

# A Slicing-Based Approach for Detecting and Patching Vulnerable Code Clones

1<sup>st</sup> Hakam W. Alomari, 2<sup>nd</sup> Christopher Vendome, 3<sup>rd</sup> Himlal Gyawali

*Department of Computer Science and Software Engineering*

*Miami University, Oxford, Ohio USA*

{alomarhw, vendomcg, gyawalh}@miamioh.edu

**Abstract**—Code cloning is a common practice in software development, but it poses significant security risks by propagating vulnerabilities across cloned segments. To address this challenge, we introduce SRCVUL, a scalable, precise detection approach that combines program slicing with Locality-Sensitive Hashing to identify vulnerable code clones and recommend patches. SRCVUL builds a database of vulnerability-related slices by analyzing known vulnerable programs and their corresponding patches, indexing each slice’s unique structural characteristics as a vulnerability slicing vector. During clone detection, SRCVUL efficiently matches slicing vectors from target programs with those in the database, recommending patches upon identifying similarities. Our evaluation of SRCVUL against three state-of-the-art vulnerable clone detectors demonstrates its accuracy, efficiency, and scalability, achieving 91% precision and 75% recall on established vulnerability databases and open-source repositories. These results highlight SRCVUL’s effectiveness in detecting complex vulnerability patterns across diverse codebases.

**Index Terms**—program slicing, vulnerable code clones, patch recommendations, security vulnerabilities, software maintenance.

## I. INTRODUCTION

As software systems grow increasingly complex, developers are more relying on code cloning to streamline and accelerate development [1], [2]. While code clones can be beneficial when used judiciously [3], [4], they are often seen as suboptimal due to their potential to propagate security vulnerabilities, which increase maintenance costs and reduce software quality [5]–[8]. Vulnerabilities in one instance of code can spread quickly across all cloned segments. Thus, once a vulnerability (e.g., a buffer overflow or use after free) is embedded in a cloned segment, it is likely to persist across multiple locations within or even across projects. This widespread duplication compromises system security, leading to numerous data breaches and network security incidents. A notable example is the OpenSSL Heartbleed vulnerability (CVE-2014-0160) [9], which affected a wide range of systems, including websites, OS distributions, and software applications [1]. This vulnerability spread rapidly because many systems either used the entire OpenSSL library or cloned segments of it, underscoring the urgent need for robust clone vulnerability detection to safeguard system security.

Each year, numerous vulnerabilities are reported by the National Vulnerability Database (NVD) [10], Common Vulnerabilities and Exposures (CVE) [11], and Common Weak-

ness Enumeration (CWE) [12] with over 240,000 CVEs as of October 2024. A significant portion of these vulnerabilities recur due to code cloning and software reuse. Studies across vulnerability databases reveal that recurring vulnerabilities often appear in systems utilizing copied code or common frameworks, where vulnerable code fragments are replicated without prompt application of patches. These reused fragments are structurally similar and retain consistent names for functions and variables, thereby amplifying security risks through cloned, unpatched code [7]. For example, in CVE-2006-3084, 42 systems reused vulnerable code related to privilege escalation, with vulnerabilities identified in 28 of these, including 20 previously unreported cases [7]. The prevalence of code duplication was first highlighted in Baker’s 1995 study [6], which found that 12% of the X Window System’s 714,479 LOC were duplicated, where some subsystems contain up to 20% duplicated code. This redundancy complicates maintenance and creates vulnerabilities when bug fixes applied to one instance of copied code are overlooked in others, resulting in unpatched code.

In this paper, we address the challenge of detecting vulnerable code clones and automatically recommending their patches using variable-level slicing that focuses on vulnerability-related variables ( $vr_{vars}$ ). These variables are identified by comparing the vulnerable code with its patched counterpart, as they are capable of carrying externally controllable data and are particularly susceptible to vulnerabilities when exploited by malicious actions. Essentially, vulnerabilities arise from a combination of insecure syntax on these critical  $vr_{vars}$ , which play a pivotal role in program behavior and can be direct triggers for vulnerabilities. For example, vulnerabilities such as “use after free” occur due to improper memory release associated with a variable [13].

Our approach utilizes decompositional program slicing at the variable level, avoiding the limitations of statement-level slicing methods, which often include all dependencies within a statement and capture both vulnerability-related and unrelated variables. In contrast, our method focuses precisely on  $vr_{vars}$  within vulnerability-related statements, ensuring thorough analysis by capturing all computations (both forward and backward) directly associated with each variable. This comprehensive slicing strategy not only minimizes irrelevant statements, reducing false positives (FPs), but also preserves the broader vulnerability context, reducing the risk of false

negatives (FNs) from overlooking relevant computations. This precision is crucial for vulnerabilities that span non-contiguous code fragments, as vulnerable lines may be distributed across various functions or areas within the codebase.

## II. BACKGROUND AND MOTIVATION

Code cloning involves creating code components that are identical or similar to existing ones, including various types such as *exact* (Type-1), *renamed* (Type-2), *near-miss* (Type-3), and *semantic* (Type-4) clones [14], [15]. Extensive research has been conducted on clone detection [14], [16]–[20], with approaches differing in their *definitions*, *similarity metrics*, and *code analysis depth*. Modern techniques for detecting vulnerable clones, building on traditional methods, continue to encounter significant challenges in both scalability and accuracy [1], [2], [21]–[26].

Scalability issues arise as these methods struggle to handle the sheer volume of code in large software repositories, where millions of LOC are common. Thus, many tools have high processing times and require significant computational resources, limiting their effectiveness for real-time or large-scale vulnerability detection tasks.

Accuracy challenges stem from two main issues. First, many methods tend to return substantial amounts of vulnerability-irrelevant information, as they often operate at granularity levels (e.g., function level) that may not capture the specifics of a vulnerability, leading to over-inclusion and reduced accuracy. Second, detecting near-miss and semantic clones remains difficult because these clones typically involve non-contiguous, reordered, or interwoven statements across different regions of the code. Their structural diversity and distribution throughout the codebase allow them to evade detection by conventional methods, which are often tuned to recognize more straightforward syntactic patterns.

Program slicing is a well-established technique for decomposing software into relevant data and control flow dependencies, facilitating an understanding of underlying program semantics [27], [28]. By tracing these dependencies, slicing identifies both direct and indirect relationships among variables and statements. Various slicing techniques exist [28]–[30], generating different types of slices. However, these methods can struggle with scalability, particularly in large software systems, due to their reliance on Program Dependence Graphs (PDGs) [31] for slice computation, a computationally intensive process that becomes resource-demanding in larger contexts.

Recent research has shown that program slicing can improve the accuracy of vulnerable code clone detection by excluding information irrelevant to vulnerabilities [2], [32], [33]. Additionally, slicing techniques have been applied to enhance general vulnerability detection processes beyond code cloning [34]. Despite its advantages, program slicing remains resource-intensive, though scalable tools like SRCSLICE have been developed to mitigate these challenges [35], [36].

Traditional program slicing considers both control and data dependencies to identify relevant code segments. However, it

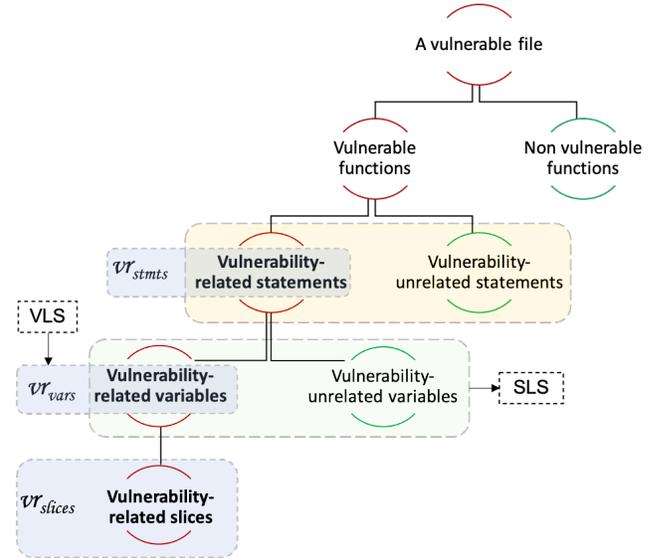


Figure 1. This figure illustrates the comparison between Statement-Level Slicing (SLS) and Variable-Level Slicing (VLS) for identifying vulnerability-related code. The bolded path represents the approach used by SR-CVUL: starting with vulnerability-related statements  $VR_{stmts}$ , we isolate the vulnerability-related variables  $VR_{vars}$  from these statements and then generate vulnerability-related slices  $VR_{slices}$ .

often operates at the statement level, referred to as Statement-Level Slicing (SLS), where dependencies across entire statements are aggregated. While this approach filters out some vulnerability-unrelated statements within function units, as illustrated in Figure 1, it may still include irrelevant information when unrelated variables are present within a statement. This aggregation can introduce noise, reducing the precision needed for detecting vulnerabilities tied to specific variables.

To address this limitation, a more fine-grained approach, Variable-Level Slicing (VLS), focuses specifically on the  $VR_{vars}$ . VLS captures precise data and control dependencies at the variable level, enabling a more targeted analysis of vulnerability-related code segments. By operating at this level of granularity, VLS provides greater accuracy and detail than SLS, improving its utility for detecting vulnerabilities in code.

