

# The Secret Life of Software Vulnerabilities: An Empirical Study

Emanuele Iannone · Roberta Guadagni ·  
Filomena Ferrucci · Fabio Palomba

Received: date / Accepted: date

**Abstract** Software vulnerabilities are weaknesses in source code that can be potentially exploited to cause loss or harm. While researchers have been devising a number of methods to deal with vulnerabilities, there is still a noticeable lack of knowledge on their software engineering life cycle, for example how vulnerabilities are introduced and removed by developers. This information can be exploited to design more effective methods for vulnerability prevention and detection, as well as to understand the granularity that these methods should aim at. To investigate the life cycle of software vulnerabilities, we focus on how, when, and under which circumstances vulnerabilities *are introduced* in software projects, as well as whether, after how long, and how they *are removed*. As subjects of our analysis we consider 1,195 vulnerabilities with public patches from the National Vulnerability Database, pertaining to 616 open source software projects on GITHUB, including OPENSLL and LINUX. In particular, we define a six-step process that involves both automated parts (for example, using the SZZ algorithm to find the vulnerability-inducing commits) and manual analyses (how the patches were removed). The investigated vulnerabilities can be classified in 62 categories, take on average 8,248 commits in a file before being introduced, and remain unfixed for a median of 1,546 commits and 1,043 days. Most vulnerabilities are *introduced while developing new features* and *removed through small changes*. Our results have implications on when and how vulnerability detectors should work to better assist developers in early detecting these issues.

---

Emanuele Iannone, Roberta Guadagni, and Filomena Ferrucci, Fabio Palomba  
University of Salerno, Italy  
E-mail: e.iannone16@studenti.unisa.it, r.guadagni1@studenti.unisa.it, fferrucci@unisa.it

## 1 Introduction

Software vulnerabilities are flaws in the design, implementation, or operation management of a software system that can be taken advantage of to break through security policies [82], possibly causing loss or harm [42, 45, 76]. The risks associated with vulnerabilities on software systems are extreme: For instance, in 2017, the WannaCry Ransomware Attack [23] exploited a vulnerability to infect more than 200,000 computers across over 150 countries within a day, with total damages in the range of hundreds of millions USD. Perhaps more importantly, it is estimated that by 2021 vulnerabilities will cost to businesses and users over six trillion USD [10] and will affect an even larger variety of software systems, ranging from public physical tests [25] to data-intensive applications [34], and more. For these reasons, keeping software security under control is still one of the main concerns of developers and organizations [26, 64].

To address this concern, software development organizations and teams are adopting the McGraw’s [60] advice to “*build security in*” rather than waiting until vulnerabilities are discovered in running software [66]. At the same time, researchers have been proposing novel methods and tools to assist developers and facilitate the identification of software vulnerabilities [27, 53, 58, 59]. Nevertheless, recent empirical studies [37, 39, 57, 67, 70] provided evidence that developers are concerned with vulnerabilities and their potential effects, but such vulnerabilities still often occur in practice.

A possible reason behind the diffuseness of vulnerabilities may be the lack of empirical knowledge of their software engineering life cycle, i.e., when developers introduce vulnerabilities, during which development activities vulnerabilities are more prone to be introduced, or how developers remove vulnerabilities in practice. An improved understanding of these aspects may inform software engineering community and tool vendors on how to better support developers in both the vulnerability identification and fixing process [87], for instance by devising novel mechanisms or adapting currently available detection approaches so that they can be ran in the moment of the development when vulnerabilities are usually introduced. Moreover, the analysis of the life cycle of vulnerabilities can uncover interesting patterns that would be helpful for the deployment of best practices [81] that help improving the development process for security [65].

To bridge such a knowledge gap, in this article we present a large-scale empirical study on the life-cycle of software vulnerabilities. To conduct our investigation, we first mine the National Vulnerability Database (NVD) [11] (the U.S. government repository of standards-based vulnerability management data) to extract a set of vulnerabilities for which their fixing commits are public. Overall, we study 1,195 vulnerabilities of 62 categories across 616 different GITHUB projects. Then, we automatically identify the corresponding vulnerability-inducing commits (also verifying the performance of the automated solution) and investigate (i) when vulnerabilities are introduced, (ii) by whom, (iii) under which circumstances, (iv) how long they remain in the considered software systems, and (v) how the fixing process takes place.

The key results of our empirical study show that developers mostly introduce vulnerabilities while developing new features. Furthermore, vulnerabilities remain in a project’s codebase for on median 1,546 commits and 1,043 days, and are removed by developers through simple code changes, such as escaping functions of HTML entities. Finally, we find that expert developers are more prone to introduce vulnerabilities.

**Structure of the paper.** We organize the rest of this article as follows. In Section 2, we describe the design of our empirical study. In Sections 3 and 4, we present and discuss the results in details and the lessons learned, respectively. Finally, we present the related work in Section 5 and the concluding remarks in Section 6.

## 2 Methodology

This section describes the research goals driving our empirical study, our terminology, the data extraction and validation procedure, as well as the research method to answer each research question. We conclude presenting the threats to the results’ validity.

