0	0	0	0	0	0
nDCG@3	0	0	0	0	0	0.5
nDCG@5	0	0	0.30	0.33	0.34	0.43
nDCG@10	0	0	0.28	0.30	0.37	0.45
nDCG@15	0.08	0.23	0.31	0.32	0.38	0.46
MAP@1	0	0	0	0	0	0
MAP@3	0	0	0	0	0	0.25
MAP@5	0	0	0.1	0.13	0.18	0.27
MAP@10	0	0	0.12	0.14	0.21	0.28
MAP@15	0.08	0.24	0.14	0.15	0.20	0.28
FR@1	n/a	n/a	n/a	n/a	n/a	n/a
FR@5	n/a	n/a	n/a	n/a	n/a	4
FR@10	n/a	n/a	10	8	6	4
FR@15	n/a	n/a	10	8	6	4
FR@20	19	16	10	8	6	4
AR@1	n/a	n/a	n/a	n/a	n/a	n/a
AR@5	n/a	n/a	n/a	n/a	n/a	4
AR@10	n/a	n/a	10	8	8	7.3
AR@15	n/a	n/a	12	10.5	9.3	8.5
AR@20	19.5	18	13.3	13.3	12	10.4
AUC	0.72	0.81	0.82	0.75	0.82	0.9

Table 5: RQ1. Method-Level VD on Reveal Dataset

	VulDee- Pecker	SySeVR	Russell <i>et al.</i>	Devign	Reveal	IVDETECT
nDCG@1	0	0	0	0	0	0
nDCG@3	0	0	0	0	0	0.63
nDCG@5	0	0	0	0	0.43	0.53
nDCG@10	0	0.30	0.32	0.34	0.39	0.52
nDCG@15	0.26	0.28	0.32	0.39	0.42	0.55
nDCG@20	0.27	0.33	0.35	0.43	0.48	0.57
MAP@1	0	0	0	0	0	0
MAP@3	0	0	0	0	0	0.33
MAP@5	0	0	0	0	0.2	0.37
MAP@10	0	0.11	0.11	0.18	0.24	0.34
MAP@15	0.07	0.12	0.16	0.23	0.25	0.36
MAP@20	0.11	0.15	0.18	0.36	0.29	0.37
FR@1	n/a	n/a	n/a	n/a	n/a	n/a
FR@3	n/a	n/a	n/a	n/a	n/a	3
FR@5	n/a	n/a	n/a	n/a	5	3
FR@10	n/a	10	9	7	5	3
FR@15	15	10	9	7	5	3
FR@20	15	10	9	7	5	3
AR@1	n/a	n/a	n/a	n/a	n/a	n/a
AR@3	n/a	n/a	n/a	n/a	n/a	3
AR@5	n/a	n/a	n/a	n/a	5	4
AR@10	n/a	10	9	8	6	6
AR@15	15	12.5	12	10.5	9.8	9.5
AR@20	18	15.5	13.3	12.7	13	11.8
AUC	0.65	0.76	0.75	0.72	0.74	0.81

Table 6: RQ1. Precision and Recall Results of Method-Level VD on Three Datasets (P: Precision; R: Recall; F: F score)

	FFMPeg+Qemu			Fan			Reveal		
	P	R	F	P	R	F	P	R	F
VulDeePecker	0.49	0.27	0.35	0.12	0.49	0.19	0.19	0.14	0.17
SySeVR	0.50	0.66	0.56	0.15	0.74	0.27	0.24	0.42	0.31
Russell <i>et al.</i>	0.55	0.41	0.45	0.16	0.48	0.24	0.26	0.12	0.16
Devign	0.52	0.63	0.57	0.18	0.52	0.26	0.33	0.32	0.32
Reveal	0.55	0.73	0.62	0.19	0.74	0.30	0.31	0.58	0.40
IVDETECT	0.60	0.72	0.65	0.23	0.72	0.35	0.39	0.52	0.45

and 2.7%) and SySeVR on Fan dataset (i.e., 2.7%). On the Reveal Dataset, IVDETECT can improve Reveal by 25.8% in terms of Precision, but decrease Recall by 10.3%. However, in terms of F1 score, IVDETECT can improve the best performed baseline Reveal by 4.8%

