



Software Engineering Group
Department of Computer Science
Nanjing University
<http://seg.nju.edu.cn>

Technical Report No. NJU-SEG-2020-IJ-008

2020-IJ-008

Automatic Buffer Overflow Warning Validation

FengJuan Gao, Yu Wang, Linzhang Wang, Ziji Yang and Xuan-Dong Li

Technical Report 2020-IJ-008

Most of the papers available from this document appear in print, and the corresponding copyright is held by the publisher. While the papers can be used for personal use, redistribution or reprinting for commercial purposes is prohibited.

Automatic Buffer Overflow Warning Validation

Feng-Juan Gao^{1,2}, Yu Wang^{1,2}, Lin-Zhang Wang^{1,2,*}, *Distinguished Member, CCF*
Zijiang Yang³, *Senior Member, IEEE*, and Xuan-Dong Li^{1,2}, *Fellow, CCF*

¹State Key Laboratory for Novel Software Technology, Nanjing University, Nanjing 210023, China

²Department of Computer Science and Technology, Nanjing University, Nanjing 210023, China

³Department of Computer Science, Western Michigan University, Kalamazoo 49008-5466, U.S.A.

E-mail: {fjgao,yuwang_cs}@smail.nju.edu.cn; lzwang@nju.edu.cn; zijiang.yang@wmich.edu; lxd@nju.edu.cn

Received April 11, 2020; revised October 22, 2020.

Abstract Static buffer overflow detection techniques tend to report too many false positives fundamentally due to the lack of software execution information. It is very time consuming to manually inspect all the static warnings. In this paper, we propose BovInspector, a framework for automatically validating static buffer overflow warnings and providing suggestions for automatic repair of true buffer overflow warnings for C programs. Given the program source code and the static buffer overflow warnings, BovInspector first performs warning reachability analysis. Then, BovInspector executes the source code symbolically under the guidance of reachable warnings. Each reachable warning is validated and classified by checking whether all the path conditions and the buffer overflow constraints can be satisfied simultaneously. For each validated true warning, BovInspector provides suggestions to automatically repair it with 11 repair strategies. BovInspector is complementary to prior static buffer overflow discovery schemes. Experimental results on real open source programs show that BovInspector can automatically validate on average 60% of total warnings reported by static tools.

Keywords buffer overflow, static analysis warning, symbolic execution, automatic repair

1 Introduction

Buffer overflow occurs when more data is written into a buffer than the buffer capacity, causing extra data being written into memory adjacent to the buffer. If the adjacent memory before being overwritten has stored information (such as the pointer to the previous frame and return address) that is critical for the OS to correctly execute programs, buffer overflow may cause unpredictable behaviors. In a buffer overflow attack, the attacker carefully crafts his/her input data to vulnerable software so that the unpredictable behavior is that the OS executes his/her malicious code embedded in the overflow data with the privilege of the vulnerable software.

Although more than 40 years have passed since the

buffer overflow technique was first documented by Anderson in 1972^[1] and almost 30 years have passed since the buffer overflow technique was first exploited by the infamous Morris worm in 1988, buffer overflow remains the most common type of software vulnerabilities, as shown in the recent studies of software vulnerability databases^[2], and it is likely to remain so for many years to come. Most existing software has buffer overflow vulnerabilities, which are unknown to their vendors and users, but could be exploited by attackers. Most future software will still be written by programmers who are not well trained in software security. The inherently unsafe languages C and C++ will remain popular for performance and backward compatibility. Although we have known how to avoid buffer overflow problems in writing programs for many years, having such know-

Regular Paper

Special Section on Software Systems 2020

A preliminary version of the paper was published in the Proceedings of ASE 2016.

This work was supported by the National Natural Science Foundation of China under Grant No. 62032010, and partially by the Postgraduate Research and Practice Innovation Program of Jiangsu Province of China.

*Corresponding Author

©Institute of Computing Technology, Chinese Academy of Sciences 2020

ledge is far from enough to thwart the rampant buffer overflow issue.

There are two general approaches to identifying buffer overflow vulnerabilities: static program analysis^[3-9] and dynamic execution analysis^[10-14]. The dynamic execution analysis approach needs a specific test case to trigger a buffer overflow vulnerability, which may not be easy to find. The static program analysis approach scans software source code to discover the code segments that are possibly vulnerable to buffer overflow attacks. The key advantage of such schemes is that buffer overflow vulnerabilities can be discovered and fixed before software deployment. However, due to the high false positive rate in static analysis, validating all the reported warnings takes huge manual effort. Therefore, people need methods to automatically validate the static warnings and fix the true warnings.

In this paper, we propose BovInspector, a framework for automatically validating static buffer overflow warnings and providing suggestions for automatic repair of true buffer overflow warnings for C programs. BovInspector is complementary to prior static buffer overflow discovery schemes. The key contribution of BovInspector is on eliminating the need to manually inspect the warnings that are actually false warnings, and providing automatic repair suggestions for validated true warnings, which will save a tremendous amount of manual efforts. The key idea of BovInspector is to use symbolic execution to automatically identify those buffer overflow warnings that are true warnings or false warnings. The advantage of symbolic execution is its capability to explore program execution states that are unavailable to static program analysis. The disadvantage of symbolic execution is the infamous path explosion issue, i.e., the number of execution paths grows exponentially with the number of branching points. In BovInspector, to avoid path explosion, we use the warning paths to guide the symbolic execution so that we only focus on these warning paths. Our experimental results on real open source programs show that BovInspector can automatically validate on average 60% of total warnings reported by the static tool Fortify^①.

In summary, this paper contributes the followings:

- an automated framework for validating static buffer overflow warnings and providing repair suggestions for the true warnings;
- an open source tool, BovInspector^②;

- an evaluation to show the effectiveness and efficiency of our tool.

Our previous demonstration paper^[15] only introduced an initial framework and preliminary experiment results. Compared with [15], this paper 1) provides theoretical foundations of buffer overflow validation by defining buffer overflow models and buffer overflow warnings, 2) extends the buffer overflow repair method by supporting eight more repair strategies, and 3) conducts more comprehensive experiments on both synthetic and real-world programs to further demonstrate the effectiveness and efficiency of our method.

The rest of the paper proceeds as follows. We first introduce the background knowledge of symbolic execution and buffer overflow in Section 2. Then we present the formal description of static buffer overflow warning validation in Section 3. In Section 4, we present the framework of BovInspector and the details of the method. In Section 5, we show implementation details and experimental results of BovInspector. Then, we review related work in Section 6. Finally, we give conclusions and future work in Section 7.

2 Background

2.1 Symbolic Execution

Symbolic execution is a classical technique for software testing and analysis^[16]. It is used to systematically test a program and generate test input with high coverage. Symbolic execution uses symbolic values as the input, instead of concrete input, to explore the execution space of a program. When symbolic execution encounters a branching condition, it forks the execution state, following both branch directions and updating the corresponding path constraints on the symbolic input. When it reaches a program exit or hits an error, the current path constraint will be solved to find a concrete test case that drives program execution to this program location. KLEE^[17] is a state-of-the-art open source symbolic execution engine and is widely used in researches.

2.2 Buffer Overflows

Buffer overflows, both on the stack and on the heap, are a major source of security vulnerabilities in C and

^①<https://www.microfocus.com/en-us/products/static-code-analysis-sast/overview>, Oct. 2020.

^②BovInspector and a technical paper with illustrative examples and more experimental details are available at <http://bovinspectortool.github.io/project/>, Sept. 2020.

C++ code. A buffer is a region of a physical memory storage used to temporarily store the data. Using a buffer while programming in C and C++ makes a lot of sense and generally speeds up the calculation process, which may also lead to unexpected buffer overflow. A buffer overflow, or buffer overrun, occurs when more data is put into a fixed-length buffer than what the buffer can handle.

There are two kinds of buffers we deal with, buffers of characters and buffers of the other values. We propose the buffer model shown in Fig.1. The symbols in the figure are defined as follows.

- S_{Buf} : the start address of buffer Buf .
- $Size(Buf)$: the available buffer space, in bytes, of buffer Buf starting from S_{Buf} .
- S_{buf} : the start address of buffer buf .
- Off_{buf} : the offset of bytes from the base address of the whole buffer of buffer buf , $Off_{buf} = S_{buf} - S_{Buf}$.
- $Size(buf)$: the available buffer space, in bytes, of buffer buf starting from S_{buf} .
- $Len(buf)$: the length of the content stored in buffer buf , which starts from S_{buf} and ends with terminator $'\0'$ ($'\0'$ is excluded). To be noticed, if there is no terminator in buffer buf , we consider its length to be infinite in our model. $Len(buf)$ is only available for buffers of characters.
- i : i is the index of buffer buf when accessing $S_{buf}[i]$, which points to the i -th element of buf . Each element in buf takes up a fixed number of bytes, marked as $typesize$.

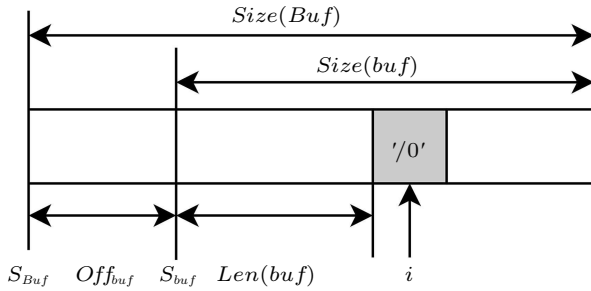


Fig.1. Buffer model.

For a buffer Buf with size $Size(Buf)$ and start address S_{Buf} , we assume a buffer operation accesses the data in the buffer buf and we know $buf = Buf + Off_{buf}$. Therefore, we can learn that the size of buf will be $Size(buf) = Size(Buf) - Off_{buf}$. The start address of buf will be $S_{buf} = S_{Buf} + Off_{buf}$. If the buffer operation is a direct buffer access, i.e., $S_{buf}[i]$, it will access the i -th element of buffer buf , namely the $(i + Off_{buf}/typesize)$ -th element of buffer

Buf at address $S_{Buf} + Off_{buf} + i \times typesize$. For a buffer of characters, a buffer operation may read the characters from it. The length of characters that will be accessed is $Len(buf)$.

3 Definition of Buffer Overflow Models and Warning Classifications

This section presents a lightweight formalism to define our buffer overflow model and the three categories of static buffer overflow warnings validated by our method.

We first define the simple imperative language (a core subset of C) in Fig.2 to represent all important features of C that are necessary for validating a buffer overflow warning. We use V , L , I , Φ and E to represent the sets of variables, labels, integers, predicates and statements, respectively. Each statement e has a unique label $l \in L$ that is used to identify e . In the following context, the statement label l for e means the line number of e . We use $buf = alloc(m \times typesize)$ to represent an allocation of a new buffer with size $m \times typesize$ in the stack or on the heap. For simplicity, we use $bufAPI(dest, src, n)$ to represent all kinds of buffer APIs, where $dest$ is a variable pointing to the destination buffer, src is a variable pointing to the source buffer and n , if provided by the API, indicates the number of bytes to be operated. If n is not supported by an API, then it means that n is unlimited. $bufAPI$ will copy the data with length n from src to $dest$. Whether these buffer operations will cause a buffer overflow depends on the remaining size of $dest$ buffer and the length of src content.

Variables	$dest, src, buf \in V$
Labels	$l \in L$
Integers	$i, m, n, typesize \in I$
Predicates	$\phi \in \Phi$
Statements	$e \in E$
$e ::= {}^l buf = alloc(m \times typesize) \mid {}^l buf[i] = x \mid {}^l e; {}^l e \mid$ ${}^l bufAPI(dest, src, n, \dots) \mid {}^l if(\phi) \text{ then } {}^l e \text{ else } {}^l e \mid$ ${}^l buf = realloc(buf, m \times typesize) \mid$ ${}^l *(buf + i) = x$	

Fig.2. Core language for validating buffer overflow warnings.