To illustrate the significance of our work, we present a code example from Song et al.’s work on detecting vulnerable code clones [32]. In the code snippet shown in Figure 2 (left), you can observe a critical vulnerability within the Linux kernel’s core sound driver, which was identified as CVE-2019-15214. In this specific case, the vulnerability relates to a “use after free” scenario occurring in a race condition between disconnection events. This condition could potentially allow a local attacker with the capability to trigger disconnection events, such as hardware removal or addition, to crash the system, corrupt memory, or escalate privileges. The problem arises from the fact that adding and deleting elements in a linked list are performed without proper locking mechanisms. This lack of synchronization becomes problematic when these operations are executed concurrently. To address this vulnerability, the recommended patch, as shown in Figure 2 (right), involves protecting the operations of adding and deleting links

```

696 static struct snd_info_entry *
697 snd_info_create_entry(const char *name,
        struct snd_info_entry *parent,
698         struct module *module)
699 {
700     struct snd_info_entry *entry;
701     entry = kzalloc(sizeof(*entry),
        GFP_KERNEL);
702     if (entry == NULL)
703         return NULL;
704     entry->name = kstrdup(name, GFP_KERNEL);
705     if (entry->name == NULL) {
706         kfree(entry);
707         return NULL;
708     }
709     entry->mode = S_IFREG | 0444;
710     entry->content = SNDRV_INFO_CONTENT_TEXT
        ;
711     mutex_init(&entry->access);
712     INIT_LIST_HEAD(&entry->children);
713     INIT_LIST_HEAD(&entry->list);
714     entry->parent = parent;
715     entry->module = module;
716     if (parent)
717         list_add_tail(&entry->list, &parent->
        children);
718     return entry;
719 }

/* Lines 720 - 787 are omitted here */

788     mutex_unlock(&info_mutex);
789 }
790
791 /* free all children at first */
792 list_for_each_entry_safe(p, n, &entry->
        children, list)
793     snd_info_free_entry(p);
794
795     list_del(&entry->list);
796     kfree(entry->name);
797     if (entry->private_free)
798         entry->private_free(entry);
799     kfree(entry);
800 }

```

```

696 static struct snd_info_entry *
697 snd_info_create_entry(const char *name,
        struct snd_info_entry *parent,
698         struct module *module)
699 {
700     struct snd_info_entry *entry;
701     entry = kzalloc(sizeof(*entry),
        GFP_KERNEL);
702     if (entry == NULL)
703         return NULL;
704     entry->name = kstrdup(name, GFP_KERNEL);
705     if (entry->name == NULL) {
706         kfree(entry);
707         return NULL;
708     }
709     entry->mode = S_IFREG | 0444;
710     entry->content = SNDRV_INFO_CONTENT_TEXT
        ;
711     mutex_init(&entry->access);
712     INIT_LIST_HEAD(&entry->children);
713     INIT_LIST_HEAD(&entry->list);
714     entry->parent = parent;
715     entry->module = module;
716     if (parent) {
717         mutex_lock(&parent->access);
718         list_add_tail(&entry->list, &parent->
        children);
719         mutex_unlock(&parent->access);
720     }
721     return entry;
722 }

/* Lines 723 - 797 are omitted here*/

798     p = entry->parent;
799     if (p) {
800         mutex_lock(&p->access);
801         list_del(&entry->list);
802         mutex_unlock(&p->access);
803     }
804     kfree(entry->name);
805     if (entry->private_free)
806         entry->private_free(entry);
807     kfree(entry);
808 }

```

Figure 2. A “use after free” vulnerability and its patch in linux-5.1.0/sound/core/info.c, associated with CVE-2019-15214. The left side shows the vulnerable code snippet spanning two functions: `snd_info_create_entry` (module size = 24 LOC) and `snd_info_free_entry` (module size = 22 LOC), with highlighted lines representing deleted lines. The right side shows the patched code snippet for the same functions: `snd_info_create_entry` (module size = 27 LOC) and `snd_info_free_entry` (module size = 27 LOC), with highlighted lines representing added lines. The line numbers correspond to those in the original vulnerable and patched files.

with the parent’s mutex.

As illustrated in Figure 2, this vulnerability affects two separate functions: “`snd_info_create_entry`” and “`snd_info_free_entry`.” In the vulnerable code shown in Figure 2 (left), deleted statements at lines #716 and #795 are identified as the *vr\_stmts*, highlighting only `parent` and `entry` as relevant *vr\_vars*. Reporting both functions as entirely vulnerable would be misleading, as it would include vulnerability-unrelated statements and unrelated variables. For example, VUDDY [1], which detects vulnerabilities at the function level, replaces parameters, variables, data types, and

function calls with unified symbols, using function token sequence lengths to filter unrelated code. Functions with matching hashes are marked as vulnerable. Although VUDDY is efficient, it is limited to detecting only minor modifications, such as variable renaming, and struggles with statement insertions or deletions, restricting its scope. Conversely, line-level tools like REDEBUG [37] focus on cloned code lines, making them suitable primarily for Type-2 clones but less effective in capturing the full context of vulnerabilities that extend across multiple lines or functions.

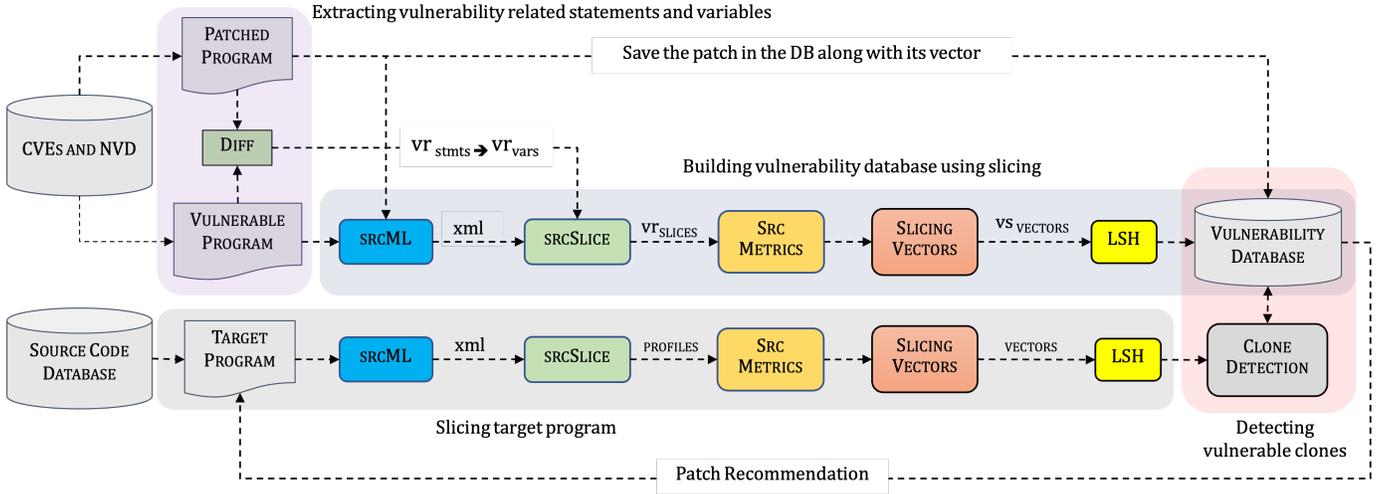


Figure 3. Architecture overview of SRCVUL, illustrating the end-to-end process for detecting vulnerabilities and recommending patches. The workflow involves extracting vulnerability-related code changes, generating slicing-based vectors, and using locality-sensitive hashing for efficient clone detection and patch retrieval from a vulnerability database.

### III. SRCVUL APPROACH

Figure 3 depicts SRCVUL’s pipeline for both vulnerability detection and patch recommendation by utilizing slicing techniques. By focusing on slice-level granularity, SRCVUL captures all relevant computations tied to vulnerabilities, ensuring accuracy without including unnecessary code. The system then optimizes slicing vectors for rapid comparison using efficient hashing methods. Additionally, SRCVUL detects vulnerable code clones, even when code spans non-contiguous sections.

#### A. Building the Vulnerability Database

1) *Identifying and Extracting  $vr\_stmts$* : The initial stage of SRCVUL involves analyzing both the vulnerable and patched versions of a program to identify specific program statements associated with the vulnerability, which we term vulnerability-related statements and denote as  $vr\_stmts$ . To identify these statements, we utilize the Unix utility *diff*, which highlights the added and/or deleted lines between the two versions. This method aligns with established industry practices and is supported by prior research [1], [2], [25], [38].

The identified statements, whether deleted from the vulnerable version or newly added in the patch, are essential for understanding both the root cause of the vulnerability and how the fix addresses it. In other words, the deleted statements include variables directly involved in the vulnerability, providing insight into the issue when we perform slicing on these variables. Conversely, slicing the variables in the added statements enables analysis of the corrected behavior, confirming that the patch resolves the issue and can be applied to similar vulnerable clones. Thus, both deleted and added lines are considered  $vr\_stmts$ . The identified  $vr\_stmts$ , with differentiation between added and deleted, are stored in XML format using SLICEDIFF [39], an extension of the SRCML representation designed to capture code differences and facilitates easy future access.

2) *Determining  $vr\_vars$* : After identifying  $vr\_stmts$ , all variables within these statements are used as criteria for slicing. These vulnerability-related variables, denoted as  $vr\_vars$  (e.g., `parent` and `entry` in Figure 2), are then utilized to derive the vulnerability-related slices, referred to as  $vr\_slices$ .