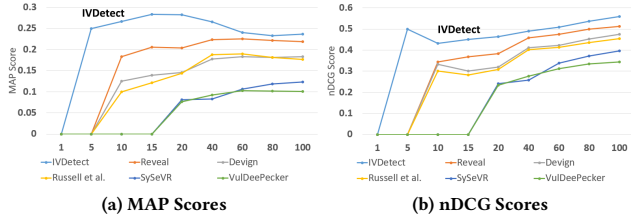


Figure 7: Scores from Top 1 to Top 100 on Fan Dataset

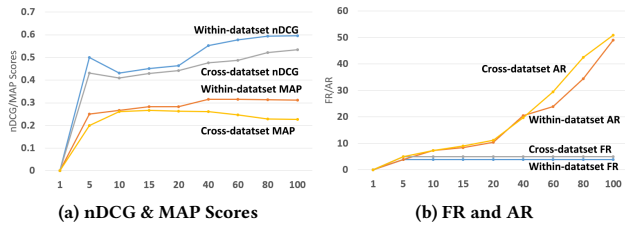


Figure 8: RQ1. Cross-Dataset Validation: Training on Reveal and FFMpeg+Qemu Datasets, testing on Fan Dataset.

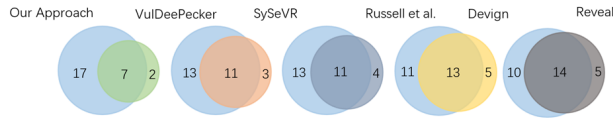


Figure 9: Overlapping Analysis

on FFMpeg+Qemu dataset, 16.7% on Fan dataset, and 12.5% on the Reveal Dataset.

Figure 7 shows that IVDetect consistently has better MAP and nDCG scores when considering top-1 to top-100 ranked lists.

For cross-dataset validation, as seen in Figure 8, the results for MAP and nDCG in within-dataset setting are better than those in cross-dataset setting. This is expected because the model might see similar vulnerable code before in the same projects in the same dataset. The FR and AR values for cross-dataset setting are one rank higher than those of within-dataset setting.

Figure 9 shows our analysis on the overlapping results between IVDetect and the baselines on Fan dataset for top-100. As seen, IVDetect can detect 17, 13, 13, 11, and 10 vulnerable methods that VulDeePecker, SySeVR, Russell, Devign, and Reveal missed, respectively, while they can detect only 2, 3, 4, 5, and 5 vulnerable methods that IVDetect missed. In summary, IVDetect can detect 15, 10, 9, 6, and 5 more vulnerable methods than the baselines.

6.2 RQ2. Comparison with other Interpretation Models for Fine-grained VD Interpretation

Table 7 shows the accuracy of different interpretation models. As seen, using GNNExplainer improves over ATT and GRAD from 12.3%–400% and 9.0%–400% in accuracy, respectively, as we vary the size of interpretation sub-graphs (i.e., the number of statements) from 1–10. Higher accuracy indicates that IVDetect can provide better fine-grained vulnerability detection interpretation at the

Table 7: RQ2. Fine-grained VD Interpretation Comparison

Interp. Model	Accuracy										MFR	MAR
	N1	N2	N3	N4	N5	N6	N7	N8	N9	N10		
ATT	0.01	0.16	0.41	0.54	0.59	0.60	0.62	0.63	0.64	0.65	4.8	6.3
GRAD	0.01	0.19	0.43	0.54	0.59	0.62	0.63	0.65	0.66	0.67	4.2	5.6
GE	0.05	0.30	0.54	0.63	0.67	0.68	0.70	0.72	0.72	0.73	3.5	5.0

GE: GNNExplainer; Nx: x is the number of nodes in the interpretation

Table 8: RQ3. Numbers of Vulnerable Code Patterns

	thres=2	thres=3	thres=4	thres=5
size=2	47	36	22	7
size=3	25	27	19	6
size=4	23	22	16	5
size=5	22	21	11	2
Total	117	102	68	21

statement level. That is, *in more cases, if IVDetect detects correctly vulnerable methods, it can point out more precisely the vulnerable statements* relevant to the vulnerabilities. For ranking vulnerable statements, using GNNExplainer improves MFR by 0.7 and 1.3 ranks, and improves MAR by 0.6 and 1.3 ranks over ATT and GRAD.

