

# On the Effectiveness of Function-Level Vulnerability Detectors for Inter-Procedural Vulnerabilities

Zhen Li<sup>\*†</sup>  
School of Cyber Science and  
Engineering, Huazhong University of  
Science and Technology  
Wuhan, China  
zh\_li@hust.edu.cn

Ning Wang<sup>\*</sup>  
School of Cyber Science and  
Engineering, Huazhong University of  
Science and Technology  
Wuhan, China  
wangn@hust.edu.cn

Deqing Zou<sup>\*†‡</sup>  
School of Cyber Science and  
Engineering, Huazhong University of  
Science and Technology  
Wuhan, China  
deqingzou@hust.edu.cn

Yating Li<sup>\*</sup>  
School of Cyber Science and  
Engineering, Huazhong University of  
Science and Technology  
Wuhan, China  
leeyating@hust.edu.cn

Ruqian Zhang<sup>\*</sup>  
School of Cyber Science and  
Engineering, Huazhong University of  
Science and Technology  
Wuhan, China  
ruqianzhang@hust.edu.cn

Shouhuai Xu  
Department of Computer Science,  
University of Colorado Colorado  
Springs, Colorado Springs  
Colorado, USA  
xsu@uccs.edu

Chao Zhang<sup>‡</sup>  
Institute for Network Sciences and  
Cyberspace, Tsinghua University  
Beijing, China  
chaoz@tsinghua.edu.cn

Hai Jin<sup>\*†</sup>  
School of Computer Science and  
Technology, Huazhong University of  
Science and Technology  
Wuhan, China  
hjin@hust.edu.cn

## ABSTRACT

Software vulnerabilities are a major cyber threat and it is important to detect them. One important approach to detecting vulnerabilities is to use deep learning while treating a program function as a whole, known as *function-level* vulnerability detectors. However, the limitation of this approach is not understood. In this paper, we investigate its limitation in detecting one class of vulnerabilities known as *inter-procedural vulnerabilities*, where the *to-be-patched statements* and the *vulnerability-triggering statements* belong to different functions. For this purpose, we create the first *Inter-Procedural Vulnerability Dataset* (InterPVD) based on C/C++ open-source software, and we propose a tool dubbed VulTrigger for identifying vulnerability-triggering statements across functions. Experimental results show that VulTrigger can effectively identify vulnerability-triggering statements and inter-procedural vulnerabilities. Our findings include: (i) inter-procedural vulnerabilities are prevalent with an

average of 2.8 inter-procedural layers; and (ii) function-level vulnerability detectors are much less effective in detecting to-be-patched functions of inter-procedural vulnerabilities than detecting their counterparts of intra-procedural vulnerabilities.

## CCS CONCEPTS

• Security and privacy → Vulnerability scanners.

## KEYWORDS

Vulnerability detection; inter-procedural vulnerability; vulnerability type; patch

## ACM Reference Format:

Zhen Li, Ning Wang, Deqing Zou, Yating Li, Ruqian Zhang, Shouhuai Xu, Chao Zhang, and Hai Jin. 2024. On the Effectiveness of Function-Level Vulnerability Detectors for Inter-Procedural Vulnerabilities. In *2024 IEEE/ACM 46th International Conference on Software Engineering (ICSE '24)*, April 14–20, 2024, Lisbon, Portugal. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3597503.3639218>

## 1 INTRODUCTION

Software vulnerabilities are arguably the most significant cyber threats. For instance, the Apache Log4j2 vulnerability (CVE-2021-44228) can be exploited to launch a remote code execution attack, which has had a huge impact because Apache Log4j2 is widely employed by many enterprises. This highlights the importance of detecting software vulnerabilities. One widely-used approach to detecting software vulnerabilities in source code is to use *Static Application Security Testing* (SAST) tools, including open-source tools [5, 9, 16] and commercial tools [4, 10]. However, the false-positive and false-negative rates of these SAST tools are high because of

<sup>\*</sup>National Engineering Research Center for Big Data Technology and System, Services Computing Technology and System Lab, Hubei Key Laboratory of Distributed System Security, Hubei Engineering Research Center on Big Data Security, Cluster and Grid Computing Lab

<sup>†</sup>JinYinHu Laboratory, Wuhan, China

<sup>‡</sup>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 the author(s) 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](mailto:permissions@acm.org).

ICSE '24, April 14–20, 2024, Lisbon, Portugal

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0217-4/24/04...\$15.00

<https://doi.org/10.1145/3597503.3639218>

their incomplete vulnerability rules or patterns defined by human experts [28, 33]. To avoid these weaknesses, researchers have investigated deep learning-based vulnerability detectors.

One promising approach is to encode a program function as a whole, known as *function-level* vulnerability detectors [3, 7, 11, 14, 22, 39, 44, 49]. This approach can effectively detect *intra-procedural vulnerabilities*, where both the *to-be-patched statements* (that must be patched to eliminate a vulnerability) and the *vulnerability-triggering statements*<sup>1</sup> (that trigger the vulnerability but do not need to be patched) belong to the same function. However, the limitation of this approach to detecting another class of vulnerabilities known as *inter-procedural vulnerabilities*, where the to-be-patched statements and the vulnerability-triggering statements belong to *different* functions, is not understood. It is important to characterize this limitation because their incapability in detecting the vulnerability-triggering statements may cause false-negatives and/or make it hard to explain why the to-be-patched statements (even if detected) represent a vulnerability. In the latter case, practitioners may simply ignore these detected vulnerabilities [13, 22, 50].

**State-of-the-art: Lack of datasets and effective detectors for inter-procedural vulnerabilities.** On one hand, although existing vulnerability datasets [1, 8, 34, 41] do provide vulnerability patches from which to-be-patched statements can be obtained, they do not provide any vulnerability-triggering statements. This explains why studies leveraging such datasets [20, 29, 42] cannot identify inter-procedural vulnerabilities, despite that they can answer the question whether or not multiple functions are involved in a vulnerability patch. This also suggests that these studies could mistakenly claim that a vulnerability involves only one vulnerable function.

On the other hand, we are not aware of any detectors that can effectively identify inter-procedural vulnerabilities. This is true despite that some SAST tools [4, 5, 9, 10, 16] can indeed identify some vulnerability-triggering statements, but their accuracy is too low to be useful because they would miss 47%-80% of vulnerabilities [28] and can achieve at most a 28.8% accuracy in identifying vulnerability-triggering statements (as shown by our study which will be presented in Table 5 and Table 6 in Section 4.2). Dynamic analysis methods (e.g., fuzzing [12, 31, 32]) have high false-negative rates because they cannot test all execution paths despite that they can identify vulnerability-triggering statements. It is also worth mentioning that dynamic analysis tools incur a very high overhead to trigger a vulnerability, meaning that they cannot scale up to deal with a large number of programs.

**Our contributions.** In this paper, we present the first study on the effectiveness of function-level vulnerability detectors in detecting inter-procedural vulnerabilities. Specifically, we make three contributions. First, we propose a novel method and accompanying tool, dubbed VulTrigger, to automatically identify vulnerability-triggering statements for known vulnerabilities with patches. VulTrigger characterizes 10 popular vulnerability types (involving 16 CWEs) and can identify 19 types of vulnerability-triggering statements by leveraging (i) vulnerability patches to obtain *critical variables* and (ii) a newly proposed program slicing method, which may

be of independent value. Experimental results show that VulTrigger significantly outperforms the SAST tools [4, 5, 9, 10, 16], which are the existing tools that have some capabilities in identifying vulnerability-triggering statements at scale, by achieving a 55.8% higher accuracy on average when applied to identify vulnerability-triggering statements. VulTrigger also outperforms the other existing methods [20, 29, 42], which have some capabilities in identifying inter-procedural vulnerabilities, with an improvement of 20.6% in the false positive rate, 32.7% in the false negative rate, and 45.1% in the overall effectiveness F1-measure when applied to identify inter-procedural vulnerabilities.

Second, we apply VulTrigger to build the first *Inter-Procedural Vulnerability Dataset* (InterPVD) based on C/C++ open-source software, which may be of independent value. The dataset involves 769 vulnerabilities, each of which is labeled with its patch statements, its vulnerability-triggering statements, its nature in terms of being inter-procedural or not, its type of inter-procedural vulnerability, and its sequence of functions starting from the to-be-patched function to the vulnerability-triggering function. We find that 24.3% of the vulnerabilities in the InterPVD are inter-procedural vulnerabilities with 2.8 layers on average, meaning that inter-procedural vulnerabilities are prevalent in open-source software.

Third, we investigate the effectiveness of 5 state-of-the-art function-level vulnerability detectors via InterPVD. Experimental results show: (i) function-level vulnerability detectors are much less effective in detecting vulnerability-triggering functions than detecting to-be-patched functions; (ii) detecting to-be-patched functions of inter-procedural vulnerabilities is more challenging than detecting their counterparts of intra-procedural vulnerabilities.

We have published our dataset InterPVD, VulTrigger source code, and other tools we evaluated at <https://github.com/CGCL-codes/VulTrigger>.