In a C/C++ program, a buffer is commonly used. There are also many buffer operations writing data into buffers in a program. As shown in the language, there are two categories of buffer operations: API call and

direct buffer access. To characterize under what condition a buffer operation will lead to a buffer overflow, we propose our Buffer Overflow Models as follows.

3.1 Buffer Overflow Models

According to the types of buffer operations, we divide buffer overflow models into two categories.

Buffer Overflow Models of API Call. For API calls, such as *strcpy*, the buffer overflow can be detected by analyzing the parameters. We propose some constraint models for the APIs operating buffers in C99^③ and Linux system call interfaces^④. To check for buffer overflow, we examine whether the data written to a buffer exceeds the buffer size. We list the constraints for different types of APIs shown in Table 1, in which $Len(src)$ and $Size(dest)$ are defined in Fig.1. Note that we select the parameter format of the first API in the group as a representative. Other APIs have a similar parameter format. We classify the APIs into four types, namely unbounded content sensitive buffer operations, bounded content sensitive buffer operations, bounded content insensitive buffer operations and direct buffer accesses.

Buffer Overflow Models of Direct Buffer Access. For direct buffer accesses in Table 1, we propose the constraint model for array and pointer accesses. Access-

ing a buffer by array access can be represented as $buf[i] = x$, while buf represents a buffer, and x is the value to be assigned to the corresponding address. Accessing a buffer by pointer access can be represented as $*(buf + i) = x$. Without loss of generality, we also take type casting into consideration, e.g., when a pointer p points to a character buffer $buf[5]$, the buffer access $*((int*)p + 1) = x$ will trigger a buffer overflow. Therefore, the buffer overflow condition is $(i + 1) \times typesize > Size(buf)$.

Based on the above buffer model and buffer overflow models, we introduce the definition of static buffer overflow warning and the principles to classify static buffer overflow warnings, which will be used in our buffer overflow validation module.

3.2 Buffer Overflow Warnings

Given the source code of a program, static analysis tools will locate all the statements that declare a buffer, track all the statements that perform operations on the buffers, and report a buffer overflow warning if the operation may violate the predefined secure coding rules.

Before validating static buffer overflow warnings, we first give the definition of static buffer overflow warning as follows.

Table 1. Buffer Overflow Models

Type	API	Parameter Format	Overflow Constraint
Unbounded content sensitive buffer operations	<i>strcpy</i>	(char* <i>dest</i> , const char* <i>src</i>)	$Len(src) \geq Size(dest)$
	<i>strcat</i>	(char* <i>dest</i> , const char* <i>src</i>)	$Len(src) + Len(dest) \geq Size(dest)$
	<i>sprintf</i> , <i>vsprintf</i>	(char* <i>str</i> , const char* <i>format</i> , ...)	$format_string_length \geq Size(str)$
	<i>scanf</i> , <i>vscanf</i>	(char* <i>format</i> , ...)	$format_string_length \geq Size(dest_i)$
	<i>sscanf</i> , <i>vsscanf</i>	(const char* <i>s</i> , char* <i>format</i> , ...)	$format_string_length \geq Size(dest_i)$
	<i>fscanf</i> , <i>vfscanf</i>	(FILE* <i>stream</i> , const char* <i>format</i> , ...)	$format_string_length \geq Size(dest_i)$
Bounded content sensitive buffer operations	<i>strncpy</i> , <i>snprintf</i> , <i>vsnprintf</i>	(char* <i>dest</i> , const char* <i>src</i> , size_t <i>n</i>)	$n > Size(dest)$
	<i>strncat</i>	(char* <i>dest</i> , const char* <i>src</i> , size_t <i>n</i>)	$\min\{Len(src), n\} + Len(dest) \geq Size(dest)$
	<i>fgets</i>	(char* <i>str</i> , int <i>num</i> , FILE* <i>stream</i>)	$num > Size(str)$
	<i>fread</i>	(void* <i>ptr</i> , size_t <i>size</i> , size_t <i>count</i> , FILE* <i>stream</i>)	$size \times count > Size(ptr)$
Bounded content insensitive buffer operations	<i>memcpy</i> ,	(int <i>fd</i> , void* <i>buf</i> , size_t <i>count</i>)	$count > Size(buf)$
	<i>memmove</i> ,	(char* <i>dest</i> , const char* <i>src</i> , size_t <i>n</i>)	$n > Size(dest)$
	<i>memset</i>		
Direct buffer accesses	$buf[i]$, $*(buf + i)$	N/A	$(i + 1) \times typesize > Size(buf)$

Note: *format_string_length* is calculated by considering the impact of formatting symbols.

^③ <https://www.iso.org/standard/29237.html>, Oct. 2020.

^④ <https://man7.org/linux/man-pages/man2/syscall.2.html>, Oct. 2020.

Definition 1 (Static Buffer Overflow Warning). A static buffer overflow warning ω is represented as a tuple: $(\langle l_1, l_2, \dots, l_n \rangle, b)$, where l_1 is the label of the statement where a buffer is first declared or the entry of the main function, each l_i ($2 \leq i \leq n$) is the label of the statement where an operation is performed on the buffer, and b is the label of the statement where an overflow may occur on the buffer, which is regarded as a buffer warning point.

Considering the source code of the example program in Fig.3, the static analysis tool reports four possible buffer overflow warnings that can be represented as $\omega_0 = (\langle 4 \rangle, 6)$, $\omega_1 = (\langle 19, 23, 9 \rangle, 12)$, $\omega_2 = (\langle 19, 23, 9, 12 \rangle, 14)$, $\omega_3 = (\langle 25 \rangle, 30)$.

```

1  #define MAX_LEN  24
2  #define MIN_LEN  4
3  void usage() {
4      char des_buff[MIN_LEN];
5      char* src_buff="this is an example";
6      strcpy(des_buff, src_buff);
7  }
8  int initialize(char* argv_string) {
9      char mapped_argv[MIN_LEN];
10     if (strlen(argv_string) >= MAX_LEN)
11         return 0;
12     strcpy(mapped_argv, argv_string);
13     if (argv_string[strlen(argv_string)
14         -1] != '-') {
15         strcat(mapped_argv, "-");
16     }
17     return MAX_LEN;
18 }
19 int main(int argc, char** argv) {
20     char* mode = (char*)malloc(MIN_LEN);
21     int len = 0;
22     if (strlen(argv[1]) < MIN_LEN)
23         len = initialize(argv[1]);
24     if (len > 0) {
25         mode = (char*)realloc(mode, len);
26     }
27     else {
28         mode[0] = argv[1][0];
29     }
30     mode[MIN_LEN-1] = '\0';
31     free(mode);
32     return 0;
33 }

```

Fig.3. Example program.

For each buffer warning point b , there may be several paths from the entrance of the program to that point. We regard these paths as a warning path set.

Definition 2 (Warning Path Set). For a static buffer overflow warning ω , a warning path set ps is a sequence (P_1, P_2, \dots, P_m) , where each P_i is a set of path segments ρ , jointly constituting a path from the program entrance to the warning point b , and each path

segment ρ records the statement label of the first statement of each basic block in the path. ps represents all the complete paths from the entrance of the program to a buffer overflow point. The path constraint of P_i is $\bigwedge_{j=1}^{j=n_i} \varphi_j^i$, where n_i is the number of branch statements in path P_i , and φ_j^i is the constraint of the j -th branch in path P_i .

Here we take as an example the warning $\omega_2 = (\langle 19, 23, 9, 12 \rangle, 14)$ for the program in Fig.3. Its warning path set is $(\langle 19, 23, 8, 12, 14 \rangle)$.

The key idea to validate a static buffer overflow warning ω is to find the evidence that there exists a test case that follows one of the paths P_i in the warning path set ps , reaches the buffer warning point b , and finally triggers the buffer overflow. Therefore, at a buffer warning point b , BovInspector solves the conjunction of path constraints and overflow constraints (denoted as function $OC(b)$), in order to validate the corresponding buffer overflow warning. $OC(b)$ is constructed and solved according to the overflow constraint in Table 1 regarding the API at the buffer warning point b . Based on the solving result, a buffer overflow warning will be validated as a true buffer overflow warning, a false buffer flow warning, or an undecided buffer overflow warning.

For a static buffer overflow warning ω , a path P_i in the warning path set ps is an overflowable path if the conjunction of its path constraints and overflow constraints is satisfiable. If it is unsatisfiable, the path will be regarded as a safe path. If the solver cannot provide the result of the constraints within the given time limit, the path will be regarded as an undecided path. For this case, we use *TIMEOUT* to represent the situation that the constraint solver cannot provide the result of the constraints within the given time limit for path P_i at warning point b .

Definition 3 (Overflowable Path). Given a warning $\omega = (\langle l_1, l_2, \dots, l_n \rangle, b)$, with the warning path set $ps = (P_1, P_2, \dots, P_m)$, $\forall P_i \in ps$, P_i is an overflowable path if (1) holds.

$$\bigwedge_{j=1}^{j=n_i} \varphi_j^i \wedge OC(b) \equiv SAT. \quad (1)$$

Definition 4 (Safe Path). Given a warning $\omega = (\langle l_1, l_2, \dots, l_n \rangle, b)$, with the warning path set $ps = (P_1, P_2, \dots, P_m)$, $\forall P_i \in ps$, P_i is a safe path if (2) holds.

$$\bigwedge_{j=1}^{j=n_i} \varphi_j^i \wedge OC(b) \equiv UNSAT. \quad (2)$$

Definition 5 (Undecided Path). Given a warning $\omega = (\langle l_1, l_2, \dots, l_n \rangle, b)$, with the warning path set $ps = (P_1, P_2, \dots, P_m)$, $\forall P_i \in ps$, P_i is an undecided path if (3) holds.

$$\bigwedge_{j=1}^{j=n_i} \varphi_j^i \wedge OC(b) \equiv TIMEOUT. \quad (3)$$

If we have found an overflowable path P_i in ps for warning ω , the warning ω will be regarded as a true warning. If we find a path that can validate the warning as a true warning, it is unnecessary to validate other paths in the warning path set ps . Otherwise, we need to traverse all paths in ps to validate a warning. If all paths in ps are safe paths, the warning ω will be regarded as a false warning. If some paths in ps are undecided paths, when there exists at least one overflowable path validating the warning as a true warning, the warning will be validated as a true warning regardless of the undecided paths. Otherwise, we cannot tell whether the warning is a false warning or not because there may be some overflowable paths hidden in the time-out paths. Therefore, we classify this kind of warnings as undecided warnings.

Definition 6 (True Buffer Overflow Warning). Given a warning $\omega = (\langle l_1, l_2, \dots, l_n \rangle, b)$, with the warning path set $ps = (P_1, P_2, \dots, P_m)$, ω is a true buffer overflow warning if (4) holds.

$$\exists P_i \in ps : \bigwedge_{j=1}^{j=n_i} \varphi_j^i \wedge OC(b) \equiv SAT. \quad (4)$$

Definition 7 (False Buffer Overflow Warning). Given a warning $\omega = (\langle l_1, l_2, \dots, l_n \rangle, b)$, with the warning path set $ps = (P_1, P_2, \dots, P_m)$, ω is a false buffer overflow warning if (5) holds.

$$\forall P_i \in ps : \bigwedge_{j=1}^{j=n_i} \varphi_j^i \wedge OC(b) \equiv UNSAT. \quad (5)$$

Definition 8 (Undecided Buffer Overflow Warning). Given a warning $\omega = (\langle l_1, l_2, \dots, l_n \rangle, b)$, with the warning path set $ps = (P_1, P_2, \dots, P_m)$, ω is an undecided buffer overflow warning if it is neither a true buffer overflow warning nor a false buffer overflow warning, or, more formally, if (6) holds.

$$\begin{aligned} & (\forall P_i \in ps : (\bigwedge_{j=1}^{j=n_i} \varphi_j^i \wedge OC(b) \equiv UNSAT) \\ & \vee (\bigwedge_{j=1}^{j=n_i} \varphi_j^i \wedge OC(b) \equiv TIMEOUT)) \\ & \wedge (\exists P_i \in ps : (\bigwedge_{j=1}^{j=n_i} \varphi_j^i \wedge OC(b) \equiv TIMEOUT)). \quad (6) \end{aligned}$$

In some cases, all the paths in the warning path set ps are infeasible, namely $\bigwedge_{j=1}^{j=n_i} \varphi_j^i \equiv UNSAT$ for all paths, which will make (5) hold. Some of these false warnings can be found before symbolic execution, which will be discussed in Subsection 4.2. Other false warnings will be validated during symbolic execution, which will be discussed in Subsection 4.4.