### 2.1 Research Goals

The *goal* of the study is to investigate the life cycle of software vulnerabilities, with the *purpose* of assessing when they are introduced and fixed by developers as well as the circumstances and reasons behind vulnerabilities introductions. The *perspective* is of both researchers and practitioners: the former are interested in better understanding the phenomenon of software vulnerabilities, from their introduction to removal, with the aim of devising novel identification and refactoring techniques that can provide an improved support to developers; the latter are interested in gathering information on the circumstances that most likely introduce vulnerabilities, so that they could inform their development process.

Our study is structured around four main research questions (**RQs**). We start analyzing *when* vulnerabilities are introduced in source code, particularly to understand whether the introduction takes place with the creation of a new file or during maintenance/evolution activities performed on existing files. Such an analysis may reveal insights on the *moment* in which vulnerability detectors should give feedback to developers [91]; for example, should the vulnerabilities be mainly introduced when the affected file is created, an approach pinpointing potential vulnerabilities as soon as a new file is introduced in the repository would be worthwhile [54, 74]. This reasoning leads to our first research question:

**RQ<sub>1</sub>.** *When are vulnerabilities introduced: during a new file creation or during maintenance/evolution activities on existing files?*

Second, we investigate under which circumstances vulnerabilities are more likely to be introduced by developers. We focus on three aspects: *commit goal*, i.e., the action that a developer is performing when introducing vulnerabilities (e.g., new feature implementation or refactoring activity); *project status*, i.e., the proximity of a vulnerability introduction commit to a new release or the startup of the project; and *developer status*, i.e., whether a programmer introducing a vulnerability has usually a high workload. These three aspects are key to describe the *context* and the *situations* in which developers are more prone to introduce vulnerabilities, two pieces of information that may be exploited by vulnerability detectors and refactoring recommenders to output more precise suggestions [55, 79], but also to provide developers with details that can help understanding the reasons leading to the vulnerability introduction [56, 62]. Hence, we ask:

**RQ<sub>2</sub>.** *Under which circumstances are vulnerabilities introduced?*

After describing the mechanisms leading to the introduction of vulnerabilities, we move toward the understanding of their *longevity* by means of a survival analysis method [63]—we investigate how long vulnerabilities remain in source code. This analysis helps describing the lifespan of vulnerabilities: on the one hand, this is useful to understand for how long such vulnerabilities can be exploited by attackers, thus challenging previous findings in the field [81]; on the other hand, we can gather insights on the *reaction time* of developers, which may suggest the need for additional instruments to alert practitioners of the presence or the potential impact of unfixed vulnerabilities [29, 33]. Thus, we ask our third research question:

**RQ<sub>3</sub>.** *What is the survivability of vulnerabilities?*

Finally, we investigate the fixing of vulnerabilities (e.g., do developers completely remove the vulnerable code or apply specific refactoring operations depending on the type of the vulnerability encountered?). Such an analysis can help researchers understanding how to better support developers with (semi-)automated techniques to remove vulnerabilities or to mitigate their effect [31, 41, 43], e.g., by defining a novel catalog of vulnerability-specific refactoring operations. Therefore, we ask:

**RQ<sub>4</sub>.** *How are vulnerabilities removed?*

## 2.2 Context Selection

The *context* of the study is composed of vulnerabilities and change history information of the affected software projects.

Table 1: Top-20 most frequent vulnerabilities in our dataset, by CWE type.

CWE	Kind	Total
CWE-119	Improper Restriction of Operations within the Bounds of a Memory Buffer	213
CWE-79	Improper Neutralization of Input During Web Page Generation	209
CWE-20	Improper Input Validation	99
NVD-CWE-Other	Other	69
CWE-200	Information Exposure	65
CWE-264	Permission/Privileges/Access Controls Issues	63
CWE-125	Out-of-bounds Read	52
CWE-89	Improper Neutralization of Special Elements used in an SQL Command	51
CWE-476	NULL Pointer Dereference	46
CWE-399	Resource Management Errors	38
CWE-352	Cross-Site Request Forgery	37
CWE-284	Improper Access Control	31
CWE-22	Improper Limitation of a Pathname to a Restricted Directory	24
CWE-190	Integer Overflow or Wraparound	22
CWE-189	Numeric Errors	20
CWE-416	Use After Free	19
CWE-787	Out-of-bounds Write	19
CWE-94	Improper Control of Generation of Code	18
CWE-310	Cryptographic Issues	18
CWE-254	7PK - Security Features	17

**Vulnerabilities:** We analyze vulnerabilities from the National Vulnerability Database (NVD) [11], which was created in 2000 by the U.S. NIST Computer Security Division [22] to collect and provide public information about known vulnerabilities affecting software systems and their causes. We rely on NVD for three main reasons: (1) It includes a comprehensive set of information on vulnerabilities, such as references to security checklists, security-related software flaws, misconfigurations, product names, and additional information such as impact metrics (Common Vulnerability Scoring System - CVSS), vulnerability types (Common Weak Enumeration - CWE), applicability statements (Common Platform Enumeration - CPE); (2) it is continuously improved, because the NVD updates the database entries (called “Common Vulnerabilities and Exposure” - CVEs) that have been changed [11]; and (3) it has been used and validated in various research communities [24, 52, 61, 93].