## 2 INTER-PROCEDURAL VULNERABILITIES

### 2.1 Vulnerability-Triggering Statements

Figure 1(a) describes an example of inter-procedural vulnerability corresponding to CVE-2015-8662. In this example, `ff_dwt_decode` is the *to-be-patched function* (i.e., the function containing some to-be-patched statement(s)); the *out-of-bounds array access* vulnerability is triggered by the *vulnerability-triggering statement* in the *vulnerability-triggering function*, which is the Line 329 in another function, namely `dwt_decode53`. For this example, the function-level vulnerability detectors may treat `ff_dwt_decode` alone as vulnerable. This contrasts the fact that whether `ff_dwt_decode` is vulnerable or not depends on the presence of Line 329 in `dwt_decode53`. Figure 1(b) describes the diff file for patching the vulnerability, where the three lines of code in green color are called *patch statements*. In this case, the patch is to add control statements. In general, patch statements can be defined as:

**DEFINITION 1 (PATCH STATEMENTS).** *Given a vulnerability  $v$  and its patch via diff file  $d$ , its patch statements are the set of statements  $P$  that are added, deleted, or modified as per  $d$ .*

According to Definition 1, a *to-be-patched function* is a function that must be patched by incorporating one or multiple patch statements. If the patch statement(s) belong to multiple functions, these functions are also to-be-patched functions. Patches can correspond

<sup>1</sup>We do not use the general term of “vulnerable statements” because some researchers treat both *to-be-patched statements* and *vulnerability-triggering statements* as vulnerable statements. It is important to distinguish between these statements in this paper.

```

326 static void dwt_decode53(DWTContext *s, int *t)
327 {
328     int lev;
329     int w = s->linelen[s->ndeclevels - 1][0];
330     int32_t *line = s->_linebuf;
331     ...
373 }
...
596 int ff_dwt_decode(DWTContext *s, void *t)
597 {
598     switch (s->type) {
599     case FF_DWT97:
600     ...
605     case FF_DWT53:
606     dwt_decode53(s, t);
607     break;
610     }
611     return 0;
612 }

```

(a) Two functions related to CVE-2015-8662

```

@@ -595,6 +595,9 @@ int ff_dwt_decode(DWTContext *s, void *t)
int ff_dwt_decode(DWTContext *s, void *t)
{
+ if (s->ndeclevels == 0)
+ return 0;
+
switch (s->type) {
case FF_DWT97:
dwt_decode97_float(s, t);

```

(b) Diff file of CVE-2015-8662

**Figure 1: An example of inter-procedural vulnerability (CVE-2015-8662) with to-be-patched function `ff_dwt_decode`, vulnerability-triggering function `dwt_decode53`, and vulnerability-triggering statement (Line 329)**

**Table 1: Types of patch statements**

No.	Description
P-1	Add/delete variable assignments
P-2	Modify variable assignments
P-3	Add/delete function calls
P-4	Modify function calls
P-5	Add/delete variable definitions
P-6	Modify variable definitions
P-7	Add/delete function definitions
P-8	Modify function definitions
P-9	Add/delete control statements
P-10	Modify control statements
P-11	Others

to adding, deleting, and modifying statements, which may involve variable assignments, system calls, variable definitions, function definitions, and control statements. This leads to the 11 types of patch statements, dubbed P-1 to P-11 and highlighted in Table 1.

**DEFINITION 2 (VULNERABILITY-TRIGGERING STATEMENTS).** *Given a vulnerability  $v$ , the vulnerability-triggering statements of  $v$  are a set of statements  $T$  such that  $v$  is triggered for the first time on an execution path; that is, an incorrect program state is manifested as a consequence of executing the vulnerability-triggering statement.*

According to Definition 2, a *vulnerability-triggering function* is a function that contains one or multiple vulnerability-triggering statements. Note that Definition 2 only considers the *first* time when a vulnerability is triggered on an execution path because the program behaves abnormally after exploitation.

We focus on 10 common vulnerability types (involving 16 CWEs [6]) listed in Table 2, because these vulnerability types can provide more vulnerabilities with diff files which are collected from the *National Vulnerability Database (NVD)* [35]. By considering the scenario that multiple vulnerabilities belong to the same CWE but may be triggered in different fashions and the scenario that multiple vulnerabilities belong to different CWEs but may be triggered in the same fashion, we define 19 types of vulnerability-triggering statements for the above-mentioned 16 CWEs according to the

**Table 2: Vulnerability types and their corresponding types of vulnerability-triggering statements**

Vulnerability type	CWE ID	Types of vulnerability-triggering statements
Buffer overflow	CWE-119	T-1, T-2, T-3
	CWE-125	T-1, T-2, T-3
	CWE-787	T-1, T-2, T-3
	CWE-120	T-1, T-2, T-3
Numeric error	CWE-189	T-1, T-2, T-3, T-4
	CWE-190	T-1, T-2, T-3, T-4
	CWE-191	T-1, T-2, T-3, T-4
Reachable assertion	CWE-617	T-5
Path traversal	CWE-22	T-6
Infinite loop	CWE-835	T-7, T-8, T-9
Missing release of resources	CWE-772	T-10, T-11
	CWE-401	T-10, T-11
Double-free	CWE-415	T-12
Use-after-free	CWE-416	T-12, T-13
NULL pointer dereference	CWE-476	T-14, T-15, T-16
Division-by-zero	CWE-369	T-17, T-18, T-19

**Table 3: Types of vulnerability-triggering statements**

No.	Description
T-1	API or system call-related function call (e.g., <code>memcpy</code> , <code>alloc</code> , <code>memset</code> ) that may trigger out-of-bounds access or memory errors
T-2	Out-of-bounds array access
T-3	Out-of-bounds pointer access
T-4	Integer overflow or underflow caused by integer addition, subtraction and multiplication
T-5	Assertion function calls (e.g., <code>assert</code> and <code>BUG</code> )
T-6	Path access-related API/system calls (e.g., <code>open</code> , <code>read</code> , <code>path_copy</code> , and <code>mkdir</code> )
T-7	Loop condition of <code>for</code> , <code>while</code> , and <code>do while</code>
T-8	<code>goto</code> statement
T-9	Recursive function call
T-10	<code>return</code> statement
T-11	The last statement of the to-be-patched function in the absence of the <code>return</code> statement
T-12	Function calls related to <code>free</code> , <code>delete</code> , <code>destroy</code> , and <code>unregister</code>
T-13	The first usage of a critical variable when there are no function calls related to <code>free</code> , <code>delete</code> , <code>destroy</code> , and <code>unregister</code>
T-14	The first usage of a member of a critical variable, which is a struct type without initialization
T-15	API or system call-related memory allocation function calls (e.g., <code>memcpy</code> and <code>alloc</code> ) when the critical variable is not a variable of struct type and its value is <code>NULL</code> or not assigned
T-16	The uninitialized function is called, mainly in the case of function pointers
T-17	Division-by-zero such as <code>/c</code> and <code>%c</code> where <code>c</code> is a critical variable
T-18	API or system call related to <code>alloc</code>
T-19	The usage of macro definition related to division

vulnerability-triggering fashion. Table 2 shows the correspondence between 16 CWEs corresponding to the 10 vulnerability types and the 19 types of vulnerability-triggering statements described in Table 3. Note that a vulnerability-triggering statement may belong to multiple types of vulnerability-triggering statements. For example, a macro definition related to division (T-19) can be also a division-by-zero statement (T-17).

Some types of vulnerabilities have no obvious vulnerability-triggering statements, in which case we select the statement(s) that are mostly related to the vulnerability as the vulnerability-triggering statements. Take vulnerabilities of the *missing release of resources* type for example, we define the last executed statements in the to-be-patched function as the vulnerability-triggering statements, which often involve the `return` statement or the last statement of the to-be-patched function in the absence of the `return` statement. This is reasonable because the resources should be released before the to-be-patched function ends. In what follows, we elaborate on two examples of vulnerability-triggering statements. **Type 1 (T-1): API or system call-related function call.** For buffer-overflow vulnerabilities, the vulnerability-triggering statement of T-1 is the API or system call-related function call that may trigger out-of-bounds access or memory error, such as `memcpy`, `alloc`, or `memset`. Consider CVE-2013-1929 (CWE-119) in Figure

```

@@ -14604,8 +14604,11 @@ static void tg3_read_vpd (struct tg3 *tp)
if (j + len > block_end)
goto partno;

- memcpy(tp->fw_ver, &vpd_data[j], len);
- strcat(tp->fw_ver, " bc ", vpdlen - len - 1);
+ if (len >= sizeof(tp->fw_ver))
+ len = sizeof(tp->fw_ver) - 1;
+ memset(tp->fw_ver, 0, sizeof(tp->fw_ver));
+ snprintf(tp->fw_ver, sizeof(tp->fw_ver), "%.*s bc ", len,
+         &vpd_data[j]);
}
...

(a) CVE-2013-1929 (CWE-119, T-1)

```

```

@@ -279,6 +279,12 @@ int ParseDsdiffHeaderConfig (FILE *infile, char
*infile_name, char *fourcc, Wavpa
...
else if (!strcmp (dff_chunk_header.ckID, "DSD ", 4)) {
+
+ if (!config->num_channels) {
+ error_line ("%s is not a valid .DFF file!", infile_name);
+ return WAVPACK_SOFT_ERROR;
+ }
+
total_samples = dff_chunk_header.ckDataSize/config->num_channels;
break;
...

(b) CVE-2019-1010315 (CWE-369, T-17)

```

**Figure 2: Diff files of vulnerabilities with vulnerability-triggering statements (highlighted in bold and italics)**

2(a) as an example. The variable `len` exceeds the appropriate length when copying `vpd_data[j]` to `tp->fw_ver` in function `memcpy`, which causes an out-of-bounds access. The vulnerability-triggering statement is “`memcpy(tp->fw_ver, &vpd_data[j], len);`”.

**Type 17 (T-17): division-by-zero.** For division-by-zero vulnerabilities, the vulnerability-triggering statement of T-17 is division-by-zero, such as `/c` and `%c` where `c` is a critical variable (i.e., a variable related to the vulnerability). Take the vulnerability CVE-2019-1010315 (CWE-369) in Figure 2(b) as an example. The value of `config->num_channels` can be zero, causing a division-by-zero error. The vulnerability-triggering statement is “`total_samples=dff_chunk_header.ckDataSize/config->num_channels;`”.