4 Automatic Buffer Overflow Validation and Repair

According to Fig.4, BovInspector consists of four modules: warning reachability analysis, guided symbolic execution, buffer overflow validation, and targeted automatic repair suggestions.

4.1 Approach Overview

Warning Reachability Analysis. The input to this module is the set of buffer overflow warnings reported by static analysis tools together with the source code. For each buffer overflow warning w , this module computes all the complete paths that start from the program entrance, go through the statements in w , and end at the warning point b . For the warnings that are

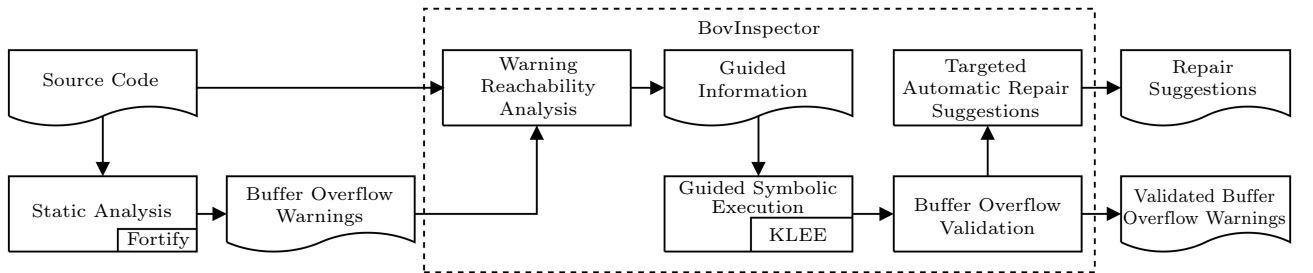


Fig.4. Overview of our automatic buffer overflow warning validation and bug repair approach.

not reachable, this module can identify and prune them; the corresponding warnings are then classified as false warnings. This module will output the warning path set ps for each static buffer overflow warning w as guided information.

Guided Symbolic Execution. This module takes as input the guided information generated by the above module. It is an extension of the traditional symbolic execution engine. We perform symbolic execution on the source code starting from the program entrance. At each branching point with k possible branches, we first make k replicates of the execution state; each replicate of the execution state represents a complete path from the program entrance to the branching point. At each branching point, this module tries to prune the state that cannot lead to the warning points by querying the guided information. When the execution reaches a warning point, we call the buffer overflow validation module to check whether buffer overflow can indeed happen.

Buffer Overflow Validation. This module takes as input the static buffer overflow warnings and outputs the validated buffer overflow warnings. During symbolic execution, each execution state maintains the necessary information, in particular path constraints. When the execution encounters a warning point, we will construct the buffer overflow constraints according to Table 1. Then the conjunction of the path constraints and the buffer overflow constraints will be fed to a constraint solver to examine whether all these conditions can be simultaneously satisfied. Based on the solution of constraint solving, we validate the path as an overflowable path, a safe path or an undecided path. If it is an overflowable path, then the warning will be validated as a true warning. A test case will be generated for it. Then the warning point will be removed from the checking list. Otherwise, the execution will continue until all the paths for the warning are executed. If all paths for a warning are safe paths, then the warning is a false warning. If there exist some undecided paths besides the false paths, the warning will be regarded as an undecided warning.

Targeted Automatic Repair Suggestions. By investigating 100 highly ranked buffer overflow vulnerabilities from 2009 to 2014 in the Common Vulnerabilities and Exposures (CVE)^⑤ and the benchmarks from prior buffer overflow detection work, we discover a total of 11 common repair strategies^[18]. For each validated

true buffer overflow warning, we automatically generate repair code as suggestions for all the usable repair strategies. Moreover, programmers can also manually configure a preferred repair strategy.

4.2 Warning Reachability Analysis

For warning reachability analysis, we first generate the Inter-procedural Control Flow Graph (ICFG)^[19]. Then, for each static buffer overflow warning $\omega = ((l_1, l_2, \dots, l_n), b)$, we calculate all the complete execution paths that start from the program entrance, go through the statements contained in the warning (l_1, l_2, \dots, l_n) , and end at the warning point b . To do that, we first map each statement label in ω to its corresponding basic block, i.e., l_1, l_2, \dots, l_n, b corresponds to $bb_{l_1}, bb_{l_2}, \dots, bb_{l_n}, bb_{l_{n+1}}$. Then, we calculate the complete paths covering these basic blocks. A straightforward solution is to perform a depth-first traversal on ICFG to find such complete paths. However, this will lead to the redundant traversal of many paths that do not contain any path segments of the buffer overflow warnings. In this work, for each warning point, we perform backward tracking on ICFG starting from the warning point. By backward tracking, we can ignore the ICFG nodes that do not lead to the warning point. For each pair of $(bb_{l_{i-1}}, bb_{l_i})$, where $2 \leq i \leq n+1$, we use depth first search to calculate all the path segments between the two basic blocks (denoted as $PS_{i-1,i}$). Similarly, we analyze the partial warning path set $PS_{i-2,i-1}$ for $(bb_{l_{i-2}}, bb_{l_{i-1}})$. By combining and connecting these two partial warning path sets $PS_{i-2,i-1}$ and $PS_{i-1,i}$, we compute the partial warning path set $PS_{i-2,i}$ for $(bb_{l_{i-2}}, bb_{l_i})$. The above steps repeat until we reach the program entrance. Finally, we get the warning path set ps of a static buffer overflow warning ω . The warning path set will be used as the guided information in the guided symbolic execution.

4.3 Guided Symbolic Execution

In this work, the program under test is symbolically executed to identify whether the possible buffer overflows can be triggered.

As shown in **Algorithm 1**, we extend the traditional symbolic execution engine to support guided symbolic execution and buffer overflow validation. The extended parts are written in bold. To be noticed, a symbolic execution engine is quite complex and here we only excerpt

^⑤<https://cve.mitre.org/>, Oct. 2020.

Algorithm 1. Guided Symbolic Execution

Input: *WarningList*, *Pathsets*
Output: *ValidationReport*

```

1  GuidedSymbolicExecution(){
2  ExecutionStatePool =  $\emptyset$ ;
3  AddedStateSet =  $\emptyset$ ;
4  RemoveStateSet =  $\emptyset$ ;
5  ExecutionStatePool.add(initialState);
6  while ExecutionStatePool.size > 0 && ! TIMEOUT do
7  |  $es = selectState(ExecutionStatePool)$ ;
8  |  $ExecuteInstruction(es)$ ;
9  |  $UpdateState(ExecutionStatePool, AddedStateSet,$ 
  |  $RemoveStateSet)$ ;
10 }
11
12  $ExecuteInstruction(ExecutionState es)$ {
13 if  $\exists \omega \in WarningList \ \&\& \ es.pc == \omega.b$  then
14 | Call Validation Module;
15 if  $es.pc.instructionType == FORK$  then
16 |  $es2 = fork(es)$ ;
17 |  $es.pc = trueBranchStmt$ ;
18 |  $es2.pc = falseBranchStmt$ ;
19 |  $AddedStateSet.add(es2)$ ;
20 |  $StatePruning(es, es2)$ ;
21 Other operations in traditional symbolic execution...
22 }
23
24  $StatePruning(ExecutionState es, ExecutionState es2)$ {
25 if  $es \neq NULL \ \&\& \ es2 \neq NULL$  then
26 |  $l1 = es.pc.stmtLabel$ ;
27 |  $l2 = es2.pc.stmtLabel$ ;
28 if  $Pathsets.contains(l1) \ \&\& \ !Pathsets.contains(l2)$  then
29 |  $RemoveStateSet.add(es2)$ ;
30 if  $!Pathsets.contains(l1) \ \&\& \ Pathsets.contains(l2)$  then
31 |  $RemoveStateSet.add(es)$ ;
32 }

```

the actions related to our algorithm. This module takes as input the static buffer overflow warnings reported by some static analysis tool and the warning path sets generated by warning reachability analysis, denoted as *WarningList* and *Pathsets*. *ExecutionStatePool* will be maintained during execution. The overall process of guided symbolic execution is executing a loop until the pool is empty or time-out (lines 6–10). In the loop, the pool will be updated by adding or removing states (line 9). This is the default move of the traditional symbolic execution engine. In each iteration, one of the states will be selected (line 7) to execute (line 8) according to a specific rule (e.g., randomly choose one). *pc* (i.e., the program counter) in the state represents the instruction to be executed. Before executing the instruction, we will first check its type. If the statement label of *pc* reaches a warning point (i.e., $\omega.b$), the buffer overflow validation module will be invoked (lines 13 and 14), which will be discussed in Subsection 4.4. If it is

a branch or loop instruction, we will fork a duplicate of the current state (i.e., *es*) and add it (i.e., *es2*) into *AddedStateSet* (lines 15–19). Then *StatePruning* will be invoked to try to prune states that cannot lead to the warning points.

Here we try to prune states right after the symbolic execution engine finished interpreting a branch instruction and only check the two execution states that correspond to the branch. This strategy has two motivations. First, if the unnecessary state can be pruned at the very beginning when it is just generated, the symbolic execution engine can explore the necessary paths more efficiently. Second, the branch information can help us to identify an execution path better. That is because we use the statement label of the first statement of a basic block to represent the block in our warning path set, referring to Definition 2 in Section 3. After a branch, *pc* of the true or false branch state is exactly the first statement of the corresponding block, i.e., *l1* for *es* and *l2* for *es2*.

In *StatePruning*, to decide whether an execution state can be removed from the execution state pool, we check whether the execution state matches the warning path sets. Namely, we will check the existence of the two statement labels (i.e., *l1* and *l2*) in *Pathsets* (lines 28–31). If only one statement label is contained in a certain warning path set, the execution state corresponding to the other statement label will be removed. If both statement labels are contained in the path set, we do not remove any of the execution states. Note that if both statement labels are not in the path set, we also do not remove any of the execution states. One reason is that the symbolic execution engine may explore the internal part of some library calls that do not appear in any warning path. Another reason is to support validating the buffer overflow warnings whose warning points are in the loop. Each path in a warning path set starts from the program entrance and ends at the warning point. Namely, the branch information after the warning point is not contained in any warning path set. But the buffer overflow may be triggered by iterating the loop. Therefore, to iterate the loop, the two states of the branch after the warning point should not be pruned. In this way, our approach allows an execution to continue after calling the buffer overflow validation module if the current path cannot trigger a buffer overflow but the warning point is in a loop. When an execution reaches the exit of the loop, which is a branch, the execution state jumps back to the entrance of the loop. Since the buffer overflow point

is in the loop, the entrance of the loop is contained in the warning path set. The execution will find a match and enter the loop again. Thus, our method can guide symbolic execution to iterate the loop until the loop ends or a buffer overflow is found within a loop time threshold.

4.4 Buffer Overflow Validation

When the symbolic execution engine encounters an instruction matching a warning point in the warning list, the buffer overflow validation module will be invoked (lines 12 and 13 in Algorithm 1). The instruction encountered at the warning point is a buffer operation, which may be the APIs listed in Table 1, e.g., a call to *strcpy*. To validate static buffer overflow warnings during symbolic execution, we extend the symbolic execution engine to monitor buffer operations, collect path constraints, construct overflow constraints and validate warnings during symbolic execution.