NVD provides a set of 2,973 vulnerabilities and their data. Among this set, we select those for which both the fixing commit and the project’s GitHub are available. As a result, we take into account 1,195 vulnerabilities of 62 distinct types. We use the fixing commit to identify the vulnerability inducing commits, which allows us to analyze how vulnerabilities are introduced, by whom, and under which circumstances. We also compute information on the vulnerability survivability. Table 1 presents the top-20 most common vulnerabilities we analyzed by type; the complete list is available in our online appendix [15].

Table 2: Top-15 projects in our dataset by number of vulnerabilities. Commits = Number of commits in GitHub; Files = Num. of files after the latest commit; KLOC = KLOC after the latest commit; Developers = Number of distinct committers; Vulnerabilities = Number of analyzed vulnerabilities.

Project	Commits	Files	KLOC	Developers	Vulnerab.
TCPdump	5,577	264	100	150	60
PHPmyadmin	115,032	1,565	1,102	1,600	49
FFmpeg	93,463	3,903	1,178	1,750	48
ImageMagick	15,441	1,017	535	71	44
WordPress	39,252	1,832	622	90	30
Radare2	21,086	2,319	665	766	30
KRB5	19,542	1,828	381	123	25
PHP-src	111,636	2,270	1,215	962	24
Jenkins	28,060	3,007	253	811	21
Kanboard	3,854	2,496	191	319	19
File	3,381	68	17	9	19
Exponent-cms	8,602	5,902	1,099	15	16
Dolibarr	74,858	6,355	1,452	380	14
phpMyFAQ	9,132	390	82	61	13
LibTiff	3,006	428	140	16	13

A total of 82 of the studied vulnerabilities are not classified by kind (i.e., they are assigned to the *NVD-CWE-Other* category).

**Systems Histories:** We gather the code history of the projects in which each subject vulnerability appeared. In total, we study vulnerabilities pertaining to 616 different projects. All the analyzed systems are hosted on GitHub [4], a hosting service for Git repositories, which allows us to access the Git history of the projects and study the life cycle of the vulnerabilities. Table 2 provides information on the top-15 projects by number of vulnerabilities. Our appendix reports the complete list of the projects [15].

### 2.3 Terminology

Next, we define the terms used in the following sections to explain our data extraction process and answer our research questions:

- $V = \{v_1, v_2, \dots, v_{1,195}\}$  is a set of NVD vulnerabilities.
- Each vulnerability  $v_i$  occurs in a GitHub project with a git repository  $r_j$ ;
- $C_{r_j} = \{c_1, c_2, \dots, c_m\}$  is the set of commits done on repository  $r_j$ ;
- $IC_{v_i} \subset C_{r_j}$  is the nonempty set of commits (*inducing commits*) that introduces the vulnerability  $v_i$ ;
- In the set  $IC_{v_i} = \{ic_{v_i,1}, ic_{v_i,2}, \dots, ic_{v_i,l}\}$ , we define the last inducing commit ( $ic_{v_i,l}$ ) as the *turning point* ( $tp_{v_i}$ ), i.e., the commit that finally activates the vulnerability  $v_i$ ;
- $fc_{v_i} \in C_{r_j}$  is the *fixing commit* that fixes  $v_i$ ;
- $F = \{f_1, f_2, \dots, f_m\}$  is the nonempty set of files that each commit ( $c$ ) in a repository ( $r$ ) modifies.

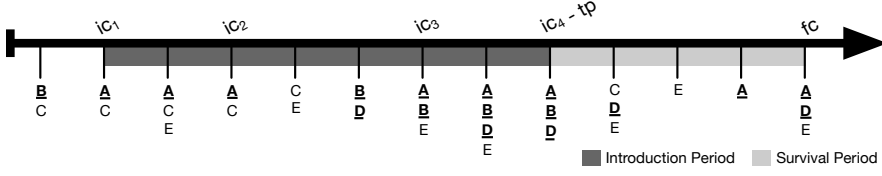


Fig. 1: Example of the life cycle of a vulnerability ( $v_i$ ).  $ic_{v_i}$  = Inducing Commit;  $\{A, B, C, D, E\}$  = files of the project (underlined and bolded if changed by the commit).

## 2.4 Vulnerability Life-cycle Example

We present a short example concerning the life cycle of a vulnerability ( $v_i$ ) to explicate from which commits and files we collect data during our process (see Section 2.5). When filtering vulnerabilities from NVD, we identify vulnerabilities with a corresponding fixing commits  $fc_{v_i}$  (Step 1) and with existing GitHub repositories (Step 2). Figure 1 presents an example of a project's repository ( $r_j$ ), where each stroke in the life-line represents a different commit. In this example, the last commit ( $fc_{v_i}$ ) fixed the  $v_i$ . We use this information as input for SZZ to identify the set of inducing commits (Step 3). As output we have  $C_{r_j} = \{ic_{v_i,1}, ic_{v_i,2}, ic_{v_i,3}, ic_{v_i,4}\}$ .

In this example,  $ic_4$  is the *turning point* ( $tp_{v_i}$ ), because it is the last inducing commit. This commit touched the files  $F = \{A, B, D\}$ . We consider that the vulnerability required file  $A$  to change 1 ( $ic_{v_i,1}$ ), 2 ( $ic_{v_i,2}$ ), 2 ( $ic_{v_i,3}$ ) and 6 ( $ic_{v_i,4}$ ) commits. As such, we compute a minimum of 1 commit, a maximum of 6 commits, an average of 2.75 commits, and a median of 2 commits on file  $A$  for the vulnerability be introduced. We compute these statistics for every  $ic_{v_i,i}$  and every  $f_m$ .