## 2.2 Characterizing Vulnerability-Triggering Statements

To obtain characteristics, we analyze the vulnerability-triggering statements corresponding to the 16 CWEs described in Table 2. We divide the vulnerability types into *two classes* according to whether the vulnerability-triggering statements are closely related to the critical variables. The *first class* of vulnerability types corresponds to the vulnerability-triggering statements that are not closely related to critical variables but do involve *missing release of resource* (CWE-772 and CWE-401) or *infinite loop* (CWE-835). Take characteristics related to missing release of resource as an example. In a to-be-patched function, if there exists a return statement in either the patch statements or their subsequent statements, the vulnerability-triggering statement is the return statement, which is one of the patch statements or the statement closest to the patch statements (T-10). The return statement becomes the last statement executed in the function when the vulnerability is triggered. Otherwise, if there is no return statement, the vulnerability-triggering statement is the last statement in the to-be-patched function (T-11).

The *second class* of vulnerability types corresponds to the vulnerability-triggering statements that depend on the critical variables, involving buffer overflow, numeric error, reachable assertion, path traversal, double free, use after free, NULL pointer dereference, and division-by-zero. Take characteristics related to a buffer overflow as an example. We create a list of keywords of APIs/system

calls, which may trigger out-of-bounds access or memory errors. The first type of vulnerability-triggering statements is the ones in the program slice that involve both the critical variable and the function call containing a keyword in the list of keywords (T-1). If the critical variable is used as an array name or an array index, the vulnerability-triggering statement is the usage of the array involving the critical variable in the program slice (T-2). If the critical variable is a pointer, the vulnerability-triggering statement is the statement in the program slice that involves the pointer that accesses memory or is involved in arithmetic operations (T-3).

## 2.3 Inter-Procedural Vulnerabilities

Based on patch and vulnerability-triggering statements, we define the inter-procedural vulnerability as follows.

**DEFINITION 3 (INTER-PROCEDURAL VULNERABILITY).** Consider a vulnerability  $v$ , its patch statements  $P$ , and vulnerability-triggering statements  $T$ . If a patch statement  $p \in P$  belongs to function  $f$ , a vulnerability-triggering statement  $t \in T$  belongs to another function  $g$  ( $f \neq g$ ), and  $t$  is data- or control-dependent on  $p$ , then we say  $v$  is an inter-procedural vulnerability. The number of inter-procedural layers is the number of different functions involved in the data-flow or control-flow from  $p$  to  $t$ .

It is worth mentioning that for the vulnerability whose patch statements exist in multiple to-be-patched functions (corresponding to one or multiple vulnerability-triggering functions), if there is a to-be-patched function which is different from its corresponding vulnerability-triggering function, the vulnerability is an inter-procedural vulnerability.

According to the relationship between the to-be-patched function  $f$  and the vulnerability-triggering function  $g$ , there are two types of inter-procedural vulnerabilities. The first type is that the to-be-patched function  $f$  directly or indirectly calls the vulnerability-triggering function  $g$ , dubbed “caller type”. Take the vulnerability CVE-2017-13000 (CWE-125) in Figure 3 as an example. The to-be-patched function `ieee802_15_4_if_print` calls function `le64addr_string`, which calls function `loopup_bytestring` at Line 567. The vulnerability-triggering statement `memcmp((const char *)bs, (const char *) (tp->bs_bytes), nlen) == 0` (Line 417) in function `loopup_bytestring` may cause out-of-bounds read when reading `nlen` bytes from `bs` and `tp->bs_byte`. The number of inter-procedural layers is 3, corresponding to three functions `ieee802_15_4_if_print`, `le64addr_string`, and `loop-up_bytestring`.

The second type of inter-procedural vulnerabilities is that the vulnerability-triggering function  $g$  may appear anywhere after the to-be-patched function  $f$  returns, dubbed “callee type”. Take the vulnerability CVE-2018-10878 (CWE-787) in Figure 4 as an example. Function `ext4_fill_super` calls to-be-patched function `ext4_check_descriptors`, then triggers the vulnerability after the to-be-patched function returns, causing out-of-bounds write at Line 4,075. The number of inter-procedural layers is 2, corresponding to functions `ext4_fill_super` and `ext4_check_descriptors`.

## 3 VULTRIGGER

As highlighted in Figure 5, we propose VulTrigger for identifying vulnerability-triggering statements of vulnerabilities. At a high level, VulTrigger takes as input some vulnerability diff files, their

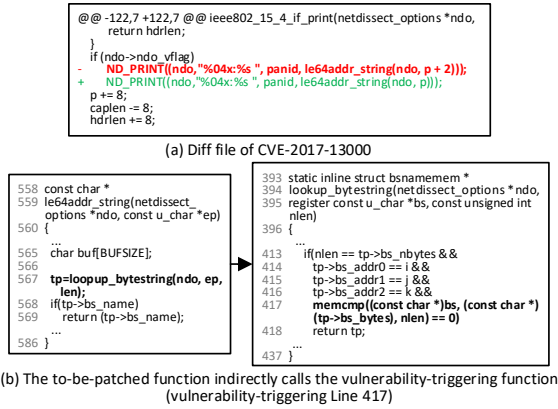


Figure 3: An inter-procedural vulnerability of caller type corresponding to CVE-2017-13000 (CWE-125)

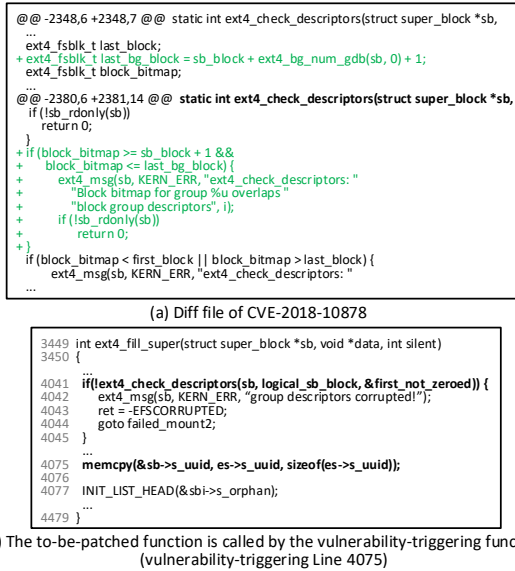


Figure 4: An inter-procedural vulnerability of callee type corresponding to CVE-2018-10878 (CWE-787)

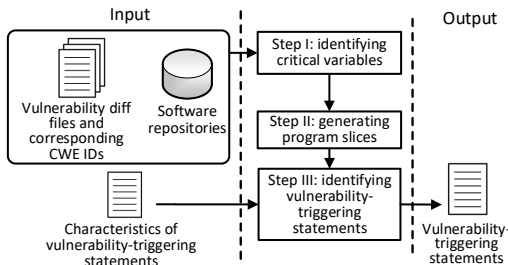


Figure 5: Overview of VulTrigger

corresponding CWE IDs, software repositories, and characteristics of vulnerability-triggering statements; it identifies critical variables based on the patch statements (Step I); then, it generates program slices corresponding to these critical variables (Step II); after that, it identifies the vulnerability-triggering statements (Step III); finally, it outputs vulnerability-triggering statements.

### 3.1 Identifying Critical Variables (Step I)

The purpose of this step is to identify critical variables based on patch statements. Given a vulnerability and its diff file, the starting point for tracing how the vulnerability can be triggered is to identify the *critical variables* related to the vulnerability. We first preprocess the diff file to remove the comments, the blank lines of code, and some common semantically equivalent statement modifications in the patch statements. Then, we identify critical variables for each patch statement according to the types of patch statements as described in Table 1.

**Variable assignments.** For P-1 and P-2 patch statements, critical variables are usually the variables to which values are assigned. For numeric error vulnerabilities, we treat all variables in the assignment statements as critical variables because numeric error vulnerabilities can be triggered at the assigned variables on the left-hand side of the assignment symbol or the arithmetic operations of variables on the right-hand side of the assignment symbol.

**Function calls.** For P-3 patch statements, the critical variables are the argument variables of the function call. For P-4 patch statements, the critical variables are the modified argument variables of the function call.

**Variable definitions.** For P-5 patch statements, the critical variables are the added/deleted variables. For P-6 patch statements, the critical variables are the variables whose definitions are modified.

**Function definitions.** For P-7 patch statements, the critical variables are the argument variables of those functions whose definitions are added/deleted. For P-8 patch statements, the critical variables are the modified argument variables of these functions.

**Control statements.** For P-9 patch statements, the critical variables are the variables in the control conditions. For P-10 patch statements, the critical variables are the modified variables in the control statements.

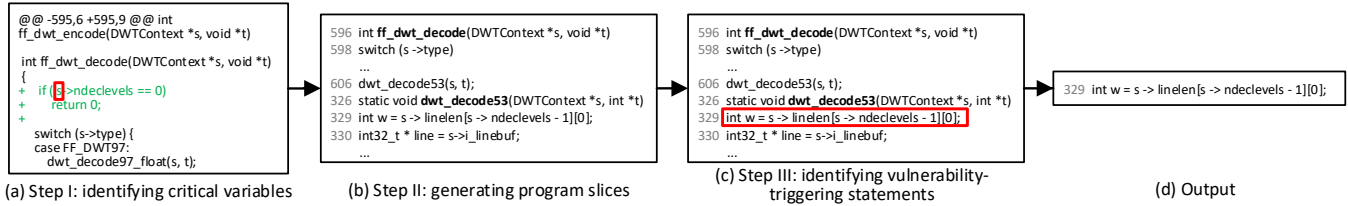
**Others.** For other types of patch statements, we take return statements as an example. If the patch statements only involve a return statement addition/deletion, the critical variables are the variables that are involved in the return statements; otherwise, the critical variables are the modified variables in the return statements.