ATT uses the edge attention in the Graph Attention Network to assign the weights for the edges, while GNNExplainer directly gives a score for the subgraph after masking. Thus, for the case in which there are more than one path from a node to another, the weight for an edge is the average weight of the weights through multiple paths, i.e., ATT might be less precise than GNNExplainer. GRAD computes the gradient of the loss function with respect to the input for computing the weight of an edge. However, such gradient-based approach may not perform well with respect to the discrete inputs (an input graph is represented as an adjacency matrix).

As the number of nodes in \mathcal{G}_M increases, the number of statements covered also increases, accuracy is higher. However, the computation time is higher and developers need to investigate more statements. As seen, when the number of statements is higher than 5, accuracy increases more slowly. Thus, we chose 5 as a default.

6.3 RQ3. Vulnerable Code Pattern Analysis

This section describes another experiment that we exploit IVDetect’s capability of providing interpretation sub-graphs to mine the patterns of vulnerable code. A vulnerable code pattern is a fragment of vulnerable code that repeats frequently, i.e., more than a certain threshold. The detected vulnerability patterns and corresponding fixes can be the good sources for developers to learn about the vulnerable code that others have frequently made, and learn to fix vulnerable code in the same patterns.

From the results in RQ2, we first collected into a set \mathcal{G} the interpretation sub-graphs \mathcal{G}_M s with the correctly detected statements as relevant to the vulnerability in the methods. In total, we obtain +700 \mathcal{G}_M s. Note that \mathcal{G}_M is a sub-graph of PDG. For each \mathcal{G}_M , we abstract out the variables’ names with a keyword VAR, and the literals with their data types. We then ran the sub-graph pattern mining algorithm [24] on \mathcal{G} with different thresholds of frequencies and collected different sizes of the sub-graph patterns. The outputs are the frequent isomorphic sub-graphs within \mathcal{G}_M s, which

```

1 // =====PATTERN 1 =====
2 if (is_link(StringLiteral) {
3     fprintf(stderr, "Error: invalid /etc/skel/.zshrc file\n"); // not in pattern
4     exit(INTLITERAL);
5 }
6 if (copy_file(StringLiteral, VAR) == INTLITERAL) { ...
7 // =====PATTERN 2 =====
8 VAR = udf_get_filename(VAR, VAR, VAR, VAR);
9 if (VAR && ...) goto LABEL;

```

Figure 10: Vulnerable Code Patterns

```

1 // ===== FIXING PATTERN 1 =====
2 - VAR = f16_update_dst(VAR, VAR, VAR);
3 + rcu_read_lock();
4 + final_p = f16_update_dst(VAR, rcu_dereference(VAR), VAR);
5 + rcu_read_unlock();
6 // ===== FIXING PATTERN 2 =====
7 - char VAR = malloc (VAR);
8 + char VAR;
9 + if (VAR < 0 || VAR > LITCONST) {
10 +     error_line (StringLiteral, VAR);
11 +     return LITCONST;
12 + }
13 + VAR = malloc (VAR);

```

Figure 11: Fixing Patterns (-: removal, +: addition)

are considered as vulnerable code patterns because we chose \mathcal{G}_M that contains correct interpretation statements relevant to the correctly detected vulnerabilities. After manual verification, we obtain a number of correct patterns (Table 8). As seen, as the frequency threshold or the size of pattern is larger, the number of patterns decreases as expected. When they are both larger than 5, we found no pattern. Let us explain a few examples.

Figure 10 shows two examples of vulnerable code patterns. The first pattern (lines 2,4, and 6) shows an API misuse in the project firejail involving `is_link(...)`, `exit`, and `copy_file(...)`. The usage is to check the validity of a link, and if yes to copy the file, or otherwise to stop the execution. This pattern appeared three times with different string literals and was fixed by developers to replace the statements. An interesting observation is that `IVDETECT` is able to eliminate the `fprintf` statement at line 2 from the interpretation sub-graph, thus, eliminating it from the pattern, even though the `fprintf` statement appears with the other statements three times in the project. This shows a benefit of `IVDETECT` because if a tool does not have statement-level VD interpretation and it mines pattern from the entire methods, it will incorrectly include `fprintf` in the pattern. The second pattern (lines 8–9) shows a pattern involving a vulnerable method call `udf_get_filename`, and the checking on its return value. The method later was fixed to add the 5th parameter.