To validate a warning, first we need to fetch the buffer information from the symbolic execution engine. It provides a memory management model to handle the information of all variables, including the addresses, size, contents, etc. By analyzing the memory object, we obtain the information of the buffer to be checked, including the starting address of the buffer, the offset from the starting address, and the size of the buffer. The symbolic execution engine updates the content in the buffer whenever there is an operation on the buffer. Therefore, before executing the instruction at the warning point, it is easy to fetch the latest content of the buffer related to the instruction. However, when the data copied to the buffer is symbolic, the length information is difficult to obtain. To address this issue, we design a model to represent the length of a symbolic string variable. Let *buf* be a symbolic string variable, $Len(buf)$ denotes the length of string *buf* in bytes, and $\backslash 0$ denotes the end of a string. In our model, we apply the following constraints based on the string length to judge buffer overflow: $Len(buf) = i$, where $\forall j \in [0, i), S_{buf}[j] \neq \backslash 0$ and $S_{buf}[i] = \backslash 0$. To unify the validation model, we apply the definition on both symbolic strings and concrete strings. As mentioned before, we use $Size(buf)$ to represent the size of the buffer assigned to the parameter *buf*, which can be obtained from the memory management model in the symbolic execution engine. When calculating the length of buffer *buf*, if the first $Size(buf)$ bytes have no $\backslash 0$, it means that the string in the buffer has no

terminator. In this case, we consider $Len(buf)$ to be infinite. To be noticed, the above assumption is only for our self-defined $Len(buf)$ in this paper. It is different from $strlen()$, which is a traditional API in C program and re-implemented symbolically in KLEE.

Based on the above buffer information, i.e., the starting address, offset, size and content, we next construct the buffer overflow constraints according to Table 1. For example, in this module, the overflow constraints $Len(src) \geq Size(dest)$ for *strcpy* will be constructed as $\bigwedge_{j=0}^{j=Size(dest)-1} S_{src}[j] \neq \backslash 0$ during symbolic execution. The buffer overflow constraints constructed here correspond to the function $OC(b)$ in the definitions in Section 3. Then, we feed the conjunction of the path constraints and overflow constraints to a constraint solver. Based on the result of constraint solving, we validate the path as an overflowable path, a safe path or an undecided path according to Definitions 3–5 in Section 3. If it is an overflowable path, then the warning will be validated as a true warning, referring to Definition 6 in Section 3. A test case will be generated for it. Then the warning point will be removed from the checking list. Otherwise, the execution will continue until all the paths for the warning being executed or an overflowable path found for the warning. If all the paths for a warning are safe paths, then the warning is a false warning, referring to Definition 7 in Section 3. If there exist some undecided paths besides the false paths, the warning will be regarded as an undecided warning, referring to Definition 8 in Section 3.

In order to improve the efficiency of BovInspector and avoid getting stuck in a loop or a path, we discuss several details of buffer overflow validation as follows. First, if we have validated a warning as a true warning, we will stop exploring other paths for this buffer warning point and report a true warning along with a test case that can follow the execution path and trigger the buffer overflow. Second, if we validate the current path as a safe path, we will check whether the current warning point is in a loop. If so, such execution will be allowed to continue until it reaches the upper bound of the loop or the path is validated as an overflowable path within a loop time threshold, as discussed in Subsection 4.3. Last, we set a time limit for the execution of each path to a buffer overflow warning point. If the procedure exceeds the time bound, it means that the solver cannot determine whether there is a solution to the constraints.

4.5 Targeted Automatic Repair Suggestions

Fixing the validated buffer overflow defects in the program is an urgent task for the developers for the sake of security. It is time consuming to manually fix the existing defects. Therefore, using automatic techniques to supplement manual software development is becoming a trend. However, automatic software repair is challenging because it is a difficult task to figure out where the bug is and how to generate a programmer's preferred fix. Our buffer overflow validation module has solved the "where the bug is" problem. It reports the true buffer overflow warnings, each of which contains the information of the buffer APIs, the size of buffer, the locations of buffer initialization, buffer operations and the overflow point, the test cases to trigger the vulnerability and the corresponding execution paths. The next step is to figure out "how to fix the bug in a way that will be adopted by programmers". Automatic software repair is a promising way to reduce the manual work for programmers. The problem is that programmers tend to lack confidence in the code repaired by some automatic tools. Improving the understandability of the repaired code will be helpful for the programmers to validate and adopt the repair suggestions. Therefore, we perform some empirical studies to survey the officially adopted or programmers' preferred fix approaches for buffer overflow vulnerabilities. By investigating 100 highly ranked buffer overflow CVEs from 2009 to 2014 and the benchmarks from prior buffer overflow detection work, we discover

a total of 11 common repair strategies^[18]. The results show that nearly half of these vulnerabilities (48%) are patched by adding boundary checks, while API replacement and using larger buffer share the second place. All the 11 kinds of repair strategies mentioned in [18] are shown in Table 2 in decreasing order of usage frequency. These strategies fix buffer overflows by adding checks or smashing overflow conditions. Inspired by these official or programmers' preferred repair habits, we adopt these 11 types of repair strategies as repair templates, assemble the final repair codes using the contexts of the buffer overflow vulnerabilities, and provide repair suggestions fully automatically or according to the repair mode selected by the developers.

The specific steps of targeted automatic repair suggestions are shown in Algorithm 2. This module takes as input the source code of the target program, the true buffer overflow warnings validated by BovInspector, and a repair mode configured by programmers, which can be empty. The repair may insert some statements in the code, which will result in a line number mismatch between the source code and the report of the other warnings. To keep the consistency of line numbers by handling true warnings from the bottom to the top, we first sort all the true warnings in each source file by the line numbers in descending order. For each validated true buffer overflow warning, we extract the necessary information for repair, such as buffer API, buffer size, definition location of the buffer, and location of the buffer API, and so on. We first query the usable repair modes for the buffer

Table 2. Buffer Overflow Repair Templates

Rank	Repair Mode	Repair Strategy	Template	Allowed Context
1	ABC	Adding boundary check	if (Overflow Constraints) exit (error);	Any context allowed
2	AR	API replacement	<i>snprintf(dest, Size(dest), format, ...)</i>	<i>strcpy/sprintf/strcat/vsprintf/scanf/vscanf/sscanf/vsscanf/fscanf/vfscanf</i>
2	ULB	Using larger buffer	static buffer: <i>TYPE buf[newSize]</i> ; dynamic buffer: <i>buf = realloc(buf, newSize)</i> ;	Static buffer: <i>TYPE buf[oldSize]</i> ; dynamic buffer: <i>TYPE* buf</i>
4	FBC	Fixing boundary check	if (Overflow Constraints){...}	Boundary check exists
5	AIC	Adding integer check	if ($i + j < i$ $i + j < j$) exit (error);	<i>buf[i + j] = x</i>
6	ASE	Adding string end	if ($i \geq \text{Size}(buf)$) <i>dest[Size(buf) - 1] = '\0'</i> ;	<i>buf[i] = x</i>
7	AMC	Adding Malloc check	if (! <i>buf</i>) exit (error);	<i>buf = malloc(n)</i>
8	SR	String reformat	<i>snprintf(str, Size(str), "...%s..."</i> , ..., <i>Size(dest) - 1, ...)</i> ; <i>sscanf(src, str, ..., dest...)</i> ;	<i>sscanf(src, "...%s..."</i> , ..., <i>dest...</i>)
9	UUV	Using unsigned value	unsigned <i>i</i> ;	Boundary check $i \geq 0$ is needed
10	LIR	Limiting index range	$i = \min(i, \text{Size}(dest) / \text{sizeof}(*dest) - 1)$; $i = \max(i, 0)$;	<i>buf[i] = x</i>
11	UCL	Using concrete length	<i>bufAPI(dest, src, ConcreteValue)</i> ;	<i>bufAPI(dest, src, n)</i>

Note: For ULB, the *newSize* of static buffer is configured by programmers or other static size analysis tools and is determined at compile time; the *newSize* of dynamic buffer is computed based on the overflow constraints at runtime.

API in Table 3, marked as *UsableRepairModeList* in line 4. If the programmer provides a none empty repair mode and the mode is usable, we clear the modes in *UsableRepairModeList* and use the configured repair mode instead (lines 5–7). Otherwise, we use the modes in the original *UsableRepairModeList* computed at line 4. For each mode in *UsableRepairModeList*, we first check whether the precondition is satisfied by analyzing the context of the buffer operation in the program. The allowed contexts are shown in the last column in Table 2. For repair mode SR, only *sscanf* is listed as an example in Table 2. Other APIs, such as *scanf*, *vscanf*, *sscanf*, *vsscanf*, *fscanf* and *vfscanf*, use the corresponding allowed contexts and templates. The allowed context to use SR to repair *sscanf* is: the buffer API *sscanf* is storing characters to the buffer *dest* with a string format like “%s”. If the allowed context is satisfied, we then query the template for the repair mode. For repair mode AR, the safe APIs used as templates are shown in Table 4. Next, we assemble the repair code based on the template and the contexts of the buffer API in the source code. For repair mode SR, the core idea is to control how many bytes are written into a buffer by a formatting string. SR will convert the string format “%s” in *sscanf* into “%ns” by applying an *snprintf* operation, where n is $Size(dest) - 1$.

To be noticed, when using BovInspector to repair a program, *sizeof(buf)* is used instead of the value of *Size(buf)* when the statement in function f is accessing a buffer buf that is a locally defined (namely defined in the same function f) static array or a globally defined static array, because by analyzing official repairs, we find “sizeof” is commonly used in such cases. When the above conditions are not satisfied, we will retrieve the size of the buffer from the output of BovInspector, i.e., *Size(buf)*. In the body of the new “if” statement added by BovInspector, we take “return 0;”, “return

false;” or “return NULL;” as the return statement depending on the return type of the function. The return statement we add may be incorrect and analyzing a more precise return statement based on the program context is left as future work.

Algorithm 2. Targeted Automatic Repair Suggestions

Input: *SrcCode*, *TrueWarningSet*, *RepairMode*
Output: *RepairedCodeSuggestions*

```

1 SortedTrueWarningSet = Sort(TrueWarningSet);
2 foreach w ∈ SortedTrueWarningSet do
3   bufAPI, bufSize, bufLoc, bufAPILoc =
   ExtractInfo(w);
4   UsableRepairModeList =
   QueryRepairSolutions(bufAPI);
5   if RepairMode != EMPTY && RepairMode ∈
   UsableRepairModeList then
6     UsableRepairModeList.clear();
7     UsableRepairModeList.add(RepairMode);
8   foreach mode ∈ UsableRepairModeList do
9     isSatisfied =
   CheckPrecondition(mode, bufAPILoc);
10    if isSatisfied then
11      template = QueryRepairTemplate(mode);
12      code = AssembleRepairCode(SrcCode,
   template, bufAPI, bufSize, bufLoc,
   bufAPILoc);
13      RepairedCodeSuggestions.add(code);
14    else
15      ReportError(w, mode);
  
```

This module provides the corresponding repaired codes for each usable repair mode and ranks them in the order of the repair modes listed in Table 2, namely in decreasing order of usage frequency. Programmers can simply take the first one as the repaired code, which is the most commonly used repair mode in official repairs, or they can configure a preferred repair mode. Note that these repair methods only ensure that the buffer overflow will not happen for the current line of code. We recommend users to manually decide whether to adopt the repair suggestions.