**Running example.** In any case, if a critical variable is a member of struct type, the struct type variable is considered a critical variable. For the vulnerability (CVE-2015-8662) shown in Figure 1, Figure 6(a) shows the critical variable identified in Step I. Since the patch statements only involve adding control statements, the type of patch statements is P-9. The critical variable is a struct type variable *s*.

### 3.2 Generating Program Slices (Step II)

The purpose of this step is to generate program slices corresponding to the critical variables identified in Step I. To accomplish this, we first extract direct and indirect dependency files of to-be-patched functions (via the system dependency graph of software), which may call the other functions or use the variables or types that are defined in other program files. Then, we generate inter-procedural program slices corresponding to the critical variables as follows. (i) In the case the patch statements in the preprocessed diff file are only about adding statements, we generate program slices for the patched version of the to-be-patched function and the functions that directly or indirectly call the to-be-patched function and that





**Figure 6: Using CVE-2015-8662 in Figure 1 to illustrate the identification of vulnerability-triggering statements, where the two red boxes respectively highlight the identified critical variable and the vulnerability-triggering statement.**

are directly or indirectly called by the to-be-patched function in the patched program. We further delete from each program slice the statements that do not appear in the vulnerable program. (ii) In the case the patch statements in the preprocessed diff file are about deleting or modifying statements, we generate program slices for the to-be-patched function and the functions that directly or indirectly call the to-be-patched function and that are directly or indirectly called by the to-be-patched function in the vulnerable program. Our method for generating program slices incorporates the following three improvements based on the method presented in [25].

**Improvement in dealing with control statements.** It is known that when the patch statement is a control statement, such as `if`, `for`, or `while`, forward slicing according to the control dependency may result in too many statements while forward slicing according to the data dependency may miss some vulnerabilities [25]. Our improvement lies in the following steps. We first obtain the critical variables in the control statement, then conduct the backward slicing till the definition of critical variables, and further conduct forward slicing according to the critical variables with respect to data dependency. Finally, we remove the statements between the statement of the critical variable definition and the control statement from the program slice to exclude the statements that are irrelevant to the vulnerability. This improvement assures that the data-dependent statements related to the critical variable (thus the vulnerability) can be captured in the program slice without introducing unrelated statements.

**Improvement in dealing with implicit return values of function calls.** For method [25], if a function call does not explicitly return values or returns an error, the resulting program slice will end with the function call statement. This is problematic because a pointer variable, when used as an argument in a function call, is an implicit return value, which should have data dependency with the following usage of the pointer variable. To resolve this problem, we identify the pointer variable arguments and conduct backward slicing till the definition of the pointer variable. Then, we conduct forward slicing according to the pointer variable with respect to data dependency. Finally, we remove the statements between the statement of the pointer variable definition and the function call statement from the program slice to exclude statements that are irrelevant to the vulnerability.

**Improvement in dealing with program slicing after the to-be-patched function returns.** The program slicing method [25] can work with the caller type of inter-procedural vulnerability. However, the method [25] stops when the to-be-patched function ends. In this paper, we need to obtain the program slices after the to-be-patched function returns. To achieve this, we generate the

function call graph of the vulnerable program to obtain a set of functions, denoted by  $H$ , which call the to-be-patched function  $f$ . For each function  $h \in H$ , we conduct the inter-procedural program slicing as in [25] by starting at the to-be-patched function call statement. Going beyond the method [25], we let program slicing continue in function  $h$  after the to-be-patched function returns. To obtain the statements of program slices after function  $h$  returns, the process is similar to that of obtaining the program slices after function  $f$  returns. This procedure is iterated until the number of inter-procedural layers reaches a given threshold  $\theta$ .

**Running example.** Corresponding to CVE-2015-8662, Figure 6(b) depicts a program slice corresponding to critical variable  $s$ . The patch statements include control statement `if (s->ndeclevels == 0)`. The program slice is generated by initially extracting statements from the patched version, and then deleting the statements that do not appear in the vulnerable program. The process of program slicing involves a function call without explicit return values, namely “`dwt_decode53 (s, t)`”. The generated program slice corresponding to the vulnerability includes two functions: `ff_dwt_decode` and `dwt_decode53`.

### 3.3 Identifying Vulnerability-Triggering Statements (Step III)

Given a vulnerability and its CWE ID, we identify vulnerability-triggering statements for each program slice generated in Step II. For the first class of vulnerability types (CWE-835, CWE-772, and CWE-401), we identify vulnerability-triggering statements by leveraging characteristics of vulnerability-triggering statements of the CWE ID. For the second class of vulnerability types, there are three substeps. (i) For the critical variables obtained from Step I (denoted by “the first-level critical variables”), we identify the vulnerability-triggering statements in the program slice according to the characteristics of the vulnerability-triggering statements corresponding to the CWE ID. (ii) If there are no vulnerability-triggering statements identified in the previous step, we transform the first-level critical variables to other related critical variables as the second-level critical variables. The transformation strategies are:

- **Assignment.** If a critical variable  $c$  is assigned to another variable  $c'$ , then  $c'$  is the second-level critical variable.
- **Function header in the function definition.** If the statement is a function header in the program slice, we trace back to the place where the function is called. If the critical variable is an argument of the function call, we transform the critical variable to the corresponding parameter variable in the function definition.
- **To-be-patched function call with return value.** If the return value of the to-be-patched function call involves the value of a

critical variable, we transform the critical variable to the variable that is assigned with the return value.

Then, we transform second-level critical variables to third-level critical variables, and so on. The smaller the level a variable has, the higher the priority. The larger-level variables are used to identify vulnerability-triggering statements only when the smaller-level critical variables fail to identify vulnerability-triggering statements. (iii) If no vulnerability-triggering statements are identified in the previous two steps, we identify vulnerability-triggering statements from the statements preceding the patch statements in the program slice according to the characteristics of vulnerability-triggering statements. This situation is assured to be encountered because of loop structures and recursive function calls.

**Running example.** Corresponding to CVE-2015-8662 shown in Figure 1, the vulnerability type is CWE-119. A member of the critical variable `s` is used as the index of an array `s->lineLen`, and the vulnerability-triggering statement is the usage of the array involving the critical variable in the program slice (Line 329, T-2). Figure 6(c) shows the identified vulnerability-triggering statement.

## 4 GENERATING AND ANALYZING INTERPVD

### 4.1 Applying VulTrigger to Construct InterPVD

We focus on 16 CWEs described in Table 2 and collect the open-source C/C++ vulnerabilities from the NVD [35] that meet the following conditions: (i) the software repositories on GitHub should involve the vulnerable version and the patched version corresponding to the diff file of the vulnerability, because VulTrigger depends on the program analysis of these versions; (ii) at least a patch statement of the vulnerability should belong to a function, meaning that the vulnerabilities whose patch statements are all outside of a function are filtered, because we focus on inter-procedural vulnerabilities. We build the InterPVD dataset for C/C++ open-source software, which involves 769 vulnerabilities belonging to 16 CWEs and 53 software products. Each vulnerability is labeled with its patch statements, its vulnerability-triggering statements, its nature in terms of being inter-procedural or not, its type of inter-procedural vulnerability, and its sequence of function calls starting from the to-be-patched function to the vulnerability-triggering function.

Our data collection process has three steps. First, we collect the vulnerabilities whose diff files can be obtained from the patch links described in the NVD [35]. Second, we filter the vulnerabilities whose diff files do not satisfy the heuristics described in [24], and then manually check the remaining ones for their relevance and correctness. Third, for each vulnerability that has multiple different commits, we decide whether to select one commit or merge multiple commits as follows. (i) If the commits are similar, indicating that they are the patches for different versions, we randomly select one commit. (ii) If the commits are not similar, indicating that they are multiple revisions for the same vulnerability, we merge these commits under the condition that the changes described in these commits are in different locations and the context of each commit is the same as the corresponding code in the versions to which the other commits belong; otherwise, we do not consider such a vulnerability because automatically merging these commits is difficult. We apply VulTrigger to identify vulnerability-triggering statements for the 769 vulnerabilities and manually check their correctness and

**Table 4: 12 vulnerabilities in the NVD with incorrect CWE IDs identified by the present study**

No.	CVE ID	CWE ID in the NVD	Correct CWE ID
1	CVE-2009-1897	CWE-119	CWE-476
2	CVE-2009-2767	CWE-119	CWE-476
3	CVE-2009-1298	CWE-119	CWE-476
4	CVE-2014-0205	CWE-119	CWE-416
5	CVE-2015-4002	CWE-119	CWE-835
6	CVE-2019-13295	CWE-125	CWE-369
7	CVE-2019-13297	CWE-125	CWE-369
8	CVE-2009-4307	CWE-189	CWE-369
9	CVE-2013-6367	CWE-189	CWE-369
10	CVE-2015-4003	CWE-189	CWE-369
11	CVE-2016-2070	CWE-189	CWE-369
12	CVE-2018-5816	CWE-190	CWE-369

types by 5 experienced researchers. Each vulnerability is checked by at least two researchers to ensure consistent results. During the process of manual check, we found 12 vulnerabilities to which NVD assigned inaccurate CWE IDs and we corrected them. We reported this to the NVD team and the team replied by stating that they would revisit it. Table 4 lists the vulnerabilities in the NVD with incorrect CWE IDs and their corresponding correct CWE IDs.

### 4.2 Effectiveness of VulTrigger

**Evaluation metrics.** To facilitate the comparison, we adapt the standard metrics [25, 36] to this setting, including *False Positive Rate* (FPR), *False Negative Rate* (FNR), accuracy, precision, and F1-measure (F1). We prefer to achieve low FNR, low FPR, high accuracy, high precision, and high F1.