Another interesting finding is that `IVDETECT` enables the discovery of not only vulnerable code patterns but also the fixing patterns for them. Figure 11 shows two fixing patterns for vulnerable code. The first vulnerability (from Linux kernel), lines 2–5, is about the method `f16_update_dst(...)`. According to the commit log, to avoid another thread changing a data record concurrently, developers need to provide mutual exclusion access and deferencing. This fixing pattern was repeated 3 times in the methods `dcpp_v6_send_response`, `inet6_csk_route_req`, and `net6_csk_route_socket`. This fixing pattern would be useful for a developer to learn the fix from one method and apply to the other two methods. The second pattern (lines 7–13) shows a fixing pattern to a vulnerability on buffer overflow with the `malloc` call in `ParseDsdiffHeaderConfig` method of `WavPack 5.0`. According to CVE-2018-7253, this problem “allows a remote attacker

Table 9: RQ4. Evaluation for the Impact of Internal Features.

	ST (A)	(A)+SST (B)	(B)+AST (C)	(C)+Var (D)	(D)+CD (E)	(E)+DD (F)
nDCG@15	0.25	0.27	0.29	0.35	0.42	0.45
nDCG@20	0.26	0.27	0.29	0.37	0.44	0.46
MAP@15	0.07	0.11	0.12	0.19	0.26	0.28
MAP@20	0.09	0.11	0.13	0.19	0.26	0.28
FR@15	14	12	11	7	5	4
FR@20	14	12	11	7	5	4
AR@15	14	13.5	11	10.3	9	8.5
AR@20	19.5	15	13.5	12.5	11.2	10.4
AUC	0.75	0.76	0.77	0.83	0.85	0.9

ST: sequence of tokens; SST: sequence of sub-tokens; AST: sub-AST; Var: variables; CD: control dependencies; DD: data dependencies; F = `IVDETECT`

to cause a denial-of-service (heap-based buffer over-read) or possibly overwrite the heap via a maliciously crafted DSDIFF file”. This fixing pattern occurred three times in the same project.

6.4 RQ4. Sensitivity Analysis for Features

Table 9 shows the changes to the metrics as we incrementally added each internal feature into our model in Figure 4. Generally, each internal feature contributes positively to the better performance of `IVDETECT`, as both the score metrics (nDCG, MAP, and AUC) and the ranking metrics (FR and AR) are improved.

When `IVDETECT` considers only the sequence of tokens (ST) in the code, the first correct detection (FR) is at the position 14, thus, $nDCG@_{\{1,5,10\}}=0$ and $MAP@_{\{1,5,10\}}=0$ (not shown). When considering the code as the sequence of sub-tokens (SST), `IVDETECT` deals with the unique tokens better because the sub-tokens appear more frequently than the tokens [31]. At top-20, FR improves 2 positions, AR improves 4.5 positions, and nDCG and MAP relatively improve 3.8% and 22.2%. When AST is additionally considered, the model can distinguish vulnerable code structures and statements. At top-20, FR and AR improve 1 and 1.5 positions, and nDCG and MAP improve 7.4% and 18.1%. However, FR is still 11 and $nDCG@_{\{1,5,10\}}=0$ and $MAP@_{\{1,5,10\}}=0$ (not shown), because tokens and AST do not help much discriminate the vulnerable statements.

The feature on variables also helps improve FR and AR from 11 to 7 and 13.5 to 12.5, and nDCG and MAP relatively improve 27.6% and 46.2% at top 20. nDCG@10 and MAP@10 improve from 0 to 0.33 and to 0.18, respectively (not shown). This feature allows the model to detect similar incorrect variable usages. By additionally integrating control dependencies (CD), FR and AR improve from 7 down to 5 and 12.5 down to 11.2, and nDCG and MAP relatively improve 18.9% and 36.8%. By adding data dependencies (DD), FR and AR improve from 5 to 4 and 11.2 to 10.4. nDCG and MAP improve 4.5% and 7.7% for top 20. This result confirms that vulnerable code often involves the statements with control and/or data dependencies [11, 37].