Studies indicate that many vulnerabilities originate from insecure operations on  $vr\_vars$  [13], [40]–[44], with system API misuse as a primary cause [45]. Traditional slicing tools, especially those based on PDGs and SLS, focus on these operations by marking them as “points of interest” (POIs) or slicing criteria to identify potential vulnerabilities. Arithmetic operations are also flagged, as they can introduce vulnerabilities through unsafe calculations, making them complementary POIs in vulnerability detection [34], [45].

In SRCVUL, the slicing criterion is based on identifying  $vr\_vars$  within  $vr\_stmts$ , assuming these variables are central to the vulnerability. This approach aligns with existing literature, which suggests that vulnerabilities often stem from unsafe operations on critical variables like pointers, system APIs, or externally controlled data. By defining  $vr\_vars$  as the slicing criteria, SRCVUL ensures that slices capture essential information related to specific vulnerability patterns, leading to comprehensive detection across a variety of security threats.

In the following section, we discuss the resulting slices in more detail and provide a rationale for how the slicing fields of these variables map to various critical types of vulnerabilities identified in the literature, highlighting their role in accurate vulnerability detection.

3) *Generating  $vr\_slices$* : After establishing the slicing criteria, SRCVUL extracts  $vr\_slices$ , a set of interconnected statements critical to the vulnerability. Using decompositional program slicing, we generate slices for each identified variable. Each slice isolates code sections tied to specific variables, preserving essential data and control dependencies, even if the statements are non-contiguous [35].

Table I  
VULNERABILITIES ADDRESSED AND LIMITATIONS OF SRCVUL.

Vulnerability Type	Slicing Fields	Description	Reference
<b>Handled by SRCVUL</b>			
Memory Management	<b>Ptrs, Def, Cfuncs</b>	Buffer overflow, use-after-free, double free	[40]
API Misuse	<b>Cfuncs, Use</b>	Unsafe API/system calls (e.g., <code>exec</code> )	[45]
Input Handling	<b>Use, Cfuncs</b>	SQL and Command Injection, XSS	[44]
Authorization Flaw	<b>Dvars, Def, Use</b>	Controls user roles, tokens for access	[37]
Arithmetic & Logic Errors	<b>Use, Dvars</b>	Logic flaws like integer overflow that bypass security	[43]
Concurrency Management	<b>Cfuncs, Ptrs, Dvars</b>	Detects race conditions and synchronization-related calls	[13]
<b>Outside SRCVUL's Scope</b>			
Semantic Flaws	Not addressed	Logic flaws unrelated to variables	—
Complex Control Flow	Partially ( <b>Cfuncs</b> )	Indirect jumps need advanced tools	—
UI Security	Not addressed	Needs UI/UX analysis for phishing, clickjacking	—

To generate the slices, we employ SRCSLICE<sup>1</sup> [36], due to its scalability. Since SRCSLICE leverages the XML representation from SRCML [46], we first utilize SRCML on the vulnerable system. Next, SRCSLICE analyzes this representation to extract information about files, functions, and variables within the system. This information is organized in a three-tiered dictionary, which includes files linked to functions, functions related to variables, and variables associated with slice profiles. The slice profiles generated by SRCSLICE offer a comprehensive view of each identifier in the system, providing insights into their impact on source code statements through data and control dependencies. Specifically, the slice profiles consist of the following fields:

- `file`, `function`, and `variable` names: identify the location of  $vr_{vars}$  within the code, tracing vulnerabilities back to specific files and functions.
- `def`: lists lines where a variable is defined or redefined, capturing allocations, initializations, and deallocations critical for memory management.
- `use`: lists computation lines where the variable's value is used without modification, aiding in identifying vulnerabilities linked to unvalidated inputs, unsafe data handling, and injection points.
- `dvars`: lists variables dependent on the slicing variable's value, supporting the detection of authorization issues and logic errors where security-sensitive variables affect other program parts.
- `ptrs`: lists variables for which the slicing variable acts as a pointer, helping detect memory management vulnerabilities by tracking potential pointer misuse.
- `cfuncs`: lists functions called with the slicing variable as an argument, key for spotting API misuse and system vulnerabilities by revealing unsafe function calls involving  $vr_{vars}$ .

Figure 4 illustrates the slice profiles for the code snippet shown in Figure 2 (left). Each line or entry in the snippet corresponds to a slice profile, structured according to the fields described earlier. After computing these slice profiles, a final pass is executed to account for dependent variables,

```

1 Linux-5.0.10/sound/core/info.c,
  snd_info_create_entry, module, def{698},
  use{715}, dvars{}, ptrs{}, cfuncs{}
2 Linux-5.0.10/sound/core/info.c,
  snd_info_create_entry, entry, def
  {700,701,704,709,710,714,715}, use
  {701,702,705,706,711,712,713,717,718},
  dvars{module,parent}, ptrs{}, cfuncs{
  list_add_tail{1}, INIT_LIST_HEAD{1},
  mutex_init{1}, kfree{1}, kzalloc{2}}
3 Linux-5.0.10/sound/core/info.c,
  snd_info_create_entry, parent, def{697},
  use{714,716,717}, dvars{}, ptrs{}, cfuncs{
  list_add_tail{2}}
4 Linux-5.0.10/sound/core/info.c,
  snd_info_create_entry, name, def{697}, use
  {704}, dvars{}, ptrs{}, cfuncs{kstrdup{1}}
5 Linux-5.0.10/sound/core/info.c,
  snd_info_free_entry, entry, def{779}, use
  {782,784,786,790,792,793,795,796}, dvars
  {}, ptrs{}, cfuncs{private_free{1}, kfree
  {1}, list_del{1}, snd_info_disconnect{1},
  list_for_each_entry_safe{3}}

```

Figure 4. Slice profiles of the code snippet in Figure 2 (left). Lines #2 and #3 represent the slice profiles for the variables `entry` and `parent`, respectively, with other irrelevant variables to  $vr_{stmt}$  omitted for clarity.

function calls, and direct pointer aliasing. SRCSLICE utilizes information from function calls and pointer aliases to calculate indirect and inter-procedural slicing details.

Table I illustrates the mapping between common vulnerability types and the fields within slice profiles, clarifying how each contributes to identifying vulnerability-related attributes. This structured approach enables SRCVUL to focus on high-risk operations effectively, improving detection accuracy and reducing extraneous information within vulnerability slices.

In Figure 2, the variables `parent` and `entry` are identified as  $vr_{vars}$ , so we focus on their slice profiles. SRCSLICE constructs the complete slice by combining the slice profile of each slicing variable, which includes its `def` and `use` sets, with the profiles of related identifiers found in dependent variables (`dvars`), called functions (`cfuncs`), and pointer (`ptrs`) fields. This combination excludes any lines preceding the variable's initial definition (i.e., the set  $1, \dots, def(v) - 1$ ). For instance, as shown in Figure 4, since `parent` has no

<sup>1</sup> <https://www.srcml.org/tools.html>

dependent variables or pointers, only its use in the call to `list_add_tail()` needs consideration. Consequently, the complete slice for `parent` can be represented as:

$$def(parent) \cup use(parent) \cup slice(list\_add\_tail()) - 1$$

yielding  $vr_{slice}(parent) = \{697, 714, 716, 717\}$ . Similarly, for the entry variable, only the instance within `snd_info_free_entry` is treated as a  $vr_{vars}$ , rather than the one in `snd_info_create_entry`. This distinction is marked in Figure 4 by path names, which specify both file and variable names to differentiate between these instances.

4) *Computing Slice-based Metrics*: Next, SRCVUL computes a range of slice-based cognitive complexity metrics using the information from the  $vr_{slices}$ . These metrics provide valuable insights into the code slices and are key to identifying potential vulnerabilities.

For each  $vr_{slice}$ , SRCVUL computes the following metrics:

- *Slice Count (SC)*: This metric quantifies the number of slices in relation to the module size. It indicates the number of slice profiles combined to form the final slice.
- *Slice Size (SZ)*: This metric represents the total number of statements within a complete final slice. It plays a role in the calculation of the *SCvg* metric.
- *Slice Coverage (SCvg)*: This metric measures the slice size relative to the module size, quantified in LOC.
- *Slice Identifier (SI)*: This metric counts the unique identifiers, including variables and method invocations, contained within a slice relative to the module's size. They are drawn from the *dvars*, *ptrs*, and *cfuncs* fields within the slice.
- *Slice Spatial (SS)*: This metric represents the spatial distance, measured in LOC, between the initial definition and the final use of the slicing variable, relative to the module size.