Table 3. Buffer Overflow Repair Solution

Type	API	ABC	AR	ULB	FBC	AIC	ASE	AMC	SR	UUV	LIR	UCL
Unbounded content sensitive buffer operations	<i>strcpy</i> , <i>sprintf</i> , <i>strcat</i> , <i>vsprintf</i>	✓	✓	✓	✓			✓				
	<i>scanf</i> , <i>sscanf</i> , <i>fscanf</i>	✓	✓	✓	✓			✓	✓			
	<i>vscanf</i> , <i>vsscanf</i> , <i>vfscanf</i>											
Bounded content sensitive buffer operations	<i>strncpy</i> , <i>snprintf</i> , <i>fgets</i> , <i>vsprintf</i> , <i>fread</i> , <i>read</i>	✓		✓	✓			✓			✓	✓
	<i>strncat</i>	✓		✓	✓			✓				✓
	<i>memcpy</i> , <i>memmove</i> , <i>memset</i>	✓		✓	✓			✓			✓	✓
Direct buffer accesses	<i>buf[i]</i> , <i>*(buf + i)</i>	✓		✓	✓	✓	✓	✓		✓	✓	

Table 4. Safe API Used in Repair Mode AR

API	Safe API Option
<i>strcpy</i>	<i>strncpy(dest, src, Size(dest) - 1), snprintf(dest, Size(dest), "%s", src)</i>
<i>strcat</i>	<i>snprintf(dest + strlen(dest), Size(dest) - strlen(dest), "%s", src)</i>
<i>sprintf, vsprintf</i>	<i>snprintf(str, Size(str), format, ...)</i>
<i>scanf, vscanf</i>	<i>scanf_s, vscanf_s</i>
<i>sscanf, vsscanf</i>	<i>sscanf_s, vsscanf_s</i>
<i>fscanf, vfscanf</i>	<i>fscanf_s, vfscanf_s</i>

5 Implementation and Evaluation

Implementation. We implemented BovInspector by extending the tool proposed in our previous work [15]. The tool is based on LLVM 2.9^⑥ and KLEE^⑦. We use the commercial Fortify 3.2 as the static buffer overflow detector. For other static analyzers, BovInspector provides a unified static buffer overflow warning report format. It is straightforward to convert the output of another static analyzer to this format.

Evaluation Goals. For evaluation, we hope to answer the following key research questions.

RQ1. Is our validation technique effective and precise for classifying static buffer overflow warnings?

RQ2. Is the performance and scalability of our technique acceptable on real-world applications?

RQ3. Is our automatic repair technique effective?

Metrics. The symbols used in the metrics are defined in Table 5.

Table 5. Symbols Used in the Metrics

Symbol	Meaning
<i>#L</i>	Number of lines of code
<i>#W</i>	Number of buffer overflow warnings
<i>#P</i>	Number of buffer overflow warning paths
<i>#true</i>	Number of true warnings validated by BovInspector
<i>#false</i>	Number of false warnings validated by BovInspector
<i>#undecided</i>	Number of undecided warnings that cannot be validated by BovInspector
<i>#FP</i>	Number of false positives
<i>#fixed</i>	Number of warnings that will not be reported by the static analysis tool Fortify after being repaired by BovInspector

For RQ1, we first define *DCR* and *UCR* to evaluate the effectiveness of our approach.

Decidable classifying ratio (*DCR*) of warning validation is the percentage of warnings that are validated as true or false warnings by BovInspector in all the static buffer overflow warnings.

$$DCR = \frac{\#true + \#false}{\#W}$$

Likewise, *DCR* of warning path validation is the percentage of warning paths that are validated by BovInspector.

Undecidable classifying ratio (*UCR*) for warning validation is the percentage of warnings that are undecided warnings that cannot be validated by BovInspector in all the static buffer overflow warnings.

$$UCR = \frac{\#undecided}{\#W}$$

Likewise, *UCR* of warning path validation is the percentage of warning paths that cannot be validated by BovInspector.

For RQ2, we use time consumptions and memory consumptions to measure the performance of our approach in the three stages of reachability analysis, symbolic execution, and buffer overflow validation.

For RQ3, the recall ratio of static analysis (RRSA) is used to measure the percentage of unfixed warnings in the validated true buffer overflow warnings.

$$RRSA = \frac{\#true - \#fixed}{\#true}$$

Experimental Setup. To prepare the benchmark, we selected eight programs from GNU COREUTILS utilities^⑧ and real-world open source programs such as sendmail-8.12.7^⑨. All experiments were conducted on a quad-core machine with an Intel Core™ Corei5-2400 3.10 GHz processor and 4 G memory, running Linux 3.11.0.

5.1 RQ1: Effectiveness and Precision of the Validation of Static Buffer Overflow Warnings by BovInspector

To answer RQ1, we performed experiments on both synthetic and real-world programs. Table 6 shows the synthetic programs used in the experiments, which are eight programs from GNU COREUTILS utilities.

^⑥ <https://releases.llvm.org/2.9/docs/ReleaseNotes.html>, Oct. 2020.

^⑦ <https://klee.github.io/>, Oct. 2020.

^⑧ <https://www.gnu.org/software/coreutils/>, Oct. 2020.

^⑨ <http://www.sendmail.org/~ca/email/sm-812.html>, Oct. 2020.

Table 6. Buffer Overflow Validation Result on Synthetic Programs

Program	BovInspector Input			BovInspector Output									
				Warning Points					Warning Paths				
	#L ($\times 10^3$)	#W	#P	#true	#FP of #true	#false	#FP of #false	DCR (%)	#true	#FP of #true	#false	#FP of #false	DCR (%)
chmod	0.6	12	12	7	0	4	0	92	7	0	4	0	92
tr	1.9	8	15	3	0	3	0	75	6	0	5	0	73
pwd	0.4	12	16	3	0	7	0	83	3	0	10	0	81
sort	3.3	22	27	8	0	7	0	68	10	0	9	0	70
su	0.6	6	8	1	0	3	0	67	2	0	3	0	63
ls	4.6	31	42	4	0	19	0	74	10	0	22	0	76
pr	2.9	29	31	11	0	15	0	90	13	0	15	0	90
df	1.0	18	18	9	0	9	0	100	9	0	9	0	100

These programs are relatively small and we can manually check the results.

To increase the number of static analysis warnings for each program, we manually inserted or removed code snippets in several random positions in the program. These code snippets will make Fortify report buffer overflow warnings, including both true and false buffer overflow warnings. True buffer overflow warnings are mainly introduced by removing the existing boundary checks for insecure buffer manipulations. False buffer overflow warnings are introduced by inserting

buffer manipulations with infeasible paths. More details can be found in the full version paper^⑩.

Table 7 shows the six real-world programs used in the experiments. The six programs were used in previous buffer overflow detection work [7, 10, 20–22]. Each real-world buffer overflow is located by referring to its CVE descriptions, comparing the current version with the repaired version and reading the log of the repaired version. The details of the real buffer overflows are shown in Table 8.

Table 7. Buffer Overflow Validation Result on Real Programs

Program	BovInspector Input			BovInspector Output									
				Warning Points					Warning Paths				
	#L ($\times 10^3$)	#W	#P	#true	#FP of #true	#false	#FP of #false	DCR (%)	#true	#FP of #true	#false	#FP of #false	DCR (%)
polymorph-0.40	0.3	11	19	0	0	11	0	100	0	0	19	0	100
bc1.06	9.7	5	15	1	0	3	0	80	1	0	11	0	80
net-tool1.46	8.1	62	1256	6	0	36	0	68	251	0	763	0	81
wwwcount2.3	8.3	20	112	1	0	5	0	30	1	0	50	0	45
gzip1.2.4	5.1	19	237	1	0	12	0	68	4	0	148	0	64
sendmail8.12.7	78.1	96	4854	1	0	51	0	54	1	0	3067	0	63

Table 8. Buffer Overflow Details in Real Programs

Program	API	Buffer API Location	Reference	Official Repair Version
polymorph-0.40	<i>strcpy</i>	polymorph.c:118	EDB-ID:22633	ϕ
bc-1.06	array access	main.c 188	[7]	ϕ
net-tools-1.46	<i>strcat</i>	netstat.c:450, 457, 602, 608, 737, 743	N/A	ϕ
wwwcount-2.3	<i>strcpy</i>	parse.c:840	CVE-1999-0021	wwwcount-2.5
gzip-1.2.4	<i>strcpy</i>	gzip.c:1009	CVE-2001-1228	gzip-1.3.9
sendmail-8.12.7	array access	headers.c:1337-1342	CVE-2002-1337	sendmail-8.12.8

Note: ϕ means we have not found an official repair for the buffer overflow. N/A means the reference of the report of the buffer overflow defect is unavailable.

^⑩<https://bovinspectorool1.github.io/project/>, Oct. 2020.

Table 6 shows the results on synthetic programs and Table 7 shows the results on real programs. The first four columns show the name of the programs and the basic information of the input of BovInspector, i.e., $\#L$, $\#W$ and $\#P$. From the fifth column to the 14th column, we show the result for warning points and warning paths. For each of them, we show five types of output of BovInspector, i.e., $\#true$ and its $\#FP$, $\#false$ and its $\#FP$, and DCR .

Result 1. BovInspector can significantly reduce the number of buffer overflow warnings to be manually checked due to its ability of validating static buffer overflow warnings. From Table 6 and Table 7, we can find that BovInspector can effectively validate most of the warning points and the warning paths. For warning points, the DCR on the synthetic programs ranges from 67% to 100% and the DCR for all the synthetic programs together is 82%; the DCR on the real programs ranges from 30% to 100% and the DCR for all the real programs together is 60%. For warning paths, the DCR on the synthetic programs ranges from 63% to 100% with an average of 81%; the DCR on the real programs ranges from 45% to 100% with an average of 67%. In general, the UCR on real-world programs is 40%. In other words, there are about 40% of static warnings, namely undecided static warnings, that need to be manually checked by programmers in real-world programs. It means BovInspector can significantly reduce the number of buffer overflow warnings to be manually checked.

Result 2. In practice, the true warnings and the false warnings identified by BovInspector are all correct. For the synthetic programs in Table 6, we manually examined all the buffer overflow warnings reported by Fortify, recorded the verified buffer overflow list and compared them with the true buffer overflow warnings and false buffer overflow warnings validated by BovInspector. As shown in Table 6, columns $\#FP$ of $\#true$ and $\#FP$ of $\#false$ record the number of false positives

for the true warnings and that for the false warnings validated by BovInspector, respectively. The values in these two columns are all 0, which means the true warnings and false warnings identified by BovInspector are all correct. Moreover, there is a one-to-one correspondence between the manually verified buffer overflow list and the warnings validated by BovInspector. By manually validating the warning paths, we find that the true and false warning paths validated by BovInspector are all correct. For the real-world programs in Table 7, after validation, BovInspector will label the warning paths as overflowable paths or safe paths (denoted as $\#true$ and $\#false$ in the “Warning Paths” column) for each warning. We manually checked all the overflowable ($\#true$) paths reported by BovInspector. For the safe ($\#false$) paths reported by BovInspector, we only randomly checked 200 of them or all of them if they are less than 200 by hand. It shows that there is no false positive for the validation of warning paths and warning points by BovInspector. Furthermore, all the true buffer overflows validated by BovInspector are included in the real buffer overflow list, referring to Table 8. We can also see that false warnings account for about 92.2% of all the decided warnings. The experiments on both synthetic and real programs show that the result of our method for decided warnings is reliable.