Figure 12 shows a detected vulnerable method: `validate_event(...)` was vulnerable and replaced with a new version with an additional parameter. We used the models (A)–(F) for detection, and observed that the rank for `validate_event(...)` in the candidate list improves from 140 (A), to 121 (B), 99 (C), 71 (D), 48 (E), and 19 (F). While the features on tokens, sub-tokens, and AST are contributing, they do not help much because the model did not see them in vulnerable methods before. However, the variable/method names, especially control/data dependencies between the surrounding statements and `validate_event(...)` help discriminate this vulnerability, and push it to

```

1  static int validate_group(struct perf_event *event)
2  {
3      ...
4      - if (!validate_event(&fake_pmu, leader))
5      + if (!validate_event(event->pmu, &fake_pmu, leader))
6          return -EINVAL;
7
8      list_for_each_entry(sibling, &leader->sibling_list, group_entry) {
9          - if (!validate_event(&fake_pmu, sibling))
10         + if (!validate_event(event->pmu, &fake_pmu, sibling))
11             return -EINVAL;
12     }
13
14     - if (!validate_event(&fake_pmu, event))
15     + if (!validate_event(event->pmu, &fake_pmu, event))
16         return -EINVAL; ...
17 }

```

Figure 12: A Detected Vulnerable Method in Android kernel

Table 10: RQ5. Sensitivity Analysis on Training Data

Train/Tune/Test	nDCG@20	MAP@20	FR@20	AR@20	AUC
40%/30%/30%	0.26	0.09	12	15.5	0.69
50%/25%/25%	0.33	0.16	8	12.3	0.74
60%/20%/20%	0.43	0.25	5	11.6	0.85
70%/15%/15%	0.44	0.26	5	11.2	0.87
80%/10%/10%	0.46	0.28	4	10.4	0.9

the top-20 list. Control dependencies (e.g., between `validate_event(...)` and `return -EINVAL`) help improve 29 ranks. Generally, the improvement in ranking shows the positive contributions of all the features.

This example also shows a fixing pattern appearing three times with different variables `leader`, `sibling`, and `event`.

6.5 RQ5. Sensitivity Analysis on Training Data

As seen in Table 10, with more training data, the performance is better as expected. Even with 60%/20%/20%, IVDetect still achieves nCDG of 0.43 and MAP of 0.25, which are still higher than those of the other baselines for top 20 (highest nDCG and MAP of the baselines are 0.38 and 0.20). With 20% less training data (60% vs 80%), IVDetect only drops AUC by 5.5%.

Time Complexity. To generate the interpretation sub-graphs for all methods, it takes about 9 days, 2 days, and 3 days to finish on Fan, Reveal, and FFMpeg+Qemu datasets, respectively. It took 23, 7, 10 hours to train IVDetect on Fan, Reveal, and FFMpeg+Qemu datasets. For VD prediction, it takes only 1-2s per method.

Threats to Validity. We only tested on the vulnerabilities in C and C++ code. In principle, IVDetect can apply to other programming languages. We tried our best to tune the baselines on same dataset for fair comparisons. We focus only on DL-based VD models.

7 RELATED WORK

Various techniques have been developed to detect vulnerabilities. The rule-based approaches were developed to leverage known vulnerability patterns to discover possible vulnerable code, such as FlawFinder [7], RATS [9], ITS4 [33], Checkmarx [1], Fortify [8] and Coverity [2]. Typically, the patterns are manually defined by human experts. The state-of-the-art vulnerability detection tools using static analysis provide the rules for each vulnerability type.

Another type is machine learning (ML)-based or metrics-based. Typically, these approaches require the human-crafted or summarized metrics as features to characterize vulnerabilities and train machine learning models on the defined features to predict whether a given code is vulnerable or not. Various ML-based approaches have

been built on top of distinct metrics, such as terms and their occurrence frequencies [28], imports and function calls [23], complexity, code churn, and developer activity [30], dependency relation [22], API symbols and subtrees [34, 35].

Recently, deep learning (DL) has been applied to detect vulnerabilities. For example, some approaches train a DL model on different code representations to detect vulnerabilities, such as the lexical representations of functions in a synthetic codebase [14], code snippets related to API calls to detect two types of vulnerabilities [20], syntax-based, semantics-based, and vector representations [19], graph-based representations [37]. None of them is designed to provide interpretations for a model in term of vulnerable statements.

8 CONCLUSION

We present IVDetect, a novel DL-based approach to provide sub-graphs in PDG, that explains the prediction results of graph-based vulnerability detection. Our empirical evaluation on vulnerability databases shows that IVDetect outperforms the existing DL-based approaches by 64%–122% and 105%–255% in top-10 nDCG and MAP ranking scores.

Our key **limitations** include 1) un-seen vulnerabilities, 2) the vulnerable statements incorrectly identified due to data/control dependencies with vulnerable ones, 3) missed vulnerable statements due to multiple edges of data/control dependencies.