**Effectiveness on vulnerability-triggering statements.** To evaluate the effectiveness of VulTrigger in identifying the vulnerability-triggering statements with respect to given vulnerability patches, we use *accuracy* to measure the proportion of vulnerabilities whose vulnerability-triggering statements are correctly identified among all vulnerabilities. We consider the vulnerability-triggering statements to be correctly identified if the identified statements include one true vulnerability-triggering statement, while noting that screening the identified vulnerability-triggering statements is left to manual analysis. Given that there are no existing methods for identifying vulnerability-triggering statements with respect to a given patch, we leverage the side-product capabilities of SAST tools in detecting vulnerability-triggering statements. We stress that these tools do not need patches as input because they are designed to detect vulnerabilities, but they have a side-product capability in detecting vulnerability-triggering statements. SAST tools can be divided into two categories according to whether one requires to compile programs or not. For SAST tools that do not require program compilation, we consider Flawfinder [9] and Checkmarx [4]. For SAST tools that require program compilation, we consider Infer [16], CodeQL [5], and Fortify [10]. Since building the compilation environments to accommodate all vulnerabilities incurs a very high labor cost, we only consider the vulnerabilities in the Magma dataset [15], which was originally created for evaluating the effectiveness of fuzzers. We consider all 72 vulnerabilities corresponding to the 16 CWEs described in Table 2 for comparison. For VulTrigger, we use Joern [45] to generate the program dependency graph, set threshold  $\theta = 3$  in the process of generating program slices of critical variables, and remove the vulnerabilities whose programs cannot pass the program analysis of Joern.

**Table 5: Comparing VulTrigger with two SAST tools which do not require program compilation based on 769 vulnerabilities**

Method	#Supported CWEs	Accuracy (%)	#Vulnerability-triggering statements
Ground truth	16	100.0	1.2
Flawfinder [9]	9	9.8	14.3
Checkmarx [4]	16	12.7	114.7
VulTrigger	16	<b>69.8</b>	<b>2.3</b>

**Table 6: Comparing VulTrigger and five SAST tools based on 72 vulnerabilities**

Method	Compiling programs?	#Supported CWEs	Accuracy (%)	#Vulnerability-triggering statements
Ground truth	-	16	100.0	1.1
Flawfinder [9]	No	9	5.9	9.1
Checkmarx [4]	No	16	31.9	275.2
Infer [16]	Yes	16	2.5	3.3
CodeQL [5]	Yes	15	8.6	2.8
Fortify [10]	Yes	13	13.0	31.8
VulTrigger	No	16	<b>66.7</b>	<b>1.2</b>

Table 5 compares VulTrigger and SAST tools which do not require program compilation via the 769 vulnerabilities corresponding to 16 CWEs. We observe that Flawfinder supports 9 CWEs involving 520 vulnerabilities and the other tools support 16 CWEs. We evaluate the accuracy of these SAST tools via the vulnerabilities whose CWEs are supported by them. We also observe that Flawfinder and Checkmarx achieve a very low accuracy for each CWE, leading to an average accuracy of 11.9%. Table 6 compares VulTrigger and SAST tools via the 72 vulnerabilities derived from the Magma dataset [15]. We observe that Flawfinder supports 9 CWEs involving 51 vulnerabilities, that CodeQL supports 15 CWEs involving 58 vulnerabilities, that Fortify supports 13 CWEs involving 49 vulnerabilities, and that the other tools supports 16 CWEs. We further observe that all these SAST tools achieve a very low accuracy. This can be understood as follows: (i) SAST tools have high FNRs in detecting vulnerabilities; and (ii) their goals are not to detect vulnerability-triggering statements of given vulnerabilities. By contrast, VulTrigger leverages the patch statements and achieves on average a 55.8% higher accuracy than that of the SAST tools.

To show how many vulnerability-triggering statements which are identified but need to be filtered by human analysis, we compare the average number of true vulnerability-triggering statements and the average number of vulnerability-triggering statements identified by SAST tools and VulTrigger. For each vulnerability, we focus on the vulnerability-triggering statements identified in the file involving the true vulnerability-triggering function(s). From Table 5 and Table 6, we observe that VulTrigger needs human experts to filter 1.1 (0.1) statements for 769 (72) vulnerabilities on average. Checkmarx demands human experts to filter, on average, 113.5 (274.1) statements for each of 769 (72) vulnerabilities, achieving a higher accuracy at the cost of filtering much more statements by human experts; the other SAST tools demand human experts to filter on average 13.1 (10.7) statements. This can be explained by the fact that SAST tools report many false vulnerable statements.

**Comparing inter-procedural vulnerability identification methods.** Existing studies related to inter-procedural vulnerabilities [20, 29, 42] (dubbed “patch-function method”) focus on whether or not multiple functions are involved in a vulnerability patch. That is, these studies consider a vulnerability inter-procedural as long as its patch involves multiple functions. As shown in Table 7, VulTrigger

**Table 7: Comparing VulTrigger and the existing methods that have some capabilities in detecting inter-procedural vulnerabilities (unit: %)**

Method	FPR	FNR	Accuracy	Precision	F1
Patch-function [20, 29, 42]	23.8	90.9	64.0	7.8	8.4
VulTrigger	<b>3.2</b>	<b>58.2</b>	<b>86.8</b>	<b>74.2</b>	<b>53.5</b>

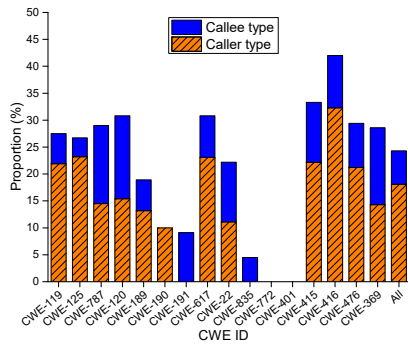
can improve FPR by 20.6%, FNR by 32.7%, accuracy by 22.8%, precision by 66.4%, and F1 by 45.1% in identifying inter-procedural vulnerabilities. Nevertheless, VulTrigger overlooks some vulnerability-triggering statements, which needs to be addressed in future work. We manually analyze these false negatives and summarize our findings as follows. (i) VulTrigger cannot identify implicit semantic associations between variables (e.g., CVE-2014-9664), which causes incorrect transformations of critical variables. (ii) If the diff file (e.g., CVE-2014-9604) only adds checking statements that are not data-dependent on critical variables, the resulting program slice cannot involve the vulnerability-triggering statement, causing the vulnerability-triggering statement to be missed. (iii) For some vulnerabilities (e.g., CVE-2015-5706), false negatives occur because there are no critical variables or the critical variables cannot be extracted correctly in the diff files. (iv) Some inter-procedural vulnerabilities are missed or incorrectly identified because some complex statements are missed or parsed incorrectly by Joern [45], leading to wrong vulnerability-triggering statements.

*INSIGHT 1. VulTrigger substantially outperforms the side-product capabilities of existing SAST tools in identifying vulnerability-triggering statements (while noting that there are no known methods specifically designed to identify vulnerability-triggering statements), by achieving a 55.8% higher accuracy on average in identifying vulnerability-triggering statements, and by improving FPR by 20.6%, FNR by 32.7%, and F1 by 45.1% in identifying inter-procedural vulnerabilities.*

### 4.3 Distribution of Inter-Procedural Vulnerabilities

**Proportion of inter-procedural vulnerabilities.** To see what proportion of different types of inter-procedural vulnerabilities, we consider all the 16 CWEs in our InterPVD dataset. Among them, CWE-119, CWE-125, and CWE-476 are the top-3 CWEs involving 57.5% vulnerabilities and 65.2% inter-procedural vulnerabilities. Figure 7 depicts the proportion of different types of inter-procedural vulnerabilities in different CWEs. We observe that 24.3% vulnerabilities are inter-procedural vulnerabilities, including 18.1% of the caller type and 6.2% of the callee type. We also observe that there are no inter-procedural vulnerabilities for CWE-772 and CWE-401. This can be explained by the fact that the types of vulnerability-triggering statements for CWE-772 and CWE-401 are usually the return statement or the last statement of the to-be-patched function (i.e., T-10 and T-11), where the vulnerability-triggering statements appear in the scope of the to-be-patched function. For CWE-191 and CWE-835, there are only inter-procedural vulnerabilities of the callee type. We speculate that this is caused by the small number of vulnerabilities of CWE-191 and CWE-835 in our dataset, involving 11 vulnerabilities of CWE-191 and 22 vulnerabilities of CWE-835. For most CWEs, we have more inter-procedural vulnerabilities of the caller type, indicating that more vulnerabilities are patched in the direct or indirect calling functions.





**Figure 7: Illustrating the proportion of different types of inter-procedural vulnerabilities corresponding to each CWE**

**Number of inter-procedural layers.** To show how many functions are involved in inter-procedural vulnerabilities, we summarize the number of inter-procedural layers from the vulnerability patch function to the vulnerability-triggering function. If a vulnerability has multiple pairs mapping from the vulnerability patch function to the vulnerability-triggering function, we take them as multiple vulnerability instances. We obtain the average number of inter-procedural layers starting at the vulnerability patch function and ending at the vulnerability-triggering function for all vulnerability instances. Table 8 summarizes the average number of layers of inter-procedural vulnerabilities for each CWE. For each CWE that has both types of inter-procedural vulnerabilities, the inter-procedural vulnerabilities of the callee type have more layers than those of the caller type. This indicates that vulnerability-triggering statements appear more often in the functions that indirectly call the to-be-patched function. The average number of layers of inter-procedural vulnerabilities is 2.8, and the inter-procedural vulnerabilities of the callee type has 0.9 layers more than those of the caller type.