Furthermore, we analyze the survivability of the vulnerability, For this, we compute the number of commits for each  $f_m$  and the days spent until the vulnerability's removal. For instance, in this example, two commits touched  $A$  until the vulnerability was removed.

## 2.5 Data Extraction

Figure 2 presents an overview of our data extraction process. First, we mine vulnerabilities from NVD. Rather than developing our own mining tool, we exploit CVE-Search [2], an open-source tool that imports the entire set of CVEs from the NVD repository into a MongoDB database for easier search and processing. To study the life cycle of these extracted issues, we need a set of vulnerabilities having a well-defined introduction and fixing time. For this reason, we filter the vulnerabilities with fixing commits (patches)  $fc_{v_i} \in C_{r_j}$  (Step 1 in Figure 2). As a result, we have a set  $V$  with 2,973 vulnerabilities

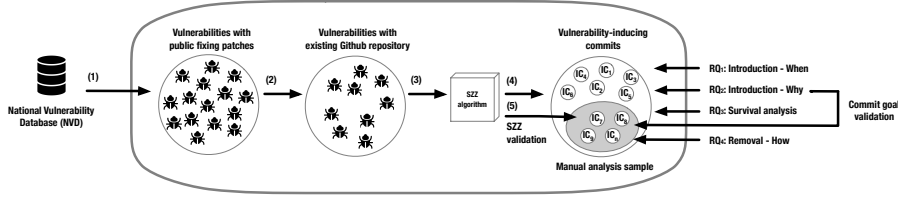


Fig. 2: Process to extract data to answer our research questions.

from projects using different programming languages, such as C and Java. To investigate the life cycle of the vulnerabilities, we mine the repositories of the projects in which they appear. Therefore, for each vulnerability  $v_i$ ,  $r_j$  must still be available on GitHub, once through this platform we have access to  $r_j$ 's history. We remove the vulnerabilities that do not have this property from  $V$  (Step 2 in Figure 2).

Subsequently, for each vulnerability  $v_i$ , we identify the introducing commit  $ic_{v_i}$  (also, *vulnerability-inducing commit*) by using the SZZ algorithm [83] (Step 3 in Figure 2). More specifically, given a vulnerability-fixing commit ( $fc_{v_i}$ ), the algorithm returns the set of commits that last changed the files touched by the commit. In summary, for each  $f_m$  in the input commit, SZZ (1) obtains the diff with respect to the previous commit, (2) retrieves the list of deleted lines, and (3) blames the file to obtain the commits where those lines were changed last. For each  $v_i$ , we run SZZ with  $fc_{v_i}$  as input and obtain the set of commits that likely introduced  $v_i$  (a vulnerability can be introduced in one or more commits). With this approach, we found vulnerability-inducing commits for 1,195 vulnerabilities, of which 651 require more than one commit to appear. We focus on this set of vulnerabilities, while we exclude the remaining 1,778 that do not have complete information.

## 2.6 Data Validation

Recent work has reported low accuracy for the SZZ algorithm [38, 77]; for this reason, we empirically verify the performance of SZZ on our dataset. Two of the authors of this paper (the *inspectors*) independently analyze a statistically significant set of 307 vulnerability-inducing commits (confidence interval=0.95). In doing so, they rely on the source code of both vulnerability-fixing and vulnerability-inducing commits: the task is performed to understand whether the change applied in the vulnerability-inducing version actually introduced the vulnerability that is fixed in the vulnerability-fixing version. After the independent inspection, the two inspectors compared their assessments, finding an agreement in 97% of the cases; they discussed the remaining 3% until an agreement was found. All in all, the precision of the SZZ algorithm was 91% (280 correct vulnerability-inducing commits): this result strongly differs from previous findings [77] and suggests that the performance of SZZ may depend on both the dataset exploited and the type of defects considered, i.e., general



defect vs vulnerability, as well as from the fact that the fixes are generally smaller than those of traditional defects (as further presented and discussed in Section 3), thus allowing SZZ to better identifies the origin of vulnerabilities.

## 2.7 RQ<sub>1</sub>. When - Research Methodology

To understand *when* vulnerabilities are introduced, we do similarly as previous work [87]. Specifically, for each vulnerability-inducing commit, we first compute the number of commits performed on the source code file  $f_m$  where  $v_i$  occurs since the first commit involving  $f_m$  up to  $ic_{v_i}$  (Step 4 in Figure 2). In cases where a vulnerability is associated to multiple vulnerability-inducing commits (i.e., the SZZ algorithm gives more the one commit as responsible for the introduction of  $v_i$ ), we consider the files modified by the last inducing commit to compute minimum, maximum, average, and median number of commits (considering each inducing commit) required by  $v_i$  to gradually affect  $f_m$ . This way, we verify whether vulnerabilities affect source code files since their introduction or are the effect of multiple maintenance and evolution operations.