With IVDetect being a ML/DL-based vulnerability detection model, we aim to raise the level of ML/DL-based approaches, which are not able to point out the statements that caused the model to predict the vulnerability. Thus, we compared IVDetect with the detection approaches of the same category, rather than with static-analysis tools. In the future, we plan to compare IVDetect with static analysis tools.

ACKNOWLEDGMENTS

This work was supported in part by the US National Science Foundation (NSF) grants CCF-1723215, CCF-1723432, TWC-1723198, CCF-1518897, and CNS-1513263.

REFERENCES

- [1] [n.d.]. *Checkmarx*. <https://www.checkmarx.com/>
- [2] [n.d.]. *Coverity*. <https://scan.coverity.com/>
- [3] [n.d.]. *CWE-120: Buffer Overflow*. <https://cwe.mitre.org/data/definitions/120.html>
- [4] [n.d.]. *CWE-290: Authentication Bypass by Spoofing*. <https://cwe.mitre.org/data/definitions/290.html>
- [5] [n.d.]. *CWE-79: Cross-site Scripting*. <http://cwe.mitre.org/data/definitions/79.html>
- [6] [n.d.]. *CWE-89: SQL Injection*. <https://cwe.mitre.org/data/definitions/89.html>
- [7] [n.d.]. *FlawFinder*. <http://www.dwheeler.com/FlawFinder>
- [8] [n.d.]. *HP Fortify*. <https://www.hpford.com/>
- [9] [n.d.]. *RATS: Rough Audit Tool for Security*. <https://code.google.com/archive/p/rough-auditing-tool-for-security/>
- [10] 2021. *The GitHub Repository for This Study*. <https://github.com/vulnerabilitydetection/VulnerabilityDetectionResearch>
- [11] Saikat Chakraborty, Rahul Krishna, Yangruibo Ding, and Baishakhi Ray. 2020. Deep Learning based Vulnerability Detection: Are We There Yet? *arXiv preprint arXiv:2009.07235* (2020).
- [12] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. 2014. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555* (2014).
- [13] Jiahao Fan, Yi Li, Shaohua Wang, and Tien Nguyen. 2020. A C/C++ Code Vulnerability Dataset with Code Changes and CVE Summaries. In *The 2020 International Conference on Mining Software Repositories (MSR)*. IEEE.
- [14] Jacob Harer, Onur Ozdemir, Tomo Lazovich, Christopher Reale, Rebecca Russell, Louis Kim, et al. 2018. Learning to repair software vulnerabilities with generative