---

**Algorithm 1:** Generate  $vs_{vectors}$  with Metrics

---

```

1: procedure
  GENERATE_VULNERABILITY_VECTORS( $vr_{slices}$ ,
   $module\_size$ )
2:   Initialize  $vs_{vectors}$  as an empty dictionary
3:   foreach  $vr_{slice} \in vr_{slices}$  do
4:     Initialize  $SC$ ,  $SZ$ ,  $SCvg$ ,  $SI$ ,  $SS$  to 0
5:      $SC \leftarrow \frac{|vr_{slice}|}{module\_size}$ 
6:      $SZ \leftarrow \sum (\text{statement count of each } vr_{slice})$ 
7:      $SCvg \leftarrow \frac{SZ}{module\_size}$ 
8:      $SI \leftarrow \frac{\text{count of unique identifiers in } \{dvars \cup ptrs \cup cfuncs\}}{module\_size}$ 
9:      $S_f \leftarrow \text{first statement line in } vr_{slice}$ 
10:     $S_l \leftarrow \text{last statement line in } vr_{slice}$ 
11:     $SS \leftarrow \frac{S_l - S_f}{module\_size}$ 
12:     $vs\_vector \leftarrow \langle SC, SCvg, SI, SS \rangle$ 
13:     $vs_{vectors}[vr_{slice}] \leftarrow vs\_vector$ 
14:   end foreach
15:   return  $vs_{vectors}$ 
16: end procedure

```

---

These metrics were originally introduced by Alqadi et al. [47] and have been adapted with specific modifications for vulnerability detection. Notably, Slice Count (SC) emphasizes constructs like variables, pointers, and function calls instead of all statements. Similarly, Slice Identifier (SI) extends Alqadi's general definition by explicitly focusing on unique identifiers, such as variables, pointers, and function calls, which are particularly relevant in analyzing vulnerabilities. Additionally, Slice Spatial (SS) focuses on variable definition and usage to better highlight patterns associated with vulnerabilities. These modifications make the metrics more effective for identifying vulnerabilities while building on Alqadi's foundational work.

Algorithm 1 depicts the algorithm SRCVUL uses to calculate these metrics. We integrate these metrics into our methodology due to their demonstrated effectiveness [48], in enhancing our understanding of program behavior compared to traditional code-based metrics. Also, these metrics play a pivotal role in the clone detection process, and the metrics are subsequently used to generate slicing vectors for the system's slices.

5) *Encoding and Archiving  $vr_{slices}$* : After computing the slice-based metrics for each  $vr_{slice}$ , SRCVUL proceeds to encode these slices into a set of vulnerability slicing vectors, denoted by  $vs_{vectors}$  using Algorithm 1. These vectors abstract essential characteristics and cognitive complexity metrics into a format that is conducive to analysis, comparison, and efficient storage. The use of  $vs_{vectors}$  enables SRCVUL to effectively compare slices from known vulnerabilities with slices from target code during the vulnerability detection process. Each  $vs_{vector}$  is represented as a single numerical vector with fixed dimensions, where each dimension corresponds to a specific metric of the  $vr_{slice}$ . A possible representation of the  $vs_{vector}$  is as follows:

$$vs_{vector}(v) = \langle SC(v), SCvg(v), SI(v), SS(v) \rangle$$

Consider our motivating example in Figure 2 (left) and the corresponding  $vr_{slices}$  shown in Figure 4. Next, we present the calculations of slice-based metrics for one variable, `parent`, and its  $vr_{vectors}$ . Given that the module spans from line 696 to line 719, the module size is 24 LOC. The calculated metrics for `parent` are as follows:

- **Slice Count (SC)**: For `parent`, with only one slice profile:

$$SC = \frac{\# \text{ Slice Profiles}}{\text{Module Size}} = \frac{1}{24} \approx 0.042$$

- **Slice Size (SZ)**: Based on the final slice for `parent` with lines  $\{697, 714, 716, 717\}$ , we have:

$$SZ = \# \text{ Statements in Complete Final Slice} = 4$$

- **Slice Coverage (SCvg)**: This metric represents the ratio of the slice size to the module size:

$$SCvg = \frac{SZ}{\text{Module Size}} = \frac{4}{24} \approx 0.167$$

- **Slice Identifier (SI)**: For `parent`, there is one unique identifier, `list_add_tail`, so:

$$SI = \frac{\# \text{ Unique Identifies in Slice}}{\text{Module Size}} = \frac{1}{24} \approx 0.042$$

- **Slice Spatial (SS):** For parent, the first statement is at line 697, and the last statement is at line 717:

$$SS = \frac{S_l - S_f}{\text{Module Size}} = \frac{717 - 697}{24} = \frac{20}{24} \approx 0.833$$

Using the computed metrics,

$$vs_{vector}(\text{parent}) = \langle 0.042, 0.167, 0.042, 0.833 \rangle$$

After this process, we’ve transformed the  $vr_{slices}$  into a collection of  $vs_{vectors}$  and stored them in a database. We would perform a similar analysis on our target system to generate vectors and compare them against those in the database to identify potential vulnerabilities and recommend patches.

### B. Slicing the Target Program & Detecting Vulnerable Clones

1) *Slicing Target Program:* For vulnerability detection, SRCVUL generates slicing vectors for the *target program* (i.e., the program with potential security vulnerabilities) through a process similar to that used for building the database. SRCVUL leverages SRCSLICE to consider all available variables within the codebase, ensuring a comprehensive analysis. This process generates program slices that encompass all variables in the code, providing broad coverage of potential vulnerability-related dependencies. Similar to the database construction phase, these slices are processed to compute slice-based metrics and are subsequently transformed into vectors similar to  $vs_{vectors}$ . The resultant slicing vectors are then used during the clone detection process.

While SRCVUL includes all variables in the analysis, the similarity-based clone detection mechanism ensures that only slices with relevant patterns are matched. Slice-based metrics abstract the behavior and roles of variables, allowing the approach to identify relevant slices even when only a subset of function variables is involved in a vulnerability. This mechanism mitigates the impact of unrelated variables on the detection process.

As shown in the evaluation section, such situations do not significantly impede the effectiveness of the approach. SRCVUL demonstrates its capability to identify relevant slices despite the inclusion of all variables, providing robustness while ensuring that potential vulnerabilities tied to less obvious dependencies are not overlooked.

To illustrate this process, we’ve selected Linux kernel version 4.14.76 with approximately 17 MLOC as our target program for analysis. To streamline our presentation, we provide a code snippet in Figure 5, which corresponds to a specific function identified by SRCVUL within this version of the kernel. This function is recognized as a clone of the vulnerable motivation example from Linux kernel version 5.1.0, in Figure 2 (left).

While both functions serve the same purpose, there are subtle distinctions. In Figure 2 (left), an additional argument, `module`, is introduced in the function `snd_info_create_entry` signature (line 698). Also, line 709, `entry->mode` is set to `S_IFREG | 0444`, indicating read permissions for a regular file. In Figure 5,

```

707 static struct snd_info_entry *
708 snd_info_create_entry(const char *name,
709                      struct snd_info_entry *parent)
710 {
711     struct snd_info_entry *entry;
712     entry = kzalloc(sizeof(*entry), GFP_KERNEL
713 );
714     if (entry == NULL)
715         return NULL;
716     entry->name = kstrdup(name, GFP_KERNEL);
717     if (entry->name == NULL) {
718         kfree(entry);
719         return NULL;
720     }
721     entry->mode = S_IFREG | S_IRUGO;
722     entry->content = SNDRV_INFO_CONTENT_TEXT;
723     mutex_init(&entry->access);
724     INIT_LIST_HEAD(&entry->children);
725     INIT_LIST_HEAD(&entry->list);
726     entry->parent = parent;
727     if (parent)
728         list_add_tail(&entry->list, &parent->
729 children);
730     return entry;
731 }

```

Figure 5. Target code snippet for the function “`snd_info_create_entry`”. This snippet is a clone of the vulnerable code shown in Figure 2 (left), with a module size of 22 LOC.

```

1 Linux-4.14.76/sound/core/info.c,
  snd_info_create_entry, entry, def{710,
  711, 714, 719, 720, 724}, use{711, 712,
  715, 716, 721, 722, 723, 726, 727},
  dvars{parent}, pointers{}, cfuncs{
  kzalloc{2}, kstrdup{1}, kfree{1},
  mutex_init{1}, INIT_LIST_HEAD{2},
  list_add_tail{1}}
2 Linux-4.14.76/sound/core/info.c,
  snd_info_create_entry, parent, def{708},
  use{724, 725, 726}, dvars{}, pointers
  {}, cfuncs{list_add_tail{2}}
3 Linux-4.14.76/sound/core/info.c,
  snd_info_create_entry, name, def{708},
  use{714}, dvars{}, pointers{}, cfuncs{
  kstrdup{1}}

```

Figure 6. Slice profiles for target code snippet in Figure 5.

`entry->mode` is set to `S_IFREG | S_IRUGO`, essentially conveying the same meaning as read-only access for a regular file. Finally, in Figure 2 (left), the line `entry->module = module;` (line 715) is specific to the code in Figure 2 (left) and is absent in Figure 5.

The slice profiles for the target code in Figure 5 are presented in Figure 6. Following the same process, the final slice-based metrics and vector for `parent` are calculated as:

$$\langle SC, SCvg, SI, SS \rangle = \langle 0.045, 0.182, 0.045, 0.818 \rangle$$

2) *Detecting Vulnerable Clones:* To detect vulnerable clones, the similarities between slicing vectors in the database and those generated from the target program must be calculated. However, given the potentially large number of  $vs_{vectors}$  in the database and the size of the target software system,

performing pairwise similarity analysis is computationally infeasible. To address this challenge, we implemented Locality Sensitive Hashing (LSH) [49], which has been used in prior clone detection approaches [15], [50]. Our implementation followed the methodology employed by the SRCCLONE clone detector for handling slice clones [15]. LSH hashes vectors into buckets, grouping items that are likely to match, thereby significantly reducing the number of required comparisons.