## 2.8 RQ<sub>2</sub>. Circumstances - Research Methodology

To understand *under which circumstances* developers introduce vulnerabilities, we automatically classify each  $ic_{v_i}$  according to one or more of the categories as described by Tufano et al. [87] (see Table 3). The categories include (i) the commit goals (i.e., the task the developer was performing when introduced a vulnerability), (ii) the project status (i.e., the development phase of the project), and (iii) the developer status (i.e., the characteristics of the developer who introduced the vulnerability). We automatically assign one or more of the categories to each  $ic_{v_i}$  as follows (Step 5 of Figure 2):

**Commit goal.** We rely on a previously proposed approach that analyzes the content of the commit messages and uses keyword-matching [72, 89]. As an example, a commit is classified as ‘*bug fixing*’ if the corresponding commit message contains keywords like ‘bug’, ‘defect’, ‘fix’, ‘repair’, ‘error’, ‘issue’. Our appendix contains the complete list of keywords we use [15]. We also proceed with a manual assessment of the accuracy of this approach in our context. To this aim, we take into account the same set of 307 commits previously used to evaluate the SZZ accuracy (see Section 2.6): two authors of this paper independently analyzed the categories assigned by the keyword-matching approach. The agreement between the inspectors was 94%. After discussion, the accuracy of the textual approach was found to be 86%, confirming the accuracy reported in previous findings [72].

**Project status.** We compute the number of days from the inducing commit’s date to the date of the nearest minor or major release and assigned the proper ‘*working on release*’ category based on this number. A similar procedure is

Table 3: Tags assigned to vulnerabilities inducing commits.

Tags	Description	Values
<b>COMMIT GOAL TAGS</b>		
Bug fixing	The commit aimed at fixing a bug	[true, false]
Enhancement	The commit aimed at implementing an enhancement in the system	[true, false]
New feature	The commit aimed at implementing a new feature in the system	[true, false]
Refactoring	The commit aimed at performing refactoring operations	[true, false]
<b>PROJECT STATUS TAGS</b>		
Working on release	The commit was performed [value] before the issuing of a major release	[one day/week/month, more than one month]
Project startup	The commit was performed [value] after the starting of the project	[one week/month/year, more than one year]
<b>DEVELOPER STATUS TAGS</b>		
Workload	The developer had a [value] workload when the commit has been performed	[low, medium, high]
Tenure	The developer tenure when the commit was performed	[newcomer, medium, expert]

followed for the ‘*project startup*’ category, computing the number of days from the date of the inducing commit to the date in which the project started (the date of the first commit).

**Developer status.** We assign the developer status considering two perspectives. First, we compute the workload of the developers who introduced the vulnerabilities (the authors of the inducing commits): similarly to previous work [87, 87], we use the number of commits made by developers as proxy metric for workload (i.e., the higher the number of commits the higher the workload of a developer in a given time window). Specifically, given a  $ic_{v_i}$  performed by a developer  $d$  during a month  $m$ , we first compute the workload distribution for all developers of the project during  $m$ . Then, we consider the workload of  $d$  to be ‘low’ if his/her number of commits is strictly lower than the first quartile of the distribution, ‘medium’ if between first and third quartile, and ‘high’ when higher than the third quartile. Secondly, we compute the ‘*project tenure*’ category [73, 88], that is, an experience metric that counts the number of months since the developer’s first event on the project where the vulnerability was introduced. Starting from the distribution of all developers’ project tenure, we consider  $d$  a ‘newcomer’ if his/her project tenure is strictly below the first quartile of the distribution and ‘expert’ if strictly higher than the third quartile; otherwise, we assign ‘medium’.

## 2.9 RQ<sub>3</sub>. Survivability - Research Methodology

By knowing the inducing and fixing commits we can analyze the survivability of each vulnerability. In the case of vulnerabilities with more than one inducing commit, we consider their last inducing commits, as they represent the turning point of the introduction process. Furthermore, given a vulnerability  $v_i$ , we

define the time interval between its inducing commit  $ic_{v_i}$  and its fixing patch  $fc_{v_i}$  as a *vulnerable interval* and determines the longevity of  $v_i$ . We can use a vulnerability longevity to compute the number of days between its introduction ( $ic_{v_i}.date$ ) and its fixing patch ( $fc_{v_i}.date$ ), as well as the number of commits between  $ic_{v_i}$  and  $fc_{v_i}$  that modified the file  $f_m$  where the vulnerability occurs (Step 6 in Figure 2). These metrics complement each other: For example, two commits may occur in a short interval of days, but with a high number of commits in-between, therefore analyzing only the date of the commits would be limiting.

After collecting the data, we perform a survival analysis [63], a statistical method that aims at analyzing and modeling the time duration until one or more events happen. The survival function  $S(t) = P\pi(T > t)$  indicates the probability that a vulnerability survives longer than a time  $t$ . The survival function does not increase as  $t$  increases; also, it is assumed that  $S(0) = 1$  at the beginning of the observation period, and, for time  $t \rightarrow \infty$ ,  $S(\infty) \rightarrow 0$ . The survival analysis aims at estimating a survival function from data and assessing the relationship of explanatory variables to survival time. In the context of our study, the population is represented by the vulnerabilities instances while the event of interest is represented by its fix. The ‘time-to-death’ of a vulnerability is represented by the observed time from its introduction to its fix. We refer to this time as the ‘lifetime’ of a vulnerability.