*INSIGHT 2. Among 769 vulnerabilities, 24.3% are inter-procedural vulnerabilities with 2.8 layers on average, including 18.1% of the caller type and 6.2% of the callee type.*

## 5 EFFECTIVENESS OF FUNCTION-LEVEL VULNERABILITIES DETECTORS

**Experimental setup.** To evaluate the effectiveness of function-level vulnerability detectors (trained by existing datasets) in detecting intra-procedural vulnerabilities and inter-procedural vulnerabilities, we consider five function-level vulnerability detectors [3, 11, 14, 44, 49] because they represent the state-of-the-art: VulBERTa [14] and LineVul [11] use Transformer-based pre-training model; Devign [49] and REVEAL [3] identify vulnerabilities by Graph Neural Network-based models; VulCNN [44] converts the graph representation of code to an image while leveraging Convolutional Neural Network-based models.

For the training set, we use the REVEAL dataset [3] because it also represents the state-of-the-art. The REVEAL dataset labels the to-be-patched functions as vulnerable, using the same labeling scheme as other datasets. For each of aforementioned five methods, we apply 10-fold cross-validation to train the model. For test set, we use the InterPVD dataset in which vulnerabilities can be divided into *inter-procedural* and *intra-procedural* vulnerability test sets. The inter-procedural vulnerability test set can be further divided

into *caller* and *callee* test sets based on whether the inter-procedural vulnerabilities are caller or callee types.

The samples in the test set contains vulnerability-triggering functions, to-be-patched functions, and their corresponding patched functions of vulnerabilities, which are labeled as vulnerable or non-vulnerable. We conduct two experimental tasks: detecting to-be-patched functions vs. vulnerability-triggering functions. In the former case, the to-be-patched functions are considered vulnerable but the patched functions and vulnerability-triggering functions are considered non-vulnerable. When detecting vulnerability-triggering functions, vulnerability-triggering functions are considered vulnerable but the patched functions and to-be-patched functions are considered non-vulnerable.

**Experimental results.** To assure fair comparison, we set FNR to 20.0% for each detector by adjusting the probability threshold when determining whether a piece of code is vulnerable or not. This value is chosen because it corresponds to the FNR of the model that achieves the highest F1. Table 9 summarizes the experimental results via the InterPVD test set. We observe that there is little variation in each evaluation metric of the five detectors in detecting to-be-patched functions and vulnerability-triggering functions. We also observe that the five detectors show significantly lower effectiveness in all metrics, except FNR, in detecting vulnerability-triggering functions when compared with detecting to-be-patched functions. The overall effectiveness, F1 score, for detecting vulnerability-triggering functions is, on average, 18.7% lower than that of detecting to-be-patched functions. This can be explained by the fact that to-be-patched functions are labeled as vulnerable in the training set, meaning that the detectors trained from this dataset tend to have lower effectiveness in detecting vulnerability-triggering functions in the test set.

Table 10 summarizes the experimental results based on the inter-procedural and intra-procedural vulnerability test sets. We observe that the effectiveness in detecting to-be-patched functions is higher than detecting vulnerability-triggering functions. This can be attributed to the fact that existing vulnerability datasets only label the to-be-patched functions as vulnerable, implying limited capability in learning features associated with vulnerability-triggering functions. In terms of detecting to-be-patched functions, the effectiveness with respect to inter-procedural vulnerabilities is consistently lower than its counterpart with respect to intra-procedural vulnerabilities. Both intra-procedural and inter-procedural vulnerabilities have some to-be-patched functions, which are labeled as vulnerable in the Reveal dataset, meaning no bias favoring detection of to-be-patched functions. Experimental results indicate that detecting to-be-patched functions of inter-procedural vulnerabilities is more challenging than detecting their counterparts of intra-procedural vulnerabilities.

For inter-procedural vulnerabilities, Table 11 compares the five detectors based on the caller and callee test sets. We observe that vulnerability detectors achieve, on average, a 4.6% higher overall effectiveness F1 in detecting to-be-patched functions for the caller test set than the callee test set, and a 7.2% lower F1 in detecting vulnerability-triggering functions for the caller test set than the callee test set. This could be attributed to the increased ratio of non-vulnerable to vulnerable samples in the callee test set, which leads to more false positives and significantly lowers the overall

**Table 8: Average number of layers of different types of inter-procedural vulnerabilities for 16 CWEs**

Type \ CWE ID	CWE ID																
	119	125	787	120	189	190	191	617	22	835	772	401	415	416	476	369	All
Caller type	2.7	2.2	2.3	3.0	2.6	3.2	-	3.3	2.0	-	-	-	2.0	2.5	2.7	2.7	2.6
Callee type	3.3	2.6	2.6	4.0	3.0	-	3.0	6.0	4.0	3.0	-	-	3.0	3.0	3.9	7.7	3.5
All	2.8	2.3	2.4	3.5	2.7	3.2	3.0	4.0	3.0	3.0	-	-	2.3	2.6	3.0	5.2	2.8

**Table 9: Effectiveness comparison via InterPVD (unit: %)**

Method	FPR	FNR	Accuracy	Precision	F1
Detecting to-be-patched functions					
VulBERTa	80.1	20.0	47.7	46.2	58.6
LineVul	77.8	20.0	48.9	46.9	59.1
Devign	81.8	20.0	46.5	45.3	57.8
REVEAL	76.0	20.0	49.7	47.1	59.3
VulCNN	80.3	20.0	47.2	45.6	58.1
Detecting vulnerability-triggering functions					
VulBERTa	79.9	20.0	35.2	25.3	38.5
LineVul	78.5	20.0	36.3	25.6	38.8
Devign	82.3	20.0	35.1	27.3	40.7
REVEAL	82.3	20.0	35.1	27.3	40.7
VulCNN	80.6	20.0	36.0	27.3	40.7

**Table 10: Effectiveness of five detectors via inter-procedural and intra-procedural vulnerability test sets (unit: %)**

Method	Test set	FPR	FNR	Accuracy	Precision	F1
Detecting to-be-patched functions						
VulBERTa	Inter-procedural	75.8	23.8	45.5	41.1	53.4
	Intra-procedural	81.9	18.8	48.5	47.9	60.3
LineVul	Inter-procedural	75.2	18.7	47.9	42.8	56.1
	Intra-procedural	78.9	20.4	49.3	48.4	60.2
Devign	Inter-procedural	84.3	20.3	41.7	39.2	52.6
	Intra-procedural	80.7	19.9	48.5	47.9	59.9
REVEAL	Inter-procedural	72.9	21.2	48.1	42.5	55.3
	Intra-procedural	77.5	19.6	50.3	49.0	60.9
VulCNN	Inter-procedural	73.7	25.3	46.0	40.9	52.9
	Intra-procedural	83.6	18.1	47.8	47.4	60.0
Detecting vulnerability-triggering functions						
VulBERTa	Inter-procedural	77.7	30.8	31.9	18.6	29.3
	Intra-procedural	80.9	17.0	36.5	27.7	41.5
LineVul	Inter-procedural	81.2	35.9	28.1	16.8	26.6
	Intra-procedural	77.4	15.8	39.3	28.8	42.9
Devign	Inter-procedural	81.6	9.5	34.0	23.5	37.3
	Intra-procedural	82.6	23.0	35.5	28.9	42.1
REVEAL	Inter-procedural	81.1	23.3	31.5	20.8	32.7
	Intra-procedural	82.9	18.9	36.6	29.9	43.7
VulCNN	Inter-procedural	74.1	25.9	36.2	21.3	33.1
	Intra-procedural	83.7	18.2	35.9	29.5	43.4

effectiveness. Specifically, for detecting to-be-patched functions, the ratio of non-vulnerable to vulnerable samples is 1.5:1 in the callee test set and 1.6:1 in the caller test set; for detecting vulnerability-triggering functions, the ratio is 4.5:1 in the callee test set and 3.9:1 in the caller test set.

**INSIGHT 3.** *Function-level vulnerability detectors are much less effective in detecting vulnerability-triggering functions than detecting to-be-patched functions; detecting to-be-patched functions of inter-procedural vulnerabilities is more challenging than detecting their counterparts of intra-procedural vulnerabilities.*

## 6 DISCUSSION

**Use cases of VulTrigger.** Two example use cases of VulTrigger are the following. First, VulTrigger can be used to label a large number of vulnerability-triggering functions, which can be leveraged to improve the effectiveness of function-level vulnerability detectors in coping with vulnerability-triggering functions. Second, VulTrigger can also be used to extract more accurate statements related to inter-procedural vulnerabilities, leading to high-quality program slices. These high-quality program slices, along with their graph

**Table 11: Effectiveness comparison on the caller and callee test sets of inter-procedural vulnerabilities (unit: %)**

Method	Test set	FPR	FNR	Accuracy	Precision	F1
Detecting to-be-patched functions						
VulBERTa	Caller	77.2	20.2	46.3	42.0	55.0
	Callee	69.8	40.5	41.9	36.2	45.1
LineVul	Caller	75.7	16.1	48.8	43.7	57.5
	Callee	73.0	30.9	43.8	38.7	49.6
Devign	Caller	83.3	21.4	42.0	39.5	52.6
	Callee	89.1	14.3	40.0	38.0	52.6
REVEAL	Caller	76.0	19.2	47.3	42.5	55.7
	Callee	58.2	31.4	52.2	42.9	52.8
VulCNN	Caller	74.1	24.6	46.1	41.2	53.3
	Callee	71.9	28.9	45.3	39.7	50.9
Detecting vulnerability-triggering functions						
VulBERTa	Caller	80.3	30.1	29.6	17.7	28.3
	Callee	65.4	33.3	41.9	23.2	34.4
LineVul	Caller	83.8	39.8	25.0	15.1	24.1
	Callee	69.1	20.8	41.9	25.3	38.4
Devign	Caller	80.7	9.8	33.9	22.6	36.1
	Callee	86.4	8.3	34.4	27.9	42.7
REVEAL	Caller	82.4	20.6	30.4	20.1	32.1
	Callee	74.2	33.3	36.7	24.6	36.0
VulCNN	Caller	74.7	25.8	35.4	20.6	32.2
	Callee	70.8	26.1	40.0	25.0	37.4

representations, can be used as input examples to train more effective slice-level vulnerability detectors in detecting inter-procedural vulnerabilities (e.g., [33]).