- adversarial networks. In *Advances in Neural Information Processing Systems*. 7933–7943.
- [15] Jacob A Harer, Louis Y Kim, Rebecca L Russell, Onur Ozdemir, Leonard R Kosta, Akshay Rangamani, Lei H Hamilton, Gabriel I Centeno, Jonathan R Key, Paul M Ellingwood, et al. 2018. Automated software vulnerability detection with machine learning. *arXiv preprint arXiv:1803.04497* (2018).
- [16] Fabian Keller, Lars Grunske, Simon Heiden, Antonio Filieri, Andre van Hoorn, and David Lo. 2017. A critical evaluation of spectrum-based fault localization techniques on a large-scale software system. In *2017 IEEE International Conference on Software Quality, Reliability and Security (QRS)*. IEEE, 114–125.
- [17] Thomas N. Kipf and Max Welling. 2016. Semi-Supervised Classification with Graph Convolutional Networks. *CoRR abs/1609.02907* (2016). [arXiv:1609.02907](https://arxiv.org/abs/1609.02907) <http://arxiv.org/abs/1609.02907>
- [18] Yi Li, Shaohua Wang, Tien N Nguyen, and Son Van Nguyen. 2019. Improving bug detection via context-based code representation learning and attention-based neural networks. *Proceedings of the ACM on Programming Languages* 3, OOPSLA (2019), 1–30.
- [19] Zhen Li, Deqing Zou, Shouhuai Xu, Hai Jin, Yawei Zhu, and Zhaoxuan Chen. 2018. Sysevr: A framework for using deep learning to detect software vulnerabilities. *arXiv preprint arXiv:1807.06756* (2018).
- [20] Zhen Li, Deqing Zou, Shouhuai Xu, Xinyu Ou, Hai Jin, Sujuan Wang, Zhijun Deng, and Yuyi Zhong. 2018. Vuldeepecker: A deep learning-based system for vulnerability detection. *arXiv preprint arXiv:1801.01681* (2018).
- [21] Microsoft. [n.d.]. Neural Network Intelligence. <https://github.com/microsoft/nni>. Last Accessed August 28th, 2020.
- [22] Stephan Neuhaus and Thomas Zimmermann. 2009. The Beauty and the Beast: Vulnerabilities in Red Hat's Packages.. In *USENIX Annual Technical Conference*.
- [23] Stephan Neuhaus, Thomas Zimmermann, Christian Holler, and Andreas Zeller. 2007. Predicting vulnerable software components. In *Proceedings of the 14th ACM conference on Computer and communications security*. 529–540.
- [24] Tung Thanh Nguyen, Hoan Anh Nguyen, Nam H. Pham, Jafar M. Al-Kofahi, and Tien N. Nguyen. 2009. Graph-Based Mining of Multiple Object Usage Patterns. In *Proceedings of the 7th Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on The Foundations of Software Engineering (Amsterdam, The Netherlands) (ESEC/FSE '09)*. Association for Computing Machinery, New York, NY, USA, 383–392. <https://doi.org/10.1145/1595696.1595767>
- [25] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. GloVe: Global Vectors for Word Representation. In *Empirical Methods in Natural Language Processing (EMNLP)*. 1532–1543. <http://www.aclweb.org/anthology/D14-1162>
- [26] Nam H Pham, Tung Thanh Nguyen, Hoan Anh Nguyen, and Tien N Nguyen. 2010. Detection of recurring software vulnerabilities. In *Proceedings of the IEEE/ACM international conference on Automated software engineering*. 447–456.
- [27] 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 *2018 17th IEEE International Conference on Machine Learning and Applications (ICMLA)*. IEEE, 757–762.
- [28] Riccardo Scandariato, James Walden, Aram Hovsepian, and Wouter Joosen. 2014. Predicting vulnerable software components via text mining. *IEEE Transactions on Software Engineering* 40, 10 (2014), 993–1006.
- [29] Min Shi, Yufei Tang, Xingquan Zhu, and Jianxun Liu. 2019. Feature-attention graph convolutional networks for noise resilient learning. *arXiv preprint arXiv:1912.11755* (2019).
- [30] Yonghee Shin, Andrew Meneely, Laurie Williams, and Jason A Osborne. 2010. Evaluating complexity, code churn, and developer activity metrics as indicators of software vulnerabilities. *IEEE transactions on software engineering* 37, 6 (2010), 772–787.
- [31] Trinh Le Son Nguyen, Hung Dang Phan and Tien N. Nguyen. 2020. Suggesting Natural Method Names to Check Name Consistencies. In *Proceedings of the 42nd International Conference on Software Engineering (ICSE '20)*. ACM Press, 12 pages.
- [32] Kai Sheng Tai, Richard Socher, and Christopher D Manning. 2015. Improved semantic representations from tree-structured long short-term memory networks. *arXiv preprint arXiv:1503.00075* (2015).
- [33] John Viega, Jon-Thomas Bloch, Yoshi Kohno, and Gary McGraw. 2000. ITS4: A static vulnerability scanner for C and C++ code. In *Proceedings 16th Annual Computer Security Applications Conference (ACSAC'00)*. IEEE, 257–267.
- [34] Fabian Yamaguchi, Felix Lindner, and Konrad Rieck. 2011. Vulnerability extrapolation: Assisted discovery of vulnerabilities using machine learning. In *Proceedings of the 5th USENIX conference on Offensive technologies*. 13–13.
- [35] Fabian Yamaguchi, Markus Lottmann, and Konrad Rieck. 2012. Generalized vulnerability extrapolation using abstract syntax trees. In *Proceedings of the 28th Annual Computer Security Applications Conference*. 359–368.
- [36] Zhitao Ying, Dylan Bourgeois, Jiaxuan You, Marinka Zitnik, and Jure Leskovec. 2019. GNNExplainer: Generating Explanations for Graph Neural Networks. In *Advances in Neural Information Processing Systems 32*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.). Curran Associates, Inc., 9244–9255.
- [37] Yaqin Zhou, Shangqing Liu, Jingkai Siow, Xiaoning Du, and Yang Liu. 2019. Devign: Effective vulnerability identification by learning comprehensive program semantics via graph neural networks. In *Advances in Neural Information Processing Systems*. 10197–10207.