#### 2.10 RQ<sub>4</sub>. Removal - Research Methodology

To understand how vulnerabilities are removed from the source code (Step 7 in Figure 2), we perform a manual analysis of the patches corresponding to the statistically significant sample of 307 vulnerability-inducing commits identified when validating the data extraction process (see Section 2.6).

We followed an open coding process [51]: The 307 vulnerability fixes were equally distributed between two authors of the paper, who independently and manually categorized the kinds of actions performed by developers (e.g., a refactoring) relying on vulnerability-fixing commit message and the Git diff (Step 8 of Figure 2). Afterwards, the authors opened a discussion aimed at clarifying and standardizing the analysis process, which led to the joint re-analysis of all classifications made so far. The final outcome of this process consisted in the definition of a set of categories that explain how vulnerabilities are removed from the code.

#### 2.11 Threats to Validity

We discuss and explain possible threats to the validity of our study.

**Construct validity.** Our results can be affected by the wrong identification of both vulnerability-fixing patches and vulnerability-inducing commits. As for the former, we relied on the vulnerability fixing patches available in

NVD: while this database is curated and improved constantly, we cannot exclude that a patch may not remove the vulnerability as intended. As for the latter, we rely on the SZZ algorithm, which analyzes the previous history of a fixed file to identify the commit that likely introduced the vulnerability. Previous studies have shown that this algorithm may frequently fail in the identification of the correct defect-inducing commit [38, 77]; to account for this aspect and control the reliability of SZZ in our context, we have manually re-assessed its performance on our dataset for a sample of 307 vulnerability-inducing commits, finding an accuracy of 91%. Although the metrics to compute the developers’ workload has been used in the past [73, 87, 87], it is only an approximation, because different commits may require a different amount of work and a developer may work on various projects.

**Internal validity.** To address **RQ<sub>2</sub>**, we implemented an automatic script to automatically characterize vulnerability-inducing commits. To study the ‘*commit goal*’ category, we re-implemented a previously proposed keyword-matching approach [46, 49, 72]: while it has been shown to be accurate [72], we further assessed its accuracy on our dataset by performing a manual analysis of the categories assigned to a statistically significant sample of vulnerability-inducing commits. Such an analysis confirmed previous accuracy, which is beyond 90% yet not perfect; this could affect our reported results. Moreover, the manual analysis performed did not reveal the existence of categories other than those proposed by Tufano et al. [87]. As for the other categories, we re-implemented previously proposed algorithms [73, 87, 88] following the exact description reported in those papers; despite the extensive testing phase, we cannot exclude possible implementation errors. For the sake of verifiability and replicability, we make all data/scripts used in this study publicly available [15].

**Conclusion validity.** While survival analysis (**RQ<sub>3</sub>**) is a standard approach for investigating lifetime expectancy [47, 48], a possible threat concerns the metrics we used to determine the vulnerability survival, i.e., the number of commits and number of days from their introduction to removal. Project activity or developers’ availability may influence the two variables. With respect to **RQ<sub>4</sub>**, we conducted a manual investigation to determine how vulnerabilities are removed: more authors of this paper have been involved in the process in order to increase the accuracy of the classification of vulnerability fixing strategies. Despite the final discussion among the authors, we cannot exclude imprecision and some degree of subjectiveness. Replications of our work are, therefore, still desirable.

**External Validity:** We took into account 1,195 of 62 different types coming from 616 open-source projects. As such, the size of our dataset is comparable with other large-scale software engineering studies (e.g., [54, 87]). However, the life cycle of vulnerabilities affecting other systems (e.g., developed in closed-source settings) may differ. Still in the same category, it is important to point out that our empirical study does not cover all instances of vulnerabilities affecting the subject systems, but rather it is limited to those for which a fixing commit was available. Similarly, we are aware that not all vulnerabil-

ities may have been reported to the CVE (especially if discovered internally by a developer). On the one hand, the large-scale nature of our study makes us confident that similar results would be achieved when considering the vulnerabilities that were missed in our analyses. On the other hand, we would highlight that replications of our study would be desirable as they may provide additional insights into the life-cycle of software vulnerabilities.

### 3 Results

In this section, we present and discuss the results concerning our research questions.

**RQ<sub>1</sub>:** *When are vulnerabilities introduced: during a new file creation or during maintenance/evolution activities on existing files?*

To address our first research question, we exploited SZZ to identify the one or more inducing commits of each of the 1,195 considered vulnerabilities.

**Number of inducing commits per vulnerability.** On total, the algorithm found 3,974 inducing commits, thus, on average, each vulnerability requires 3.32 commits to manifest itself (median=2). Therefore vulnerabilities are mostly introduced in source code after only a few commits, thus not necessarily requiring a large amount of evolutionary activities. Moreover, 651 vulnerabilities (i.e., 55%), have more than one inducing commit. The vulnerability with the highest number of inducing commits (125) is an SQL Injection in the Navigate CMS project (CVE-2018-17552), which was gradually introduced because of the large amount of modifications done by developers to the `login.php` module and, in particular, to the way external users can access the application. It seems reasonable to think that these continuous modifications had the effect of focusing the developers on the functionality and distracting them from checking for anomalous inputs, thus leading to the introduction of vulnerable code.