The LSH implementation in SRCVUL uses cosine similarity, a measure well-suited for capturing structural similarities in slicing vectors [51], with a predefined similarity threshold of 0.8. This threshold means that two slicing vectors must share at least 80% similarity in their MinHash signatures to be considered part of the same clone group or pair. The choice of 0.8 was informed by preliminary experiments, where thresholds ranging from 0.6 to 0.9 were tested. Lower thresholds (e.g., 0.6) improved recall but increased false positives, while higher thresholds (e.g., 0.9) enhanced precision but missed semantically similar vectors. The threshold of 0.8 struck an optimal balance, achieving high precision while maintaining effective vulnerability detection.

Thus, SRCVUL minimizes the number of direct pairwise comparisons required by grouping similar vectors together based on this threshold, ensuring both computational efficiency and accurate vulnerability detection.

To demonstrate this process, consider our motivating example in Figure 2 (left) along with the vulnerable code clone identified in Figure 5. Focusing on the variable of interest, `parent`, we determined that its slicing vector is  $vs_{vector}(\text{parent}) = \langle 0.042, 0.167, 0.042, 0.833 \rangle$ . Now, when examining the slicing vector of `parent` in the vulnerable code clone from Figure 5, we find that  $vs_{vector}(\text{parent}) = \langle 0.045, 0.182, 0.045, 0.818 \rangle$ . Using Cosine similarity, we calculate a similarity between  $v_1$  and  $v_2$  of approximately 0.9994, indicating a very high similarity. This similarity demonstrates that the two vectors are nearly identical in direction, supporting the notion that the two code segments likely represent clones with similar vulnerability characteristics. With such a high similarity, the LSH process effectively hashes these vectors into the same or nearby buckets, enabling SRCVUL to quickly identify and compare them without requiring exhaustive pairwise comparisons.

### C. Patch Recommendation

Previous work on vulnerability detection lacks integrated patch recommendation capabilities [1], [2], [25], [34], [37], [44]. In contrast, SRCVUL combines vulnerability detection with patch recommendation, enabling a streamlined approach to remediating vulnerabilities. By associating each ( $vs_{vector}$ ) with a specific patch, SRCVUL not only identifies vulnerable clones but also provides an efficient patch retrieval.

Each patch in SRCVUL is indexed by its associated slicing vector, which serves as a unique identifier. The vector is derived from the vulnerability-related slices ( $vr_{slices}$ ), representing the structural and semantic characteristics of vulnerable code fragments. For instance, if a vulnerable program’s

---

### Algorithm 2: Detect Vulnerable Clones and Patches

---

```

1: procedure
   DETECT_CLONES_WITH_PATCH(target_program,
   vul_database, LSH_index, similarity_threshold)
2:   Initialize detected_vul_clones as an empty list
3:   target_slices  $\leftarrow$  srcSlice(target_program)
4:   vs_vectors  $\leftarrow$ 
   Generate_Vectors(target_slices, module_size)
5:   foreach (target_slice, vs_vector)  $\in$  vs_vectors do
6:     similar_vectors  $\leftarrow$ 
   LSH_index.Query(vs_vector, similarity_threshold)
7:     foreach db_vector  $\in$  similar_vectors do
8:       Retrieve db_slice and patch from
   vulnerability_database corresponding to db_vector
9:       Calculate
   similarity  $\leftarrow$  Cosine_Similarity(vs_vector, db_vector)
   if similarity  $\geq$  similarity_threshold then
10:    Add (target_slice, db_slice, patch) to
   detected_vulnerable_clones
11:    Recommend patch for target_slice as a fix
12:  end foreach
13: end foreach
14: return detected_vulnerable_clones
15: end procedure

```

---

slicing vector is  $\langle 0.042, 0.167, 0.042, 0.833 \rangle$ , this vector is stored in the database as a lookup key, mapping directly to the associated patch. Because the vulnerabilities are mapped to the patches, SRCVUL is able to retrieve the patches after performing the vulnerability clone detection.

Algorithm 2 summarizes this process of detecting vulnerable clones and recommending patches, where the unique slicing vectors streamline the search and retrieval of specific patches. This integrated approach in SRCVUL provides a practical solution to vulnerability detection and remediation.

## IV. EXPERIMENTAL METHODOLOGY

For vulnerability detection, we assess SRCVUL across three key dimensions: accuracy (precision and recall), performance (execution time), and scalability. Additionally, we evaluate SRCVUL’s capability to recommend patches for detected vulnerabilities. To provide a thorough comparative analysis, SRCVUL is benchmarked against other leading vulnerable clone detectors. For slice-based approaches, we evaluated against VULSLICER [2]; relies on PDGs for slicing source code and is designed for vulnerable code clone detection. We also evaluated against non-slice-based approaches, VUDDY [1] and REDEBUG [37], for a comprehensive assessment.

To systematically investigate SRCVUL’s effectiveness, we address the following research questions:

- **RQ1:** How accurately does SRCVUL identify vulnerable code clones compared to other detection methods?
- **RQ2:** How does SRCVUL perform in terms of scalability in comparison with existing approaches?

- **RQ3:** How effective is SRCVUL in recommending appropriate patches for identified vulnerabilities?

**RQ1** investigates whether the vulnerable code clones detected by SRCVUL are comparable in accuracy and scope to those identified by other methods. **RQ2** investigate the SRCVUL’s scalability and efficiency compared to existing approaches. **RQ3** evaluates the effectiveness of SRCVUL’s patch recommendation, specifically in its ability to suggest suitable fixes for detected vulnerabilities.

**Replication Package:** The replication package can be found at: <https://github.com/alomarhw/srcVul>. It includes the SRCVUL system for detecting and patching vulnerable code clones, organized into modules for slicing, clone detection, and datasets. Detailed instructions for generating slicing vectors, analyzing vulnerabilities, and reproducing results are provided. The evaluation was performed on a computer with a 3.2 GHz 8-Core Intel processor, 32 GB of DDR4 RAM, and running MacOS.

### A. Data Collection

In order to construct a comprehensive and representative database for our assessment, we sourced data from the National Vulnerability Database (NVD) [10], GitHub, VULSLICER [52] and VUDDY [53], selecting relevant vulnerabilities from publicly disclosed Common Vulnerabilities and Exposures (CVEs) [11], [54]. For each CVE, we gathered associated source code files and patch files containing `diff` data, indicating code changes. These files were gathered by systematically accessing the reference links provided within each CVE, which map to source code and patch information. We utilized a web crawler to navigate through reference links and locate pages that provide code diffs or patch information. Each CVE entry in our dataset includes the following metadata: i) the CVE ID and description to contextualize the vulnerability type, ii) the vulnerable code and corresponding patch, and iii) versioning details to align the data with specific project versions.

From this collection process, we collected a total of 65,195 `diff` files. Of these, 60,145 files contain only C and C++ code as either added or deleted lines (excluding comments). Specifically, 13,553 files contain added code lines only, 714 files contain only deleted lines, and 45,878 files contain both added and deleted lines. Overall, these `diff` files represent 3,666 unique CVEs across 512 projects and sub-projects, with multiple `diff` files per CVE where necessary (e.g., for vulnerabilities affecting multiple files or project versions).

In total, 1,145,200 LOC were deleted in the patches, and 2,869,167 LOC were added. This resulted in a total of 2,350,306 `vr_stmts`. This number is not equal to the summation of the added and deleted lines, as lines deleted and then re-added are only counted once. From the `vr_stmts`, we found a total of 2,887,118 `vr_vars`.

Table I depicts the categories of vulnerabilities covered by SRCVUL. We grouped the data by analyzing comments and descriptive content within each CVE’s `diff` file, using

keywords specific to well-known vulnerabilities. The categorization resulted in the following vulnerability mapping: 2,335 *Memory Management*, 686 *Authorization Flaws*, 1,085 *Concurrency Management*, 211 *Input Handling*, 13,470 *API Misuse*, 992 *Logic Errors*, and 41,366 *Uncategorized* (these include vulnerabilities outside the defined categories, which may require further review for precise classification). These 18,779 `diff` files correspond to 1,438 unique CVEs associated with the selected target programs, as discussed in the next section. The slices for the `vr_vars` associated with each CVE are then computed and `vs_vectors` are generated as described in Section III. These vectors are associated with CVE information (e.g., CVE ID, original and patched code, commit hash, etc.)

### B. Design of Experiments

To evaluate the effectiveness (**RQ1**), we evaluated SRCVUL using the same open-source systems that were analyzed in prior work for the tools that we compared against, ensuring a fair and consistent comparison. From VULSLICER, REDEBUG, and VUDDY, we evaluated our approach on the Linux kernel, `libgd`, `Libvirt`, and `Samba`.

To evaluate the accuracy of SRCVUL, we measured the number of true positives (TPs), false positives (FPs), false negatives (FNs), precision, recall, and F1-score with respect to the detection of vulnerable clones within a controlled set of target programs that are not presented in our known vulnerability database. This approach ensures that each detection made by SRCVUL can be directly compared against established vulnerabilities, providing an unbiased and reliable measure of detection accuracy. We selected the following target programs for testing: `libvirt-1.1.0`, `linux-4.14.76`, `libgd-2.3.0`, and `samba-4.0.26`. The number of CVEs from the 1,438 identified in our database for each of these systems is provided in Table II.