5.2 RQ2: Performance and Scalability of BovInspector

To evaluate the performance and the scalability of our method on real-world programs, we recorded the time and memory consumptions of BovInspector on the six real-world programs used in Table 7. Moreover, to show the benefit of introducing the warning reachability analysis and guided symbolic execution, we conducted two sets of experiments shown in Table 9. For the first set, we list the statistics of BovInspector. We first list the statistics of the two stages of BovInspector, i.e., $\#RA$ for reachability analysis and $\#GSE + BV$ for

Table 9. Time and Memory Consumption on Real Programs

Program	Time Consumption (s)				Memory Consumption (MB)			
	BovInspector			Unguided	BovInspector			Unguided
	$\#RA$	$\#GSE + BV$	$\#All$		$\#RA$	$\#GSE + BV$	$\#Peak$	
polymorph-0.40	0.1	101.2	101.3	294.2	2.4	36.1	36.1	87.1
bc1.06	3.1	125.9	129.0	201.1	22.2	77.4	77.4	283.4
net-tool1.46	39.8	282.7	322.5	2547.7	68.0	394.5	394.5	512.3
wwwcount2.3	9.1	159.3	168.4	352.6	19.4	225.1	225.1	435.9
gzip1.2.4	16.7	169.2	185.9	453.7	15.1	64.2	64.2	166.2
sendmail8.12.7	145.8	5758.2	5904.0	7124.1	109.2	1435.1	1435.1	3046.9

guided symbolic execution with buffer overflow validation. We then take the time consumptions of $\#RA$ and $\#GSE + BV$ as the time consumptions of BovInspector and list the values in column $\#All$. We recorded the peak memory consumption of $\#RA$ and $\#GSE + BV$ as the memory consumption of BovInspector and list the values in column $\#Peak$. To be noticed, for the item $\#GSE + BV$, we only considered those validated true or false buffer overflow warnings, because the undecided buffer overflow warnings are limited by a time bound, which may disturb the accuracy of the statistics. Symbolic execution will be terminated when all the buffer overflow warning points and warning paths are traversed. For the second set, a comparison experiment was performed to study whether it is also capable of validating static warnings fast without the guidance of warning reachability analysis, namely the column labeled as “Unguided”. In this case, we performed buffer overflow validation by using unguided symbolic execution. Namely, we omitted the warning reachability analysis and simply took as input the buffer overflow warnings reported by Fortify. Actually, to be fair, only the decided warnings validated by BovInspector were fed to the symbolic execution engine in the “Unguided” experiment. Since no path information is available for the “Unguided” method, symbolic execution would be terminated when all the buffer overflow warning points were traversed. In that case, the “Unguided” method may miss a lot of executions to trigger buffer overflow. Therefore, the actual time and memory consumptions of the “Unguided” method would be even larger than the values shown in Table 9.

Result 3. Our method is capable of handling large-scale real-world programs. Table 7 shows the buffer overflow validation results on real programs. Most of the programs have almost 10 000 lines of source code while the largest one reaches 78 100. We observe that BovInspector performs well on most of these programs and validate a large number of false warnings, especially for the warning paths. About 60% of the static buffer overflow warnings are decidable for BovInspector. According to Table 9, BovInspector can finish validating all the decidable static buffer overflow warnings in these real-world programs with an acceptable time consumption and memory consumption.

Result 4. Warning reachability analysis and guided symbolic execution are effective for improving the performance and scalability of BovInspector. By comparing the statistics of BovInspector and the “Unguided” method, we find that the “Unguided” method always

consumes more time and memory than BovInspector. The gap becomes more obvious especially when the scale of the program increases. This means the warning reachability analysis and the guided symbolic execution perform well in saving time and memory by reducing the exploration space of symbolic execution. By comparing the time consumption of BovInspector and “Unguided”, we can see that the guidance of the warning reachability analysis helps save about 17.1% to 87.3% of time, with an average of 37.9%. The result shows that BovInspector consumes about 23.0%–72.7% with an average of 50.7% less memory than the “Unguided” method. Despite we only made a relatively conservative statistics on the “Unguided” method, namely the data for it would be even larger than the value shown in Table 9, BovInspector still costs much less time and memory than the “Unguided” method. In that case, we can conclude that applying warning reachability analysis and guided symbolic execution can effectively reduce the size of search space, which further significantly increases the performance and scalability of BovInspector.

5.3 RQ3: Effectiveness of Automatic Repair

After validating the buffer overflow warnings of the programs in Table 7, we continued to evaluate the effectiveness of the repair module of BovInspector on these programs. After repair, we used Fortify to re-scan the repaired code to see whether the buffer overflow is fixed. To study the reliability of the repair of BovInspector, we made a comparison between the repair codes generated by BovInspector and the official repairs. According to Table 8, there is no official repair version found for the first three programs. To further enrich the repair examples, we repaired another seven real-world programs with known buffer overflows using BovInspector.

Result 5. The targeted automatic repair module can fix buffer overflow in most cases. Table 10 shows the results of Fortify on the original programs and the programs repaired by BovInspector. From the column “Fortify (before)” and BovInspector, we can see that all the buffer overflow warnings are validated as true buffer overflow warnings. The column “Fortify (after)” lists the Fortify results for the four repair modes, i.e., ABC, AR, ULB and LIR. We observe that by applying the ABC repair mode of BovInspector to the programs, three buffer overflow warnings are degraded to a dangerous function, which is a very low-risk warning in Fortify, and the other eight buffer overflow warnings

Table 10. Validation and Repair Results on Real Programs

Program	API	#W	Fortify (Before)	BovInspector	Fortify (After)			
					ABC	AR	ULB	LIR
polymorph-0.40	<i>strcpy</i>	1	Bov	Bov	Dangerous function	ϕ	N/A	N/A
bc-1.06	array access	1	Bov	Bov	ϕ	N/A	Bov	Bov
net-tool-1.46	<i>strcat</i>	6	Bov	Bov	ϕ	ϕ	N/A	N/A
wwwcount-2.3	<i>strcpy</i>	1	Bov	Bov	Dangerous function	ϕ	N/A	N/A
gzip-1.2.4	<i>strcpy</i>	1	Bov	Bov	Dangerous function	ϕ	N/A	N/A
sendmail-8.12.7	array access	1	Bov	Bov	ϕ	N/A	N/A	N/A

Note: Bov means buffer overflow. ϕ means there is no report from Fortify after applying the corresponding repair. N/A means the corresponding repair mode is unavailable for the bug.

are eliminated from the report of Fortify. Therefore, *RRSA* in the ABC mode is 27.3%. The AR mode performs better. All the nine buffer overflow warnings are not reported by Fortify after being repaired by BovInspector in the AR mode and therefore *RRSA* is 0%. All these repaired codes are verified to be correct by manual inspection. For the repaired code of bc-1.06, by applying the ULB and the LIR mode, Fortify still reports buffer overflow warning on it due to the bug finding schemes of Fortify. Through re-validation by

BovInspector and manual inspection, we found that the warning for ULB is a true warning, but the warning for LIR is a false warning.

Result 6. Our repair method is similar to the human programmers' repair habit.

Table 11 shows the results of the BovInspector's repair and the official repair. The "BovInspector Repair" column shows the repair results for the usable repair modes. ABC and AR are usable for *strcpy* and *sprintf*. ABC and LIR are usable for array access. We list all

Table 11. Results of the BovInspector's Repair and the Official Repair

Program	Location	API	BovInspector Repair	Official Repair
wwwcount-2.3	parse.c 840	<i>strcpy</i>	ABC:if(<i>strlen</i> (<i>qs</i>) \geq <i>sizeof</i> (<i>query_string</i>)) return 1; AR: <i>strncpy</i> (<i>query_string</i> , <i>qs</i> , <i>sizeof</i> (<i>query_string</i>) - 1); <i>query_string</i> [<i>sizeof</i> (<i>query_string</i>) - 1] = '\0';	AR: <i>safeStrcpy</i> (<i>query_string</i> , <i>qs</i> , <i>sizeof</i> (<i>query_string</i>) - 1);
gzip-1.2.4	gzip.c 1009	<i>strcpy</i>	ABC:if(<i>strlen</i> (<i>iname</i>) \geq <i>sizeof</i> (<i>iname</i>)) return 1; AR: <i>strncpy</i> (<i>iname</i> , <i>iname</i> , <i>sizeof</i> (<i>iname</i>) - 1); <i>iname</i> [<i>sizeof</i> (<i>iname</i>) - 1] = '\0';	ABC: if(<i>sizeof iname</i> - 1 \leq <i>strlen</i> (<i>iname</i>))
sendmail-8.12.7	headers.c 1337	array	ABC:if(<i>strlen</i> (<i>bp</i>) \geq (MAXNAME + 1)) return <i>bp</i> ;	ABC: if(<i>realqmode</i> && <i>bp</i> < <i>bufend</i>)
man-1.5i2	man.c 299	<i>strcpy</i>	ABC:if(<i>strlen</i> (<i>name0</i>) \geq <i>sizeof</i> (<i>ulname</i>)) return NULL; AR: <i>strncpy</i> (<i>ulname</i> , <i>name0</i> , <i>sizeof</i> (<i>ulname</i>) - 1); <i>ulname</i> [<i>sizeof</i> (<i>ulname</i>) - 1] = '\0';	ABC: if(<i>strlen</i> (<i>name0</i>) \geq <i>sizeof</i> (<i>ulname</i>)) {return <i>name0</i> ;}
wu-ftpd-2.5.0	ftpd.c 1210	<i>strcpy</i>	ABC:if(<i>strlen</i> (<i>mapped_path</i>) \geq 1024) return NULL; AR: <i>strncpy</i> (<i>path</i> , <i>mapped_path</i> , 1023); <i>path</i> [1023] = '\0';	AR: <i>strncpy</i> (<i>path</i> , <i>mapped_path</i> , <i>size</i>); <i>path</i> [<i>size</i> - 1] = '\0';
xmp-2.5.1	dtl_load.c 79	array	ABC:if((<i>i</i> +1)* <i>sizeof</i> (* <i>pofs</i>) > <i>sizeof</i> (* <i>pofs</i>)) return 1; LIR: <i>i</i> = min(<i>i</i> , <i>sizeof</i> (* <i>pofs</i>)/ <i>sizeof</i> (* <i>pofs</i>) - 1);	ABC: if (<i>i</i> < 256)
mapserver-5.2.0 Beta4	mapserv.c 1334	<i>sprintf</i>	ABC:#include "MY_vsnprintf.h" if(<i>MY_vsnprintf</i> ("%s%s%s.map", <i>mapserv</i> -> <i>map</i> -> <i>web.imagepath</i> , <i>mapserv</i> -> <i>map</i> -> <i>name</i> , <i>mapserv</i> -> <i>Id</i>) \geq <i>sizeof</i> (<i>buffer</i>)) return 1; AR: <i>snprintf</i> (<i>buffer</i> , <i>sizeof</i> (<i>buffer</i>), "%s%s%s.map", <i>mapserv</i> -> <i>map</i> -> <i>web.imagepath</i> , <i>mapserv</i> -> <i>map</i> -> <i>name</i> , <i>mapserv</i> -> <i>Id</i>);	AR: <i>snprintf</i> (<i>buffer</i> , <i>sizeof</i> (<i>buffer</i>), "%s%s %s.map", <i>mapserv</i> -> <i>map</i> -> <i>web</i> . <i>imagepath</i> , <i>mapserv</i> -> <i>map</i> -> <i>name</i> , <i>mapserv</i> -> <i>Id</i>);
cgminer-4.3.4	util.c 1883	<i>sprintf</i>	ABC:#include "MY_vsnprintf.h" if(<i>MY_vsnprintf</i> ("%s:%s", <i>url</i> , <i>port</i>) \geq <i>sizeof</i> (<i>address</i>)) return false; AR: <i>snprintf</i> (<i>address</i> , <i>sizeof</i> (<i>address</i>), "%s:%s", <i>url</i> , <i>port</i>);	AR: <i>snprintf</i> (<i>address</i> , 254, "%s:%s", <i>url</i> , <i>port</i>);

Note: In the last two table rows, we use our self-defined function *MY_vsnprintf*() in header *MY_vsnprintf.h*. It calls *vsprintf*(NULL, 0, format, ...) to compute the length of the format string.

the repair results of the usable repair modes to further show how BovInspector repairs buffer overflow vulnerabilities. The last column shows the results of the official repairs. We also present the repair modes used by official repair. As we can see, the code repaired by BovInspector, when using the same repair mode with the official repair, is very similar to the code repaired by official developers. The repair of `man-1.5i2` using repair mode ABC and the repair of `mapserver=5.2.0Beta4` using repair mode AR of BovInspector are exactly the same as the official repairs. For `wwwcount-2.3`, in the AR repair mode, BovInspector repairs it by replacing `strcpy` with `strncpy` and setting the last element to `'\0'`. Official repair uses a `safeStrcpy`, which actually is a wrapper of our repair code. Therefore, they are essentially the same. The repair of `cgminer-4.3.4` using repair mode AR of BovInspector is also essentially the same as the official repair. The repair of `wu-ftpd-2.5.0` using repair mode AR is different from the official repair. In our repair, for the buffer `mapped_path` coming from the parameter, we can only find all its call sites and take the minimum value as the buffer size, which will be too conservative. The developer of `wu-ftpd` changed the definition of the function to make the size of the buffer available directly from a new parameter, i.e., `size`. However, the last argument of `strncpy` is mistakenly set to `size` instead of `size - 1`, which may lead to an off-by-one buffer overflow. For `gzip-1.2.4`, `sendmail-8.12.7` and `xmp-2.5.1`, BovInspector uses “if($a \geq b$)” and the official repairs use “if($a < b$)” when adding boundary checks. It is up to the developers to choose which kind of boundary checks to be added in the programs.