**The affected files' perspective.** We identified a total of 4,343 files touched by vulnerability-inducing commits. Most of the times vulnerabilities are introduced in commits where developers apply changes on a few files. On average, the number of files modified per inducing commit is 3.63 (median=1), while in 81 cases the inducing commits modified more than five files.

On the one hand, this is somehow in contrast with the findings of Hindle et al. [50], who reported that commits involving more than ten files are more prone to introduce defects. On the other hand, our findings corroborate the hypothesis that vulnerabilities are different from other types of software defects, as also pointed out by Morrison et al. [65].

How many commits touching a certain file are required to introduce a vulnerability into the file? To answer this question, for each file affected by a vulnerability, we compute average and median number of commits between

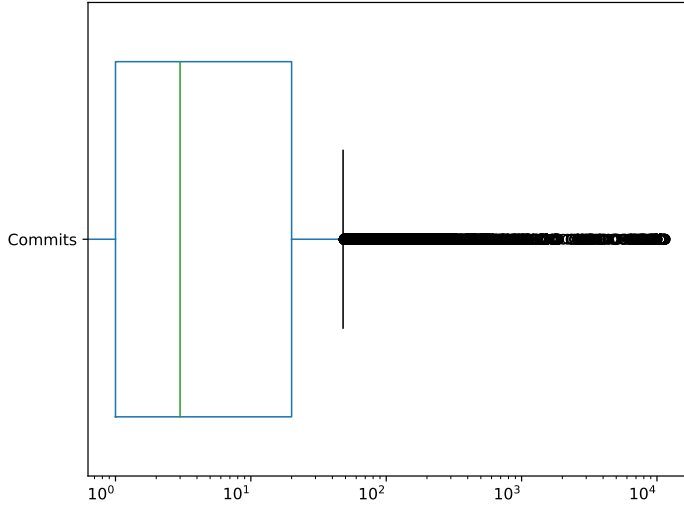


Fig. 3: Number of commits that touched the files up to the inducing commits.

the creation of the file and the first vulnerability inducing commit, as well as between the creation of the file and the turning point. The median number of commits between creation and first *ic* is 16; the median number of commits between creation and turning point is 19.

We found that 1,980 files (46%) were created within a *ic*. Also, we found three cases in which a file was created during a single inducing commit, thus was born vulnerable in the first place. On the other hand of the spectrum, we find files that underwent substantial modifications before they are affected by a vulnerability. A set of 570 files (13%) are changed in more than 500 commits before they start becoming vulnerable. Figure 3 shows the box plot of the number of commits that touched the files up to the inducing commits. We observed the longest case in the WordPress project, where a cross-site scripting vulnerability (CVE-2018-5776) was introduced in the `version.php` file after 11,511 commits on it.

#### Summary for RQ<sub>1</sub>

Source code files become affected by vulnerabilities just after 2 commits, on median, meaning that developers are more prone to introduce vulnerabilities during the initial activities on a newly created file. As such, vulnerability detection strategies should provide a closer support during the first development phases of new files.

Table 4: Categories of vulnerabilities inducing commits.

<b>Commit Goal</b>	New Feature	30%
	Bug Fixing	27%
	Enhancement	24%
	Refactoring	19%
<b>Working on release</b>	One day	67%
	One week	26%
	One month	6%
	More than one month	1%
<b>Project Startup</b>	One week	5%
	One month	1%
	One year	11%
	More than one year	83%
<b>Workload</b>	High	82%
	Medium	14%
<b>Tenure</b>	Low	4%
	Expert	94%
	Medium	5%
	Newcomer	1%

**RQ<sub>2</sub>:** *Under which circumstances are vulnerabilities introduced?*

We investigated the circumstances in which vulnerabilities are introduced in source code, by classifying the vulnerability inducing commits with respect to their goals, project status, and developer status.

Table 4 reports the percentage of vulnerability inducing commits assigned to each category. A commit may have several goals and we analyze the commits messages to match one or more of the following categories: *bug fixing*, *enhancement*, *new feature*, and *refactoring*. Overall, we classified 1,045 vulnerabilities commits with 2,298 commit goal categories. The messages of the remaining 246 inducing commits did not match any of our categories: a missing categorization was due, for instance, to empty commit messages like in the case of a *cross-site scripting* vulnerability (CVE-2017-12474) of the Bento4 project [1].

We found that most inducing commits introduce a *new feature* (30%)—we classified 689 commits in this category. For example, we found an *out-of-bounds* vulnerability (CVE-2013-7456) in the LibGD project [7] with an inducing commit [8] having the following message: “add new interpolation method affine methods scale and rotation”. In this case, we matched *add* in the *new feature* category. Furthermore, 614 (27%) inducing commits were classified as *bug fixing* activities. For instance, an *NULL pointer dereference* vulnerability (CVE-2017-7458) in the ntopng project [12] was introduced by a commit that fixed a security issue. The commit [13] involved two files, added 16 lines of code and deleted 42. Thus, our findings suggest that developers may introduce a security-related problem in the source code during maintenance activities, such as bug fixing, especially when the activity is complex and requires changes to several code entities and/or lines of code: this is in line with previous findings achieved when understanding how maintainability is-

sues are introduced in open-source projects [87,89]. Furthermore, we observed that 550 (24%) of the inducing commits are related to enhancement activities. More surprisingly, *refactoring* operations pertain to the introduction of 445 vulnerabilities (19%): this means that, while trying to improve source code quality, developers often fall into error and introduce security-related problems. On the one hand, this confirms previous findings on the hardness to apply refactoring [31, 68, 90]; on the other hand, our findings indicate that besides maintainability problems [30], the side effect of refactoring includes security issues. This is something that would deserve further investigations.