**Threats to validity.** There are two threats to the validity of our study. First, we use manual analysis to establish the ground truth of vulnerability-triggering statements in C/C++ open-source software. Although each vulnerability is checked by at least two researchers, it is still possible that we miss some vulnerability-triggering statements because our manual check may still miss some execution paths. This is true because identifying the complete set of execution paths and thus the complete set of vulnerability-triggering statements is difficult even by manual analysis. To mitigate the threat, a more capable static analysis tool is needed to accurately identify all execution paths. Second, we create and maintain a list of keywords of APIs/system calls, which may trigger out-of-bounds access or memory errors. If a user-defined function name contains a keyword in the above list and directly or indirectly makes a system call, the user-defined function call is considered the vulnerability-triggering statement. This may reduce the number of inter-procedural layers. **Limitations.** The present study has four limitations. First, we focus on inter-procedural vulnerabilities in C/C++ open-source software. Future studies need to consider other programming languages. Second, we only consider vulnerability-triggering statements related to 16 CWEs. Future studies need to investigate more CWEs. Third, we identify as many vulnerability-triggering statements as possible for a given vulnerability. Future studies need to consider how to identify all vulnerability-triggering statements in all execution paths with respect to a vulnerability. Fourth, future studies need to reduce false negatives of VulTrigger as discussed in Section 4.2 and propose effective inter-procedural vulnerability detection approaches.

## 7 RELATED WORK

**Prior studies on analyzing open-source software vulnerabilities.** These studies are from various perspectives, such as: (i) *distributions* of project-level vulnerabilities [29], of vulnerabilities in crowd-sourced software [38], and of vulnerabilities in cloned vs. non-cloned code [17, 19]; (ii) *patches* [20, 40], patch explanations [27], and security impact of patches [43]; (iii) *characteristics* of project dependencies related to vulnerabilities [21], of features indicating vulnerabilities [48], and of vulnerable code changes [2]. To our knowledge, we are the first to systematically study inter-procedural vulnerabilities of open-source software, despite the previous studies [20, 29, 42] that investigate whether or not multiple functions are involved in a vulnerability patch and cannot accurately identify inter-procedural vulnerabilities (Table 7).

**Prior studies on identifying to-be-patched vs. vulnerability-triggering statements.** (i) To-be-patched statements are the root cause of a vulnerability [28] and can be identified by various methods, such as clone-based detection [18, 24], deep learning-based detection [23], statistical localization enabled by exploits [37]. (ii) Vulnerability-triggering statements can be identified by static and dynamic methods. Static methods are demonstrated by SAST tools [4, 5, 9, 10, 16], which have high false-positive rates and high false-negative rates. For example, the best static analysis tool misses 47%-80% vulnerabilities [28], which is also confirmed by our experiments (cf. Tables 5 and 6). VulTrigger is a static method to identify vulnerability-triggering statements. Dynamic methods (e.g., fuzzing [12, 31, 32]) can identify vulnerability-triggering statements but incur high false-negative rates because it is difficult to test all execution paths. Moreover, they incur a very high overhead when attempting to trigger a vulnerability (i.e., not scalable).

**Prior studies on detecting intra/inter-procedural vulnerabilities.** *Intra-procedural* vulnerability detection methods often take to-be-patched functions and their patched versions to train a detector, while leveraging machine learning [46, 47] or deep learning [3, 7, 11, 14, 22, 39, 44, 49]. *Inter-procedural* vulnerability detection methods can detect vulnerabilities whose to-be-patched statements and vulnerability-triggering statements belong to different functions. These methods extract inter-procedural program slices to train a deep learning model [25, 26, 33] or leverage inter-procedural analyses to detect different types of vulnerabilities [4, 5, 10, 16, 30]. Whereas, VulTrigger determines whether a vulnerability is inter-procedural or not, which allows us to present the first study on the effectiveness of function-level vulnerability detectors in detecting inter-procedural vulnerabilities.

## 8 CONCLUSION

We have investigated the effectiveness of function-level vulnerability detectors in coping with inter-procedural vulnerabilities in C/C++ open-source software, by proposing the innovative VulTrigger to identify vulnerability-triggering statements for known vulnerabilities with patches. We show that inter-procedural vulnerabilities are prevalent in C/C++ open-source software, and that detecting to-be-patched functions and vulnerability-triggering functions, especially the latter, for inter-procedural vulnerabilities is significantly more challenging than for intra-procedural vulnerabilities. This explains why existing function-level vulnerability

detectors cannot effectively cope with inter-procedural vulnerabilities. The limitations discussed in Section 6 offer interesting problems for future research.

## 9 ACKNOWLEDGEMENT

We thank the anonymous reviewers for their comments which guided us in improving the paper. The authors affiliated with Huazhong University of Science and Technology were supported by the National Natural Science Foundation of China under Grant No. 62272187. Shouhuai Xu was partly supported by the National Science Foundation under Grants #2122631, #2115134, and #1910488 as well as Colorado State Bill 18-086. Any opinions, findings, conclusions, or recommendations expressed in this work are those of the authors and do not reflect the views of the funding agencies in any sense.

## REFERENCES

- [1] Guru Prasad Bhandari, Amara Naseer, and Leon Moonen. 2021. 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 (PROMISE)*, Athens Greece. ACM, 30–39.
- [2] Amiangshu Bosu, Jeffrey C. Carver, Munawar Hafiz, Patrick Hillel, and Derek Janni. 2014. Identifying the Characteristics of Vulnerable Code Changes: An Empirical Study. In *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering (FSE)*, Hong Kong, China. ACM, 257–268.
- [3] Saikat Chakraborty, Rahul Krishna, Yangruo Ding, and Baishakhi Ray. 2022. Deep Learning Based Vulnerability Detection: Are We There Yet? *IEEE Transactions on Software Engineering* 48, 9 (2022), 3280–3296.
- [4] Checkmarx. 2023. <https://checkmarx.com/>.
- [5] CodeQL. 2023. <https://codeql.github.com/>.
- [6] CWE. 2023. Common Weakness Enumeration. <https://cwe.mitre.org>.
- [7] Xu Duan, Jingzheng Wu, Shouling Ji, Zhiqing Rui, Tianyue Luo, Mutian Yang, and Yanjun Wu. 2019. VulSniper: Focus Your Attention to Shoot Fine-Grained Vulnerabilities. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence (IJCAI)*, Macao, China. 4665–4671.
- [8] Jiahao Fan, Yi Li, Shaohua Wang, and Tien N. Nguyen. 2020. A C/C++ Code Vulnerability Dataset with Code Changes and CVE Summaries. In *Proceedings of the 17th International Conference on Mining Software Repositories (MSR)*, Seoul, Republic of Korea. ACM, 508–512.
- [9] Flawfinder. 2023. <https://dwheeler.com/flawfinder/>.
- [10] Fortify. 2023. <https://www.microfocus.com/en-us/cyberres/application-security>.
- [11] Michael Fu and Chakkrit Tantithamthavorn. 2022. LineVul: A Transformer-based Line-Level Vulnerability Prediction. In *Proceedings of the 19th IEEE/ACM International Conference on Mining Software Repositories (MSR)*, Pittsburgh, PA, USA. ACM, 608–620.
- [12] Shuitao Gan, Chao Zhang, Xiaojun Qin, Xuwen Tu, Kang Li, Zhongyu Pei, and Zuoning Chen. 2018. CollAFL: Path Sensitive Fuzzing. In *Proceedings of 2018 IEEE Symposium on Security and Privacy (S&P)*, San Francisco, California, USA. IEEE Computer Society, 679–696.
- [13] Wenbo Guo, Dongliang Mu, Jun Xu, Purui Su, Gang Wang, and Xinyu Xing. 2018. LEMNA: Explaining Deep Learning Based Security Applications. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security (CCS)*, Toronto, ON, Canada. ACM, 364–379.
- [14] Hazim Hanif and Sergio Maffei. 2022. VulBERTa: Simplified Source Code Pre-Training for Vulnerability Detection. In *Proceedings of International Joint Conference on Neural Networks (IJCNN)*, Padua, Italy. IEEE, 1–8.
- [15] Ahmad Hazimeh, Adrian Herrera, and Mathias Payer. 2020. Magma: A Ground-Truth Fuzzing Benchmark. *Proceedings of the ACM on Measurement and Analysis of Computing Systems* 4, 3 (2020), 49:1–49:29.
- [16] Infer. 2023. <https://fbinfer.com/>.
- [17] Md Rakibul Islam, Minhaz F. Zibran, and Aayush Nagpal. 2017. Security Vulnerabilities in Categories of Clones and Non-Cloned Code: An Empirical Study. In *Proceedings of 2017 ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM)*, Toronto, ON, Canada. IEEE Computer Society, 20–29.
- [18] Jiyong Jang, Abeer Agrawal, and David Brumley. 2012. ReDeBug: Finding Unpatched Code Clones in Entire OS Distributions. In *Proceedings of 2012 IEEE Symposium on Security and Privacy (S&P)*, San Francisco, California, USA. IEEE Computer Society, 48–62.
- [19] Seulbae Kim and Heejo Lee. 2018. Software Systems at Risk: An Empirical Study of Cloned Vulnerabilities in Practice. *Computers & Security* 77 (2018), 720–736.