SRCVUL was run on each target program to generate the `vs_vectors` and identify clone candidates corresponding to known vulnerability vectors in the database. The ground truth for vulnerabilities in the target systems was established by comparing the detected vulnerabilities against entries in the SRCVUL database, which was built from known vulnerabilities in the CVE/NVD datasets. Matches between `vs_vectors` in the target systems and the database entries were treated as true positives (TPs). If a detected clone did not correspond to any CVE in the database, it was classified as a false positive (FP). Conversely, if a known vulnerability from the database was not detected in the target system, it was considered a false negative (FN). For cases where the similarity between a detected clone and an entry in the database fell within the predefined similarity threshold (e.g., 80%), manual verification was conducted. This manual process involved evaluating whether the detected slice corresponded to the same vulnerability described in the database entry, ensuring accurate classification of true and false positives.

Other state-of-the-art approaches rely on their own matching criteria (e.g., function signatures or line-level granularity) for

Table II

ACCURACY OF SRCVUL IN COMPARISON TO OTHER TECHNIQUES. #CVEs = NUMBER OF CVEs FOUND IN THE VULNERABILITY DATABASE. P = PRECISION, R = RECALL, T = TIME IN SECONDS. LOC = CODE COUNTED USING CLOC NO COMMENTS NOR BLANK.

Target	Files	#CVEs	LOC	SRCVUL					VULSLICER <sup>+</sup>		REDEBUG <sup>-</sup>		VUDDY*	
				TPs	FNs	FPS	P.	R.	P.	R.	P.	R.	P.	R.
kernel	54,745	1311	17,373,016	1000	311	90	0.92	0.76	0.95	0.74	0.56	0.62	0.65	0.56
samba	6,759	71	2,622,156	45	26	8	0.85	0.63	0.88	0.37	0.74	0.32	0.63	0.55
libvirt	2,814	43	2,101,222	30	13	5	0.86	0.70	0.46	0.14	0.78	0.67	0.65	0.60
libgd	530	13	106,498	7	6	3	0.70	0.54	1.00	0.15	0.63	0.40	0.73	0.33
<b>Total</b>	<b>64,848</b>	<b>1438</b>	<b>22,202,892</b>	<b>1082</b>	<b>356</b>	<b>106</b>	<b>0.91</b>	<b>0.75</b>	<b>0.91</b>	<b>0.56</b>	<b>0.65</b>	<b>0.38</b>	<b>0.79</b>	<b>0.42</b>

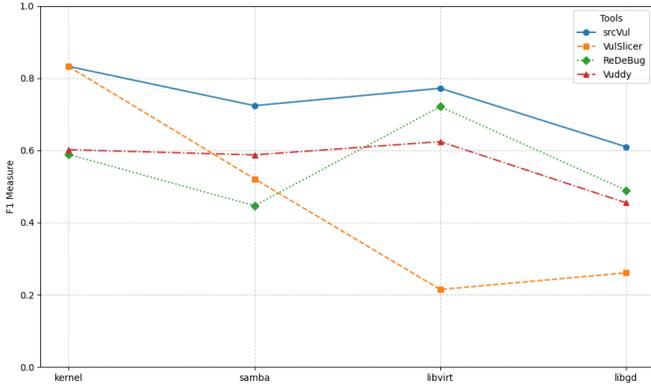


Figure 7. F1 Measure values are shown for SRCVUL and other detectors across target systems.

evaluation. Ground truth vulnerabilities for these approaches are derived from their respective datasets or the CVE/NVD repositories. Metrics such as true positives (TP), false positives (FP), and false negatives (FN) are computed similarly to SRCVUL, based on whether the detected vulnerabilities align with the known vulnerabilities in the respective ground truth datasets. This consistent methodology ensures a fair comparison of detection accuracy across all approaches.

A notable threat to validity in this evaluation is the potential lack of representativeness of the selected target programs relative to the broader set of CVEs in our database. This could mean that some FNs might arise not because SRCVUL failed to detect a clone, but rather because the specific versions of the target programs may not exhibit characteristics or code patterns associated with other CVEs for that system within the database. With these counts, we calculated precision, recall, and the F1-score (the harmonic mean of precision and recall) for each target system. The comprehensive results of SRCVUL, along with comparisons to other clone detection tools, are presented in Table II and Figure 7. As shown, the precision and recall metrics highlight SRCVUL as the most effective tool, with a high precision of 91% and recall of 75%, indicating strong accuracy and coverage in detecting vulnerabilities. In contrast, while VULSLICER achieves similar precision (91%), its recall is lower at 56%, suggesting it may miss more vulnerabilities. REDEBUG and VUDDY have lower precision and recall scores, with REDEBUG at 65% precision and 38% recall, and VUDDY at 79% precision but only 42% recall.

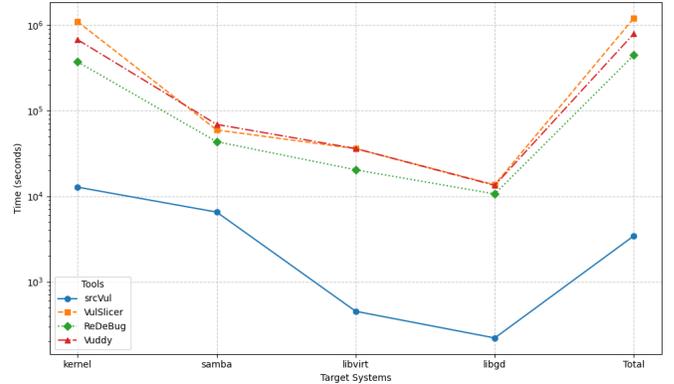


Figure 8. Runtime comparison of SRCVUL and others.

To answer **RQ2**, we measured the runtime for each approach across the target systems. The differences in system size led to variability in the number of entities tracked by each method, such as slices for SRCVUL, function signatures for VUDDY, statements for REDEBUG and VULSLICER. We present the runtime as compared to the number of files, LOC, and CVEs to demonstrate scalability as shown in Figure 8. SRCVUL demonstrates lower runtime across all systems, achieving efficient scalability compared to other techniques.

Lastly, we evaluated **RQ3**, which investigates SRCVUL's ability to recommend patches after identifying vulnerabilities in the target systems. The recommended patches were manually inspected and evaluated. To illustrate this, the  $vs_{vector}$  for parent variable in the target system linux-4.14.76 as shown in Figure 5 was calculated as:  $\langle 0.045, 0.182, 0.045, 0.818 \rangle$ . We found that this vector is similar to the diff file presented in Figure 9. A closer look at the patch, we can observe that `down_read(&key->sem)` and `up_read(&key->sem)` are used to manage concurrent access to a key object using a semaphore. This prevents conflicts when multiple threads attempt to read or modify the key data. In the target code in Figure 5, `mutex_lock(&parent->access)` and `mutex_unlock(&parent->access)` are used to ensure exclusive access to the parent structure while adding an entry to the `parent->children` list. The mutex protects against concurrent modifications that could cause corruption or unexpected behavior in the list. Both vectors deal with synchronization and resource protection in concurrent environments, indicating similar types of race condition vulnerabilities

```

{"CVE_ID": "CVE-2015-7550", "Commit_Hash": "0c47f5e06108c877597de831903d653bcd03a2c2.diff", "Project": "M
↳ -Bab__linux-kernel-amdgpu", "Vector_Hash": "72d4709d92231669d063f88c4a2aa4350767d8346f07bb64-17
↳ fe154a92d3a83f", "Slicing_Vector": {"Slice_Count (SC)": 0.046, "Slice_Coverage (SCvg)": 0.181, "
↳ Slice_Identifier (SI)": 0.041, "Slice_Spatial (SS)": 0.82}, "Patch_Content": "- ret = key_validate
↳ (key);- if (ret == 0) {- ret = -EOPNOTSUPP;- if (key->type->read) {- down_read(&key->sem)
↳ ;+ ret = -EOPNOTSUPP;+ if (key->type->read) {+ down_read(&key->sem);+ ret = key_validate(key
↳ );+ if (ret == 0) ret = key->type->read(key, buffer, buflen);- up_read(&key->sem);- }+
↳ up_read(&key->sem);"}

```

Figure 9. CVE-2015-7550 diff file: improper semaphore use in `keyctl_read_key` (Linux kernel security/keys/keyctl.c pre-4.3.4).

or data corruption risks. Recognizing this pattern in the patch can highlight potential areas to investigate in the target code.

## V. RELATED WORK

Several studies have focused on detecting vulnerable code clones using program slicing techniques. VULSLICER [2] identifies vulnerability-relevant statements (VRS) for known vulnerabilities, slices the target program, and matches the generated slices against a database of vulnerability slices. This approach can detect the first three types of vulnerable code clones. Song et al. [32] also leverage program slicing to extract vulnerable statements within functions, using a hashed database for vulnerability detection. This method is limited to C/C++ programs and can only identify the first three types of clone vulnerabilities.

Some tools detect code clone vulnerabilities without relying on slicing. REDEBUG [37] operates at the line level, scanning source code using a fixed-window size to detect vulnerable patterns. VUDDY [1] identifies vulnerabilities based on function signatures by normalizing variables, function calls, data types, and parameters. Both REDEBUG and VUDDY can detect clones that exactly match known vulnerabilities but may miss those with different syntax or structure. Other tools, like JOANAUDIT [55], employ slicing techniques specifically for injection vulnerabilities, such as SQL. Thome et al. [56] introduced a lightweight slicing-based tool for detecting injection vulnerabilities in Java programs.

Additional static analysis frameworks, such as SFLOW [57], FLOWTWIST [58], ANDROMEDA [59], SOOT [60], LAPSE+ [61], TAJ [62], and MOVERY [22], are also designed to detect security vulnerabilities in Java applications. Recent work has explored thin slicing for vulnerability screening. Dashevskiy et al. [63] adapted it from Sridharan et al. [64] to focus on producer statements that directly contribute to vulnerable computations, reducing slice size while preserving key vulnerability dependencies. However, it struggles with inter-procedural dependencies and may miss vulnerabilities spanning multiple functions and files.