5.4 Discussion

Experiments show that our method works well on buffer overflow detection and false alarm elimination.

Precision of Warning Validation. According to the warning classification rules shown in Subsection 3.2, for a warning, if we find there exists an overflowable path, then a test case that can trigger the buffer overflow will be generated, and thus the warning must be a true warning. If a warning is validated as a false warning, it indicates every reachable path for the warning is a safe path. Based on the warnings with suspicious buffer access information provided by static tools, we only check all the paths that cover these suspicious buffer operations. The manual inspection results have shown that all the validated false warnings by BovInspector

are correct. But BovInspector can also conservatively ignore all the buffer operations in the warning (i.e., $\omega = (\langle \rangle, b)$). Then it will analyze all possible paths from the entry of a program to b . In that case, false warnings validated by BovInspector would be trustworthy. If there exists any path that exceeds the time limit and no overflowable path has been found, the warning will be classified as an undecided warning, which will need further manual inspection.

Our method can automatically validate about 60% of the buffer overflow warnings reported by Fortify for real-world programs, but there are still a lot of undecided ones. By analyzing the undecided cases, we find some main reasons.

Function Pointers. In the warning reachability analysis module, for function pointers, we simply skip the analysis inside the called functions. In that case, all sub-paths of these functions will be explored during symbolic execution, which may lead to path explosion. When the execution times out before exploring all the paths of a warning point, an undecided warning will be reported.

Loops and Branches. Another scenario is that some buffer overflow points are contained in a loop. During guided symbolic execution, our approach allows an execution to continue if the buffer overflow point is in a loop and the current overflowable path constraints cannot be solved. Then the executions can explore the rest of the loop. When an execution reaches the exit of the loop, there are two directions to be selected: one leads to the outer part of loop body, while the other jumps back to the loop entrance. In the warning reachability module, since the buffer overflow point is in the loop, the entrance of the loop is contained in the warning path. In that case, the execution will continue to explore the loop again, until the loop reaches the loop upper bound (which means the path constraints of the entrance branch cannot be solved) or the time upper bound, or it finds a solution to the overflowable path constraints at the buffer overflow point (which means it is a true warning). Some buffer overflows can only be triggered by a specific number of loop iterations, which will lead to the path explosion problem. Moreover, the very existence of branches and loops will also increase the complexity of constraints. Therefore, sometimes, it may be impossible to solve constraints within a reasonable amount of time. In these situations, our method may not be able to decide whether corresponding warnings are true or false within the time limit and then an undecided warning will be reported.

Library Calls. The warning reachability module skips the analysis inside library functions, and without the guidance of warning reachability analysis, the symbolic execution engine will have to explore all the paths in the library functions if the library functions have been modeled by the symbolic execution engine. The number of library functions to be explored will grow exponentially with the number of branches, which will further increase the number of branches and may eventually lead to undecided warnings. Besides, a warning point may not be reachable by the symbolic execution engine due to unknown library calls. When an execution path contains calls to the library functions that have not been modeled by the symbolic execution engine, the engine will terminate the path that cannot be explored further. Besides, because there is a limitation in KLEE for accessing some parameters during symbolic execution or processing multi-threads and multi-processes systems, some buffer operations are not supported in BovInspector. In these situations, BovInspector will also treat the corresponding warning as an undecided one.

6 Related Work

Prior work on identifying buffer overflow vulnerabilities falls into two categories: static program analysis^[3-7,9] and dynamic execution analysis^[10-14]. Besides, there are some researches on guided symbolic execution for test case generation^[23-26] and automatic bug repair^[27-30].

6.1 Static Program Analysis

Most static program analysis approaches scan software source code to discover the code segments that are possibly vulnerable to buffer overflow attacks. Each vulnerability warning needs to be manually inspected to check whether it is indeed a true vulnerability. ITS4^[3] scans C or C++ source code, breaks the codes into lexical tokens, and then matches patterns in the token stream to find possible vulnerabilities. Other similar tools include FlawFinder^①, which has more detailed report and supports more source code types. These schemes only consider the lexical information although they are simple and can be easily applied to large-scale programs.

Some schemes perform semantic analysis. BOON^[4] focuses on string operations. By checking whether an

operation can make the range of a variable outside the boundaries, BOON can report possible buffer overflow vulnerabilities. Splint^[5] requires users to add source annotation to apply inter-procedural analysis and produces warnings for all library functions susceptible to buffer overflow vulnerabilities. These schemes share the same drawback of lacking run-time information so that they often report a large number of false alarms.

To improve accuracy, some schemes introduce path-sensitive analysis. ARCHER^[6] adopts path-sensitive inter-procedural symbolic analysis on program source code, which reduces false alarms since some false positives are caused by infeasible paths. However, it cannot understand C string operations and spends too much resource for checking the paths not related to buffer overflow vulnerabilities. Marple^[7] uses path-sensitive analysis to improve the detection accuracy, and classifies paths into five types: infeasible, safe, vulnerable, overflow-input-independence, and do not know. The main drawback of Marple is that the path-sensitive analysis is static, which means that it cannot identify the buffer overflows that need run-time information. Yamaguchi *et al.*^[9] designed a novel representation of source code called code property graph and with the comprehensive view of source code, their method is able to model different common vulnerabilities more precisely. However, since this scheme does not interpret code, it cannot find vulnerabilities induced by runtime behaviors. AEC^[8] mixes binary analysis with source code analysis to find exploitable vulnerabilities. It uses preconditioned symbolic execution to find bugs at the source code level. Then, it performs dynamic analysis at the binary level with the input generated by symbolic execution to verify whether the vulnerability is exploitable. However, the precision is limited by the preconditions for symbolic execution.

Although many static analysis schemes for buffer overflow have been proposed, it remains an open problem for scaling to large real systems and manually inspecting amounts of warnings. Our work is complementary to such schemes.

6.2 Dynamic Program Analysis

The dynamic program analysis approach inserts special code into software so that buffer overflow occurrences can be detected and properly processed such as terminating software execution. Some dynamic analysis tools, such as StackGuard^[10], add canaries before

① <https://dwheeler.com/flawfinder/>, Oct. 2020.

return addresses in the stack layout to protect entire distributions from stack smashing buffer overflow attacks. Some other methods^[11–13] assume the boundary of variables should not be exceeded by all accesses and monitor the variables to find buffer overflow.

Splat^[14] is a tool for automatically generating test cases for detecting buffer overflows and it performs directed random testing guided by symbolic execution. It uses symbolic length abstractions techniques to prune the state space without losing the buffer overflow detection ability. UndefinedBehaviorSanitizer (UBSan)^⑫ is a fast undefined behavior detector. It modifies the program at compile time to catch various kinds of undefined behaviors during program execution. All dynamic execution approaches of buffer overflows are challenged by generating high-quality test cases to trigger the bugs with limited efforts. Our approach is based on static analysis and symbolic execution; thus it is directed, automatic and cost-effective.

6.3 Guided Symbolic Execution

Guided symbolic execution techniques focus on controlling the procedure of symbolic execution and different methods are proposed to steer the exploration to various parts of the program to tackle the problem of path explosion. Generally, most guided symbolic execution techniques aim to improve the coverage of programs. The control-flow guided search strategy^[23] constructs a weighted control flow graph (CFG), guiding the exploration to the nearest currently uncovered parts based on the distance in the CFG when the concolic testing needs to choose branches to negate. The fitness-guided search strategy^[26] calculates fitness values from explored paths to target predicates and fitness gains for the branches to be flipped, and then selects proper paths and branches to cover the target predicates.

There are also some approaches which guide symbolic execution to specific parts of the program based on various purposes. The tool eXpress^[24] introduces dynamic symbolic execution for regression test generation and prunes paths that do not expose behavioral differences while exploring new program versions. Babić *et al.*^[25] proposed to exploit static analysis to guide dynamic automated test generation for binary programs. This work is similar to ours, but it depends on the visibly pushdown automaton (VPA) generated by the seed tests, which may be not complete.

6.4 Automatic Bug Repair

Automatic bug repair is a promising way of reducing the cost and many methods have been proposed recently. There are general techniques and fault-specific techniques existing for automatic bug repair^[31,32]. GenProg^[27] uses genetic programming to guide the generate-and-validate process to repair defects in C programs. With new representation and mutation and crossover operators, GenProg can scale to large, open-source programs by taking advantages of cloud computing. RSRepair^[28] presents a new automated repair technique using random search instead of genetic search. AE^[29] also focuses on generate-and-validate repair methods. It uses a formal cost model which suggests an improved algorithm for defining and searching the space of patches and the order in which tests are considered. Qi *et al.*^[30] presented a generate-and-validate patch generation system Kali which can achieve the same effect as prior work by just deleting functionality. CodePhage^[33] automatically transfers correct code from donor applications into recipient applications that process the same inputs to successfully eliminate errors in the recipient. SearchRepair^[34] exploits a database of human-written patches encoded as SMT formulas. These techniques try to automatically fix universal defects. Compared with these studies, our method only repairs those validated true buffer overflow warnings, and is more targeted and simpler. Currently, it is specific to buffer overflow vulnerabilities. By introducing more vulnerability models and investigating successful repair practices, our method can be extended for various defects. There are also some approaches that focus on repairing buffer overflow defects. DIRA^[35] automatically instruments a network service program to detect control hijacking and record enough runtime information to generate the corresponding patch. It will extend a buffer according to its runtime information but the patch may be useless under another test case. TAP^[36] is an automatic buffer overflow and integer overflow discovery and patching system. Its application is limited to those programs which contain incorrect checks. ClearView^[37] reallocates the compromised local array as a global array and sandwiches it in a pair of write-protected pages. But its patches are not similar to those that human programmers write. BovInspector generates automatic repair suggestions according to the 11 repair strategies preferred by programmers only for the true buffer overflow warnings.

⑫ <https://clang.llvm.org/docs/UndefinedBehaviorSanitizer.html>, Oct. 2020.

6.5 False Positives Elimination

There are several ways of eliminating false positives reported by static analysis. Ruthruff *et al.* [38] proposed logistic regression models to generate binary classifications of static warnings. Their results show that the generated models are over 85% accurate in predicting false positives. However, this method cannot ensure the validity of the true warnings it predicts. On the contrast, BovInspector is able to further provide a test case that can trigger the buffer overflow for each validated true warning.

The second way is using verification methods [39,40] such as model checking, symbolic execution to verify warnings reported by static analysis. Junker *et al.* [39] presented an abstraction refinement technique to automatically find and eliminate false positives. The analysis starts with a syntactic model according to the static analysis. Then, it iteratively computes the infeasible sub-paths using SMT solvers and refines the model using additional automata. Muske and Khedker [40] proposed an observation-based method to avoid redundancy when using bounded model checking to verify false positives.

The last kind of methods leverages a precise static method to validate warnings reported by an imprecise but fast static method [41–43]. Smoke [41] shows a two-stage method that uses a precise static analysis method to eliminate false positives from the previous, imprecise but fast static method. Kim *et al.* [42] proposed to only perform a more precise analysis on the small fragments of the code that are more relevant to a buffer overflow alarm by invoking an SMT solver. In this way, their method only tries to remove false alarms as many as possible but cannot determine the truth of alarms. Finally, the most related work was proposed by Arzt *et al.* [43], which is a post-analysis step based on symbolic execution to prune infeasible paths from the result of data flow analysis. However, it only verifies one path among all possible paths for each warning. In contrast, our method verifies all possible paths that cover the buffer operations reported by a static tool to avoid false negatives. Moreover, our method models buffer operations to verify not only whether a path is feasible, but also whether its overflowable constraints are satisfiable.