- [20] Frank Li and Vern Paxson. 2017. A Large-Scale Empirical Study of Security Patches. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security (CCS)*, Dallas, TX, USA. ACM, 2201–2215.
- [21] Qiang Li, Jinke Song, Dawei Tan, Haining Wang, and Jiqiang Liu. 2021. PDGraph: A Large-Scale Empirical Study on Project Dependency of Security Vulnerabilities. In *Proceedings of the 51st Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*, Taipei, Taiwan. IEEE, 161–173.
- [22] 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)*, Athens, Greece. ACM, 292–303.
- [23] Zhen Li, Deqing Zou, Shouhuai Xu, Zhaoxuan Chen, Yawei Zhu, and Hai Jin. 2022. VulDeeLocator: A Deep Learning-Based Fine-Grained Vulnerability Detector. *IEEE Transactions on Dependable and Secure Computing* 19, 4 (2022), 2821–2837.
- [24] Zhen Li, Deqing Zou, Shouhuai Xu, Hai Jin, Hanchao Qi, and Jie Hu. 2016. VulPecker: An Automated Vulnerability Detection System Based on Code Similarity Analysis. In *Proceedings of the 32nd Annual Conference on Computer Security Applications (ACSAC)*, Los Angeles, CA, USA. ACM, 201–213.
- [25] Zhen Li, Deqing Zou, Shouhuai Xu, Hai Jin, Yawei Zhu, and Zhaoxuan Chen. 2022. SysVR: A Framework for Using Deep Learning to Detect Software Vulnerabilities. *IEEE Transactions on Dependable and Secure Computing* 19, 4 (2022), 2244–2258.
- [26] 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. In *Proceedings of the 25th Annual Network and Distributed System Security Symposium (NDSS)*, San Diego, California, USA. The Internet Society.
- [27] Jingjing Liang, Yaozong Hou, Shurui Zhou, Junjie Chen, Yingfei Xiong, and Gang Huang. 2019. How to Explain a Patch: An Empirical Study of Patch Explanations in Open Source Projects. In *Proceedings of the 30th IEEE International Symposium on Software Reliability Engineering (ISSRE)*, Berlin, Germany. IEEE, 58–69.
- [28] Stephan Lipp, Sebastian Banescu, and Alexander Pretschner. 2022. An Empirical Study on the Effectiveness of Static C Code Analyzers for Vulnerability Detection. In *Proceedings of the 31st ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA)*, Virtual Event, South Korea. ACM, 544–555.
- [29] Bingchang Liu, Guozhu Meng, Wei Zou, Qi Gong, Feng Li, Min Lin, Dandan Sun, Wei Huo, and Chao Zhang. 2020. A Large-Scale Empirical Study on Vulnerability Distribution within Projects and the Lessons Learned. In *Proceedings of the 42nd International Conference on Software Engineering (ICSE)*, Seoul, South Korea. ACM, 1547–1559.
- [30] Changhua Luo, Penghui Li, and Wei Meng. 2022. TChecker: Precise Static Inter-Procedural Analysis for Detecting Taint-Style Vulnerabilities in PHP Applications. In *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security (CCS)*, Los Angeles, CA, USA. ACM, 2175–2188.
- [31] Valentin J. M. Manès, HyungSeok Han, Choongwoo Han, Sang Kil Cha, Manuel Egele, Edward J. Schwartz, and Maverick Woo. 2021. The Art, Science, and Engineering of Fuzzing: A Survey. *IEEE Transactions on Software Engineering* 47, 11 (2021), 2312–2331.
- [32] Ruijie Meng, Zhen Dong, Jialin Li, Ivan Beschastnikh, and Abhik Roychoudhury. 2022. Linear-time Temporal Logic guided Greybox Fuzzing. In *Proceedings of the 44th IEEE/ACM International Conference on Software Engineering (ICSE)*, Pittsburgh, PA, USA. ACM, 1343–1355.
- [33] Yisroel Mirsky, George Macon, Michael Brown, Carter Yagemann, Matthew Pruett, Evan Downing, Sukarno Mertoguno, and Wenke Lee. 2023. VulChecker: Graph-based Vulnerability Localization in Source Code. In *Proceedings of the 32nd USENIX Security Symposium*, Anaheim, CA, USA. USENIX Association, 6557–6574.
- [34] Georgios Nikitopoulos, Konstantina Dritsa, Panos Louridas, and Dimitris Mitropoulos. 2021. CrossVul: A Cross-Language Vulnerability Dataset with Commit Data. In *Proceedings of the 29th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE)*, Athens, Greece. ACM, 1565–1569.
- [35] NVD. 2023. <https://nvd.nist.gov/>.
- [36] Marcus Pendleton, Richard Garcia-Lebron, Jin-Hee Cho, and Shouhuai Xu. 2017. A Survey on Systems Security Metrics. *ACM Computing Surveys* 49, 4 (2017), 62:1–62:35.
- [37] Shiqi Shen, Aashish Kolluri, Zhen Dong, Prateek Saxena, and Abhik Roychoudhury. 2021. Localizing Vulnerabilities Statistically from One Exploit. In *Proceedings of 2021 ACM Asia Conference on Computer and Communications Security (ASIACCS)*, Virtual Event, Hong Kong. ACM, 537–549.
- [38] Morteza Verdi, Ashkan Sami, Jafar Akhondali, Foutse Khomh, Gias Uddin, and Alireza Karami Motlagh. 2022. An Empirical Study of C++ Vulnerabilities in Crowd-Sourced Code Examples. *IEEE Transactions on Software Engineering* 48, 5 (2022), 1497–1514.
- [39] Huanting Wang, Guixian Ye, Zhanyong Tang, Shin Hwei Tan, Songfang Huang, Dingyi Fang, Yansong Feng, Lizhong Bian, and Zheng Wang. 2021. Combining Graph-Based Learning With Automated Data Collection for Code Vulnerability Detection. *IEEE Transactions on Information Forensics and Security* 16 (2021), 1943–1958.
- [40] Xinda Wang, Kun Sun, Archer L. Batcheller, and Sushil Jajodia. 2019. Detecting ‘0-Day’ Vulnerability: An Empirical Study of Secret Security Patch in OSS. In *Proceedings of the 49th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*, Portland, OR, USA. IEEE, 485–492.
- [41] Xinda Wang, Shu Wang, Pengbin Feng, Kun Sun, and Sushil Jajodia. 2021. PatchDB: A Large-Scale Security Patch Dataset. In *Proceedings of the 51st Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*, Taipei, Taiwan. IEEE, 149–160.
- [42] Xinda Wang, Shu Wang, Kun Sun, Archer L. Batcheller, and Sushil Jajodia. 2020. A Machine Learning Approach to Classify Security Patches into Vulnerability Types. In *Proceedings of the 8th IEEE Conference on Communications and Network Security (CNS)*, Avignon, France. IEEE, 1–9.
- [43] Qiushi Wu, Yang He, Stephen McCamant, and Kangjie Lu. 2020. Precisely Characterizing Security Impact in a Flood of Patches via Symbolic Rule Comparison. In *Proceedings of the 27th Annual Network and Distributed System Security Symposium (NDSS)*, San Diego, California, USA. The Internet Society.
- [44] Yueming Wu, Deqing Zou, Shihan Dou, Wei Yang, Duo Xu, and Hai Jin. 2022. VulCNN: An Image-Inspired Scalable Vulnerability Detection System. In *Proceedings of the 44th IEEE/ACM International Conference on Software Engineering (ICSE)*, Pittsburgh, PA, USA. ACM, 2365–2376.
- [45] Fabian Yamaguchi, Nico Golde, Daniel Arp, and Konrad Rieck. 2014. Modeling and Discovering Vulnerabilities with Code Property Graphs. In *Proceedings of 2014 IEEE Symposium on Security and Privacy (S&P)*, Berkeley, CA, USA. IEEE Computer Society, 590–604.
- [46] Fabian Yamaguchi, Felix “FX” Lindner, and Konrad Rieck. 2011. Vulnerability Extrapolation: Assisted Discovery of Vulnerabilities Using Machine Learning. In *Proceedings of the 5th USENIX Workshop on Offensive Technologies (WOOT)*, San Francisco, CA, USA. USENIX Association, 118–127.
- [47] 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 (ACSAC)*, Orlando, FL, USA. ACM, 359–368.
- [48] Mengyuan Zhang, Xavier de Carné de Carnavalet, Lingyu Wang, and Ahmed Ragab. 2019. Large-Scale Empirical Study of Important Features Indicative of Discovered Vulnerabilities to Assess Application Security. *IEEE Transactions on Information Forensics and Security* 14, 9 (2019), 2315–2330.
- [49] Yaqin Zhou, Shangqing Liu, Jing Kai Siow, Xiaoning Du, and Yang Liu. 2019. Devign: Effective Vulnerability Identification by Learning Comprehensive Program Semantics via Graph Neural Networks. In *Proceedings of 2019 Annual Conference on Neural Information Processing Systems (NeurIPS)*, Vancouver, BC, Canada. 10197–10207.
- [50] Deqing Zou, Yawei Zhu, Shouhuai Xu, Zhen Li, Hai Jin, and Hengkai Ye. 2021. Interpreting Deep Learning-based Vulnerability Detector Predictions Based on Heuristic Searching. *ACM Transactions on Software Engineering and Methodology* 30, 2 (2021), 23:1–23:31.