Beyond slicing-based methods, other approaches aim to detect vulnerabilities without slicing but are often limited by different aspects. For example, Livshits et al. [65] and Shankar et al. [66] proposed methods specific to particular programming languages. VULPECKER [25] automates vulnerability detection but faces scalability challenges, while VCCFINDER [67] generates a high rate of FPs, a problem also encountered by REDEBUG and VUDDY. Yamaguchi et al. [68]

introduced machine learning-based approaches for detecting vulnerabilities in C/C++, providing an alternative to traditional methods.

## VI. CONCLUSION

This paper presents SRCVUL, a vulnerability detection approach for large-scale systems that combines variable-level program slicing with semantic representation through slicing vectors ( $vs_{vectors}$ ) to identify vulnerable clones, even with syntactic differences. The key idea of SRCVUL lies in its ability to capture the precise behavior and context of vulnerability-related variables through detailed slice profiles, ensuring improved recall and detection accuracy.

SRCVUL achieves its effectiveness by focusing on slices derived from specific variables rather than function- or line-level granularity. This variable-level granularity allows it to capture all relevant computations and dependencies, reducing the likelihood of missing vulnerabilities. The generation of detailed slice profiles, including data and control dependencies (e.g., def, use, dvars, ptrs, cfuncs), ensures that subtle and complex vulnerability patterns are accurately identified. Additionally, by leveraging Locality-Sensitive Hashing (LSH) for semantic similarity matching, SRCVUL can detect semantically similar slices across diverse codebases, overcoming the limitations of exact-match techniques. Furthermore, the approach validates vulnerability matches by analyzing associated patches, ensuring the identified clones are meaningful and reducing false negatives.

Evaluation across various target systems demonstrates SRCVUL's efficiency, accuracy, and speed, surpassing existing tools. Additionally, SRCVUL uniquely explores patch recommendations by leveraging embedded information in CVEs. While this feature represents a promising step toward automated vulnerability remediation, further research is required to enhance its reliability and effectiveness. Future work will focus on expanding the patch recommendation capabilities by refining similarity-matching techniques and exploring machine-learning models to predict more precise patches.

## REFERENCES

- [1] S. Kim, S. Woo, H. Lee, and H. Oh, "Vuddy: A scalable approach for vulnerable code clone discovery," in *2017 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2017, pp. 595–614.
- [2] S. Salimi and M. Kharrazi, "Vulslicer: Vulnerability detection through code slicing," *Journal of Systems and Software*, vol. 193, p. 111450, 2022.
- [3] C. J. Kapser and M. W. Godfrey, "Cloning Considered Harmful" Considered Harmful: Patterns of Cloning in Software," *Empirical Software Engineering*, vol. 13, no. 6, p. 645, 2008.

- [4] M. Kim, L. Bergman, T. Lau, and D. Notkin, "An ethnographic study of copy and paste programming practices in oopl," in Proceedings. 2004 International Symposium on Empirical Software Engineering, 2004. ISESE'04. IEEE, 2004, pp. 83–92.
- [5] J. Mayrand, C. Leblanc, and E. Merlo, "Experiment on the Automatic Detection of Function Clones in a Software System Using Metrics," in icsm, vol. 96. Monterey, CA, USA: IEEE, 1996, p. 244.
- [6] B. S. Baker, "On Finding Duplication and Near-duplication in Large Software Systems," in Proceedings of 2nd Working Conference on Reverse Engineering. Toronto, Ontario, Canada: IEEE, 1995, pp. 86–95.
- [7] N. H. Pham, T. T. Nguyen, H. A. Nguyen, and T. N. Nguyen, "Detection of recurring software vulnerabilities," in Proceedings of the 25th IEEE/ACM International Conference on Automated Software Engineering, 2010, pp. 447–456.
- [8] H. Li, H. Kwon, J. Kwon, and H. Lee, "Cloriff: software vulnerability discovery using code clone verification," Concurrency and Computation: Practice and Experience, vol. 28, no. 6, pp. 1900–1917, 2016.
- [9] A. C. D. Agency. (2024) Openssl heartbleed vulnerability. <https://www.cisa.gov/news-events/alerts/2014/04/08/openssl-heartbleed-vulnerability-cve-2014-0160> [Accessed: (2024)].
- [10] (2024) National vulnerability database. <https://nvd.nist.gov/> [Accessed: (2024)].
- [11] (2024) Common vulnerabilities and exposures. <https://cve.mitre.org> [Accessed: (2024)].
- [12] (2024) Common weakness enumeration. <https://cwe.mitre.org> [Accessed: (2024)].
- [13] F. Schuster, T. Tendyck, C. Liebchen, L. Davi, A.-R. Sadeghi, and T. Holz, "Counterfeit object-oriented programming: On the difficulty of preventing code reuse attacks in c++ applications," in 2015 IEEE Symposium on Security and Privacy. IEEE, 2015, pp. 745–762.
- [14] C. K. Roy and J. R. Cordy, "A survey on software clone detection research," Queen's School of Computing TR, vol. 541, no. 115, pp. 64–68, 2007.
- [15] H. W. Alomari and M. Stephan, "Clone detection through srcclone: A program slicing based approach," Journal of Systems and Software, vol. 184, p. 111115, 2022.
- [16] S. Bellon, R. Koschke, G. Antonioli, J. Krinke, and E. Merlo, "Comparison and Evaluation of Clone Detection Tools," IEEE Transactions on software engineering, vol. 33, no. 9, pp. 577–591, 2007.
- [17] D. Rattan, R. Bhatia, and M. Singh, "Software Clone Detection: A Systematic Review," Information and Software Technology, vol. 55, no. 7, pp. 1165–1199, 2013.
- [18] A. Sheneamer and J. Kalita, "A Survey of Software Clone Detection Techniques," International Journal of Computer Applications, vol. 137, no. 10, pp. 1–21, 2016.
- [19] J. Svajlenko and C. K. Roy, "Evaluating Modern Clone Detection Tools," in 2014 IEEE International Conference on Software Maintenance and Evolution. Victoria, BC, Canada: IEEE, 2014, pp. 321–330.
- [20] H. W. Alomari and M. Stephan, "Srcclone: Detecting code clones via decompositional slicing," in Proceedings of the 28th International Conference on Program Comprehension, 2020, pp. 274–284.
- [21] F. Yamaguchi, K. Rieck et al., "Vulnerability extrapolation: Assisted discovery of vulnerabilities using machine learning," in 5th USENIX Workshop on Offensive Technologies (WOOT 11), 2011.
- [22] S. Woo, H. Hong, E. Choi, and H. Lee, "{MOVERY}: A precise approach for modified vulnerable code clone discovery from modified {Open-Source} software components," in 31st USENIX Security Symposium (USENIX Security 22), 2022, pp. 3037–3053.
- [23] F. Yamaguchi, N. Golde, D. Arp, and K. Rieck, "Modeling and discovering vulnerabilities with code property graphs," in 2014 IEEE Symposium on Security and Privacy. IEEE, 2014, pp. 590–604.
- [24] F. Yamaguchi, A. Maier, H. Gascon, and K. Rieck, "Automatic inference of search patterns for taint-style vulnerabilities," in 2015 IEEE Symposium on Security and Privacy. IEEE, 2015, pp. 797–812.
- [25] Z. Li, D. Zou, S. Xu, H. Jin, H. Qi, and J. Hu, "Vulpecker: an automated vulnerability detection system based on code similarity analysis," in Proceedings of the 32nd annual conference on computer security applications, 2016, pp. 201–213.
- [26] A. Petukhov and D. Kozlov, "Detecting security vulnerabilities in web applications using dynamic analysis with penetration testing," Computing Systems Lab, Department of Computer Science, Moscow State University, pp. 1–120, 2008.
- [27] M. Weiser, "Program Slicing," IEEE Transactions on software engineering, vol. 0, no. 4, pp. 352–357, 1984.
- [28] F. Tip, A Survey of Program Slicing Techniques. Centrum voor Wiskunde en Informatica Amsterdam, 1994.
- [29] J. Silva, "A Vocabulary of Program Slicing-based Techniques," ACM computing surveys (CSUR), vol. 44, no. 3, p. 12, 2012.
- [30] D. W. Binkley and M. Harman, "A Survey of Empirical Results on Program Slicing," Advances in Computers, vol. 62, no. 105178, pp. 105–178, 2004.
- [31] J. Ferrante, K. J. Ottenstein, and J. D. Warren, "The Program Dependence Graph and its Use in Optimization," ACM Transactions on Programming Languages and Systems (TOPLAS), vol. 9, no. 3, pp. 319–349, 1987.
- [32] X. Song, A. Yu, H. Yu, S. Liu, X. Bai, L. Cai, and D. Meng, "Program slice based vulnerable code clone detection," in 2020 IEEE 19th International Conference on Trust, Security and Privacy in Computing and Communications (TrustCom). IEEE, 2020, pp. 293–300.
- [33] H. Xue, G. Venkataramani, and T. Lan, "Clone-slicer: Detecting Domain Specific Binary Code Clones Through Program Slicing," in Proceedings of the 2018 Workshop on Forming an Ecosystem Around Software Transformation. Canada: ACM, 2018, pp. 27–33.
- [34] T. Wu, L. Chen, G. Du, D. Meng, and G. Shi, "Ultravcs: Ultra-fine-grained variable-based code slicing for automated vulnerability detection," IEEE Transactions on Information Forensics and Security, 2024.
- [35] H. W. Alomari, M. L. Collard, J. I. Maletic, N. Alhindawi, and O. Meqdadi, "srcSlice: Very Efficient and Scalable Forward Static Slicing," Journal of Software: Evolution and Process, vol. 26, no. 11, pp. 931–961, 2014.
- [36] C. D. Newman, T. Sage, M. L. Collard, H. W. Alomari, and J. I. Maletic, "srslice: A tool for efficient static forward slicing," in Proceedings of the 38th International Conference on Software Engineering Companion, 2016, pp. 621–624.
- [37] J. Jang, A. Agrawal, and D. Brumley, "Redebug: finding unpatched code clones in entire os distributions," in 2012 IEEE Symposium on Security and Privacy. IEEE, 2012, pp. 48–62.
- [38] X. Du, B. Chen, Y. Li, J. Guo, Y. Zhou, Y. Liu, and Y. Jiang, "Leopard: Identifying vulnerable code for vulnerability assessment through program metrics," in 2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE). IEEE, 2019, pp. 60–71.
- [39] H. W. Alomari, M. L. Collard, and J. I. Maletic, "A Slice-based Estimation Approach for Maintenance Effort," in 2014 IEEE International Conference on Software Maintenance and Evolution. IEEE, 2014, pp. 81–90.
- [40] C. Cowan, C. Pu, D. Maier, J. Walpole, P. Bakke, S. Beattie, A. Grier, P. Wagle, Q. Zhang, and H. Hinton, "Stackguard: automatic adaptive detection and prevention of buffer-overflow attacks," in USENIX security symposium, vol. 98. San Antonio, TX, 1998, pp. 63–78.
- [41] C. Cowan, F. Wagle, C. Pu, S. Beattie, and J. Walpole, "Buffer overflows: Attacks and defenses for the vulnerability of the decade," in Proceedings DARPA Information Survivability Conference and Exposition. DISCEX'00, vol. 2. IEEE, 2000, pp. 119–129.
- [42] Y. Sui and J. Xue, "Svf: interprocedural static value-flow analysis in llvm," in Proceedings of the 25th international conference on compiler construction, 2016, pp. 265–266.
- [43] L. Szekeres, M. Payer, T. Wei, and D. Song, "Sok: Eternal war in memory," in 2013 IEEE Symposium on Security and Privacy. IEEE, 2013, pp. 48–62.
- [44] Z. Li, D. Zou, S. Xu, X. Ou, H. Jin, S. Wang, Z. Deng, and Y. Zhong, "Vuldeepecker: A deep learning-based system for vulnerability detection," arXiv preprint arXiv:1801.01681, 2018.
- [45] X. Cheng, H. Wang, J. Hua, G. Xu, and Y. Sui, "Deepwukong: Statically detecting software vulnerabilities using deep graph neural network," ACM Transactions on Software Engineering and Methodology (TOSEM), vol. 30, no. 3, pp. 1–33, 2021.
- [46] M. L. Collard, M. J. Decker, and J. I. Maletic, "Lightweight Transformation and Fact Extraction with the srcML Toolkit," in 2011 IEEE 11th international working conference on source code analysis and manipulation. Williamsburg, VI, USA: IEEE, 2011, pp. 173–184.
- [47] B. Alqadi, "The Relationship Between Cognitive Complexity and the Probability of Defects," in 2019 IEEE International Conference on Software Maintenance and Evolution (ICSME). IEEE, 2019, pp. 600–604.