7 Conclusions

In this paper, we made three key contributions. First, we proposed the framework called BovInspector

for automatically validating buffer overflow warnings outputted by existing static program analysis tools and providing suggestions to repair the true warnings. Second, we proposed the method for warning reachability analysis, guided symbolic execution, buffer overflow validation and targeted automatic repair. Third, we implemented BovInspector and evaluated its performance on both synthetic programs and real-world open source programs. The experimental results showed that BovInspector can significantly reduce the number of false alarms in buffer overflow warnings outputted by static program analysis tools.

In the future, we will extend our method to validate and repair other static warnings. More accurate guidance of symbolic execution is also a promising way to improve the efficiency of the validation.

References

- [1] Anderson J P. Computer security technology planning study. Technical Report, Air Force Electronic Systems Division, 1972. <https://apps.dtic.mil/sti/citations/AD0758206>, Oct. 2020.
- [2] Shahzad M, Shafiq M Z, Liu A X. A large scale exploratory analysis of software vulnerability life cycles. In *Proc. the 34th Int. Conference on Software Engineering*, Jun. 2012, pp.771-781.
- [3] Viega J, Bloch J T, Kohno Y, McGraw G. ITS4: A static vulnerability scanner for C and C++ code. In *Proc. the 16th Annual Computer Security Applications Conference*, Dec. 2000, pp.257-267.
- [4] Wagner D A, Foster J S, Brewer E A, Aiken A. A first step towards automated detection of buffer overrun vulnerabilities. In *Proc. the Network and Distributed System Security Symp.*, Feb. 2000.
- [5] Evans D, Larochelle D. Improving security using extensible lightweight static analysis. *IEEE Software*, 2002, 19(1): 42-51.
- [6] Xie Y, Chou A, Engler D. ARCHER: Using symbolic, path-sensitive analysis to detect memory access errors. In *Proc. the 9th European Software Engineering Conference Held Jointly with the 11th ACM SIGSOFT Int. Symp. Foundations of Software Engineering*, Sept. 2003, pp.327-336.
- [7] Le W, Soffa M L. Marple: A demand-driven path-sensitive buffer overflow detector. In *Proc. the 16th ACM SIGSOFT Int. Symp. Foundations of Software Engineering*, Nov. 2008, pp.272-282.
- [8] Avgerinos T, Cha S, Hao B, Brumley D. AEG: Automatic exploit generation. In *Proc. the Network and Distributed System Security Symp.*, Feb. 2011, pp.59-66.
- [9] Yamaguchi F, Golde N, Arp D, Rieck K. Modeling and discovering vulnerabilities with code property graphs. In *Proc. the 2014 IEEE Symp. Security and Privacy*, May 2014, pp.590-604.

- [10] Cowan C, Pu C, Maier D, Walpole J, Bakke P, Beattie S, Grier A, Wagle P, Zhang Q, Hinton H. StackGuard: Automatic adaptive detection and prevention of buffer-overflow attacks. In *Proc. the 7th USENIX Security Symp.*, Jan. 1998, pp.63-78.
- [11] Jones R W, Kelly P H. Backwards-compatible bounds checking for arrays and pointers in C programs. In *Proc. the 3rd International Workshop on Automated Debugging*, May 1997, pp.13-26.
- [12] Wagner D, Dean R. Intrusion detection via static analysis. In *Proc. the 2001 IEEE Symp. Security and Privacy*, May 2001, pp.156-168.
- [13] Haugh E, Bishop M. Testing C programs for buffer overflow vulnerabilities. In *Proc. the Network and Distributed System Security Symp.*, Feb. 2003.
- [14] Xu R G, Godefroid P, Majumdar R. Testing for buffer overflows with length abstraction. In *Proc. the 2008 ACM/SIGSOFT Int. Symp. Software Testing and Analysis*, Jul. 2008, pp.27-38.
- [15] Gao F, Wang L, Li X. BovInspector: Automatic inspection and repair of buffer overflow vulnerabilities. In *Proc. the 31st IEEE/ACM Int. Conference on Automated Software Engineering*, Sept. 2016, pp.786-791.
- [16] Clarke L A. A system to generate test data and symbolically execute programs. *IEEE Trans. Software Engineering*, 1976, 2(3): 215-222.
- [17] Cadar C, Dunbar D, Engler D R et al. KLEE: Unassisted and automatic generation of high-coverage tests for complex systems programs. In *Proc. the 8th USENIX Symp. Operating Systems Design and Implementations*, Dec. 2008, pp.209-224.
- [18] Ye T, Zhang L, Wang L, Li X. An empirical study on detecting and fixing buffer overflow bugs. In *Proc. the IEEE Int. Conference on Software Testing, Verification and Validation*, Apr. 2016, pp.91-101.
- [19] Sinha S, Harrold M J, Rothermel G. Interprocedural control dependence. *ACM Trans. Software Engineering and Methodology*, 2001, 10(2): 209-254.
- [20] Larochelle D, Evans D. Statically detecting likely buffer overflow vulnerabilities. In *Proc. the 10th USENIX Security Symp.*, Aug. 2001, pp.177-190.
- [21] Zitser M, Lippmann R, Leek T. Testing static analysis tools using exploitable buffer overflows from open source code. In *Proc. the 12th ACM SIGSOFT Int. Symp. Foundations of Software Engineering*, Oct. 2004, pp.97-106.
- [22] Lu S, Li Z, Qin F, Tan L, Zhou P, Zhou Y. *BugBench*: Benchmarks for evaluating bug detection tools. In *Proc. the Workshop on the Evaluation of Software Defect Detection Tools*, Jun. 2005.
- [23] Burnim J, Sen K. Heuristics for scalable dynamic test generation. In *Proc. the 23rd IEEE/ACM Int. Conference on Automated Software Engineering*, Sept. 2008, pp.443-446.
- [24] Taneja K, Xie T, Tillmann N, de Halleux J. eXpress: Guided path exploration for efficient regression test generation. In *Proc. the 20th Int. Symp. Software Testing and Analysis*, Jul. 2011, pp.1-11.
- [25] Babić D, Martignoni L, McCamant S, Song D. Statically-directed dynamic automated test generation. In *Proc. the 20th Int. Symp. Software Testing and Analysis*, Jul. 2011, pp.12-22.
- [26] Xie T, Tillmann N, De Halleux J, Schulte W. Fitness-guided path exploration in dynamic symbolic execution. In *Proc. the 2009 IEEE/IFIP Int. Conference on Dependable Systems and Networks*, Jun. 2009, pp.359-368.
- [27] le Goues C, Dewey-Vogt M, Forrest S, Weimer W. A systematic study of automated program repair: Fixing 55 out of 105 bugs for \$8 each. In *Proc. the 34th Int. Conference on Software Engineering*, Jun. 2012, pp.3-13.
- [28] Qi Y, Mao X, Lei Y, Dai Z, Wang C. The strength of random search on automated program repair. In *Proc. the 36th Int. Conference on Software Engineering*, May 2014, pp.254-265.
- [29] Weimer W, Fry Z P, Forrest S. Leveraging program equivalence for adaptive program repair: Models and first results. In *Proc. the 28th IEEE/ACM Int. Conference on Automated Software Engineering*, Nov. 2013, pp.356-366.
- [30] Qi Z, Long F, Achour S, Rinard M. An analysis of patch plausibility and correctness for generate-and-validate patch generation systems. In *Proc. the 2015 Int. Symp. Software Testing and Analysis*, Jul. 2015, pp.24-36.
- [31] Gazzola L, Micucci D, Mariani L. Automatic software repair: A survey. *IEEE Trans. Software Engineering*, 2017, 45(1): 34-67.
- [32] Monperrus M. Automatic software repair: A bibliography. *ACM Computing Surveys*, 2018, 51(1): Article No. 17.
- [33] Sidiroglou-Douskos S, Lahtinen E, Long F, Rinard M. Automatic error elimination by horizontal code transfer across multiple applications. In *Proc. the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation*, Jun. 2015, pp.43-54.
- [34] Ke Y, Stolee K T, le Goues C, Brun Y. Repairing programs with semantic code search (T). In *Proc. the 30th IEEE/ACM Int. Conference on Automated Software Engineering*, Nov. 2015, pp.295-306.
- [35] Smirnov A, Chiueh T C. DIRA: Automatic detection, identification and repair of control-hijacking attacks. In *Proc. the Network and Distributed System Security Symp.*, Feb. 2005.
- [36] Sidiroglou-Douskos S, Lahtinen E, Rinard M. Automatic discovery and patching of buffer and integer overflow errors. Technical Report, Massachusetts Institute of Technology, Cambridge, 2015. <https://dspace.mit.edu/handle/1721.1/97087>, Oct. 2020.
- [37] Perkins J H, Kim S, Larsen S et al. Automatically patching errors in deployed software. In *Proc. the 22nd ACM SIGOPS Symp. Operating Systems Principles*, Oct. 2009, pp.87-102.
- [38] Ruthruff J, Penix J, Morgenthaler J, Elbaum S, Rothermel G. Predicting accurate and actionable static analysis warnings. In *Proc. the 30th ACM/IEEE Int. Conference on Software Engineering*, May 2008, pp.341-350.
- [39] Junker M, Huuck R, Fehnker A, Knapp A. SMT-based false positive elimination in static program analysis. In *Proc. the 14th Int. Conference on Formal Engineering Methods*, Nov. 2012, pp.316-331.
- [40] Muske T, Khedker U P. Efficient elimination of false positives using static analysis. In *Proc. the 26th IEEE Int. Symp. Software Reliability Engineering*, Nov. 2015, pp.270-280.

- [41] Fan G, Wu R, Shi Q, Xiao X, Zhou J, Zhang C. Smoke: Scalable path-sensitive memory leak detection for millions of lines of code. In *Proc. the 41st IEEE/ACM Int. Conference on Software Engineering*, May 2019, pp.72-82.
- [42] Kim Y, Lee J, Han H, Choe K M. Filtering false alarms of buffer overflow analysis using SMT solvers. *Information and Software Technology*, 2010, 52(2): 210-219.
- [43] Arzt S, Rasthofer S, Hahn R, Bodden E. Using targeted symbolic execution for reducing false-positives in dataflow analysis. In *Proc. the 4th ACM SIGPLAN Int. Workshop on State of the Art in Program Analysis*, Jun. 2015, pp.1-6.



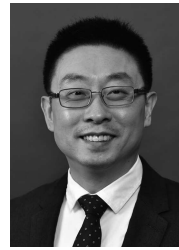
Feng-Juan Gao is a Ph.D. candidate in Nanjing University, Nanjing. She received her B.S. degree in computer science from University of Electronic Science and Technology of China, Chengdu, in 2014. Her research is in software engineering, with focus on symbolic execution.



Yu Wang is a Ph.D. candidate in Nanjing University, Nanjing. He received his B.S. degree in computer science from University of Electronic Science and Technology of China, Chengdu, in 2014. His research is in software engineering, with focus on analyzing concurrent software defects.



Lin-Zhang Wang is a professor in State Key Laboratory for Novel Software Technology, and Department of Computer Science and Technology, Nanjing University, Nanjing. He received his Ph.D. degree in computer science from Nanjing University, Nanjing, in 2005. His research interests include software engineering and software testing.



Zijiang Yang is a professor at Western Michigan University, Kalamazoo. He received his Ph.D. degree in computer science from the University of Pennsylvania, Philadelphia, in 2003. His research interests are in the broad areas of software engineering and formal methods.



Xuan-Dong Li is a professor in State Key Laboratory for Novel Software Technology, and Department of Computer Science and Technology, Nanjing University, Nanjing. He received his B.S. degree in 1985, M.S. degree in 1991, and Ph.D. degree in 1994, all in computer science from Nanjing University, Nanjing. His research interests include software modeling and analysis, software testing and verification.