- [48] B. S. Alqadi and J. I. Maletic, "Slice-Based Cognitive Complexity Metrics for Defect Prediction," in 2020 IEEE 27th International Conference on Software Analysis, Evolution and Reengineering (SANER). IEEE, 2020, pp. 411–422.
- [49] P. Indyk and R. Motwani, "Approximate Nearest Neighbors: Towards Removing the Curse of Dimensionality," in Proceedings of the thirtieth annual ACM symposium on Theory of computing. Texas, Dallas, USA: ACM, 1998, pp. 604–613.
- [50] L. Jiang, G. Misserghy, Z. Su, and S. Glondu, "Deckard: scalable and accurate tree-based detection of code clones," in Proceedings of the 29th international conference on Software Engineering, 06 2007, pp. 96–105.
- [51] G. Salton, A. Wong, and C.-S. Yang, "A vector space model for automatic indexing," Communications of the ACM, vol. 18, no. 11, pp. 613–620, 1975.
- [52] (2024) Vulslicer database. <https://zenodo.org/records/6059924> [Accessed: (2024)].
- [53] (2024) Vuddy database. <https://github.com/squizz617/vuddy> [Accessed: (2024)].
- [54] (2024) Cvedetails. <https://www.cvedetails.com> [Accessed: (2024)].
- [55] J. Thomé, L. K. Shar, D. Bianculli, and L. C. Briand, "Joanaudit: A tool for auditing common injection vulnerabilities," in Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering, ser. ESEC/FSE 2017. New York, NY, USA: Association for Computing Machinery, 2017, p. 1004–1008. [Online]. Available: <https://doi.org/10.1145/3106237.3122822>
- [56] J. Thomé, L. K. Shar, and L. Briand, "Security slicing for auditing xml, xpath, and sql injection vulnerabilities," in 2015 IEEE 26th International Symposium on Software Reliability Engineering (ISSRE), 2015, pp. 553–564.
- [57] W. Huang, Y. Dong, and A. Milanova, "Type-based taint analysis for java web applications," in Fundamental Approaches to Software Engineering, S. Gnesi and A. Rensink, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2014, pp. 140–154.
- [58] J. Lerch, B. Hermann, E. Bodden, and M. Mezini, "Flowtwist: Efficient context-sensitive inside-out taint analysis for large codebases," in Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering, ser. FSE 2014. New York, NY, USA: Association for Computing Machinery, 2014, p. 98–108. [Online]. Available: <https://doi.org/10.1145/2635868.2635878>
- [59] O. Tripp, M. Pistoia, P. Cousot, R. Cousot, and S. Guarnieri, "Andromeda: Accurate and scalable security analysis of web applications," in Fundamental Approaches to Software Engineering: 16th International Conference, FASE 2013, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2013, Rome, Italy, March 16-24, 2013. Proceedings 16. Springer, 2013, pp. 210–225.
- [60] R. Vallée-Rai, P. Co, E. Gagnon, L. Hendren, P. Lam, and V. Sundaresan, "Soot - a java bytecode optimization framework," in Proceedings of the 1999 Conference of the Centre for Advanced Studies on Collaborative Research, ser. CASCON '99. IBM Press, 1999, p. 13.
- [61] P. M. Pérez, J. Filipiak, and J. M. Sierra, "Lapse+ static analysis security software: Vulnerabilities detection in java ee applications," in Future Information Technology, J. J. Park, L. T. Yang, and C. Lee, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 148–156.
- [62] O. Tripp, M. Pistoia, S. J. Fink, M. Sridharan, and O. Weisman, "Taj: Effective taint analysis of web applications," in Proceedings of the 30th ACM SIGPLAN Conference on Programming Language Design and Implementation, ser. PLDI '09. New York, NY, USA: Association for Computing Machinery, 2009, p. 87–97. [Online]. Available: <https://doi.org/10.1145/1542476.1542486>
- [63] S. Dashevskiy, A. D. Brucker, and F. Massacci, "A screening test for disclosed vulnerabilities in foss components," IEEE Transactions on Software Engineering, vol. 45, no. 10, pp. 945–966, 2019.
- [64] M. Sridharan, S. J. Fink, and R. Bodík, "Thin slicing," ACM SIGPLAN Notices, vol. 42, no. 6, pp. 112–122, 2007.
- [65] V. B. Livshits and M. S. Lam, "Finding security vulnerabilities in java applications with static analysis," in Proceedings of the 14th Conference on USENIX Security Symposium - Volume 14, ser. SSYM'05. USA: USENIX Association, 2005, p. 18.
- [66] U. Shankar, K. Talwar, J. S. Foster, and D. Wagner, "Detecting format string vulnerabilities with type qualifiers," in 10th USENIX Security Symposium (USENIX Security 01). Washington, D.C.: USENIX Association, Aug. 2001. [Online]. Available: <https://www.usenix.org/conference/10th-usenix-security-symposium/detecting-format-string-vulnerabilities-type-qualifiers>
- [67] H. Perl, S. Dechand, M. Smith, D. Arp, F. Yamaguchi, K. Rieck, S. Fahl, and Y. Acar, "Vccfinder: Finding potential vulnerabilities in open-source projects to assist code audits," in Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security, ser. CCS '15. New York, NY, USA: Association for Computing Machinery, 2015, p. 426–437. [Online]. Available: <https://doi.org/10.1145/2810103.2813604>
- [68] F. Yamaguchi, C. Wressnegger, H. Gascon, and K. Rieck, "Chucky: Exposing missing checks in source code for vulnerability discovery," in Proceedings of the 2013 ACM SIGSAC Conference on Computer & Communications Security, ser. CCS '13. New York, NY, USA: Association for Computing Machinery, 2013, p. 499–510. [Online]. Available: <https://doi.org/10.1145/2508859.2516665>