A vulnerability could match more than one category. In this respect, we categorized 681 vulnerabilities with more than one tag. Overall, 184 vulnerabilities matched all categories. For instance, a *cross-site script* vulnerability (CVE-2017-15278) of the TeamPass project [18] has 20 inducing commits, we found the following messages among them: (i) *“Improved install/upgrade. Continuing code review”* [19]; (ii) *“Fixed an issue with DUOSecurity login. Small improvement on searching on Items page. Added the possibility to search by Tag only. Several small fixes”* [20]; and, (iii) *“Moved import and export features from Home to Items page. - Fixed bug on folders access (list to restraint)”* [21]. We matched the words *upgrade*, *fixed*, *added* and *moved* to *enhancement*, *bug fixing*, *new feature* and *refactoring*, respectively. On average, we matched a vulnerability to 1.78 categories (median=2).

Besides the commit goal tags, we discovered that most of the vulnerabilities are introduced the last days before issuing a release. Indeed, the percentage of vulnerabilities introduced more than one month prior to issuing a release is low (1%). Additionally, most of the vulnerabilities are introduced after more than one year of the project startup (83%), while only 6% are introduced in projects newer than a month. Next, we looked at the developer status in the moment they introduced vulnerabilities. We found that the majority of developers have a high workload when performing vulnerability inducing commits (82%), while 14% and 4% of them have medium and low workloads, respectively. Regarding the tenure of developers, we found that most vulnerabilities are introduced by experts (94%), again confirming findings reported in literature on the source code drawbacks introduced by expert developers [40, 87]. Furthermore, we found that most vulnerabilities (83%) are introduced in advanced phases of the development process, i.e., after more than one year from the project startup, and in most of the cases the day before a deadline (67%). These results may be due to the fact that developers have high workload during these periods, which complies with our findings that most developers (82%) have high workload when performing the vulnerability inducing commits. On the other side, developers who introduced the vulnerabilities are experts, which may be due to the fact that more experienced developers tend to perform more complex and critical tasks [92]. As such, our findings suggest the need for improved methods on how to schedule resources in software projects able to take into account the potential security drawbacks that developers may introduce.



### Summary for RQ<sub>2</sub>

Developers mostly introduce vulnerabilities when introducing new features in the source code, but in a non-negligible number of cases refactoring is the main cause for that. Moreover, the majority of the vulnerabilities are introduced up to a month before a deadline and after the first year from the project startup. Finally, expert developers with high workloads are more prone to introduce vulnerabilities.

Table 5: Descriptive statistics on the number of commits/days before a vulnerability is removed from a system.

Survival	Min.	1st Qu.	Median	Mean	3rd Qu.	Max
Days	0	410.8	1043	1,535.1	2,238.8	7,912
Commits	2	298.2	1,546.5	7,082.8	6,111	351.16

### RQ<sub>3</sub>: What is the survivability of vulnerabilities?

We analyzed the vulnerabilities in terms of the interval between their turning point and their fixing patches. We consider the turning point, because this is when the code becomes fully vulnerable. Figure 4 and Figure 5 show the box plot of the distribution of the number of days and commits, respectively, necessary to fix some of the kinds of vulnerabilities. Further box plots are also available online [15].

Table 5 shows the descriptive statistics of the survival distribution of the number of commits days when aggregating all vulnerability kinds considered in our study. In particular, the median value of such distributions is 1,043 and 1,546 for days and commits, respectively. We found a *double free* vulnerability (high severity) in the code of the *krb5* project (CVE-2017-11462) with 17 inducing commits, the last one occurred on 1995. This vulnerability was fixed in 2017, 7,912 days later. In total, 23 (1.2%) vulnerabilities were introduced and fixed the same day. Moreover, an *out-of-bounds write* vulnerability of the *LibreOffice* [9] project (CVE-2017-7870) took the highest number of commits to be fixed (351,160). In total, 17 (1.1%) vulnerabilities were fixed after only two commits, most of them the same day they were introduced.

For each vulnerability, we analyzed the interval delimited by its last vulnerability inducing commit and its fixing patch (removing commit). Figure 4 shows the box plot of the distribution of the number of days needed to fix some vulnerabilities, the box plots of the remaining vulnerabilities are available online [15]. The box plots, depicted in log-scale, show that most vulnerabilities are fixed after a long period of time (i.e., over 500 days). Figure 5 shows the box plot of the distribution of the number of commits needed to fix some of the kinds of vulnerabilities. The box plots of the remaining kinds are available online [15]. Some kinds of vulnerabilities are removed after a small number of commits, while others are fixed after thousands commits. Aggregating all kinds of vulnerabilities, developers took 1,546 median commits to fix the vulnerabilities.

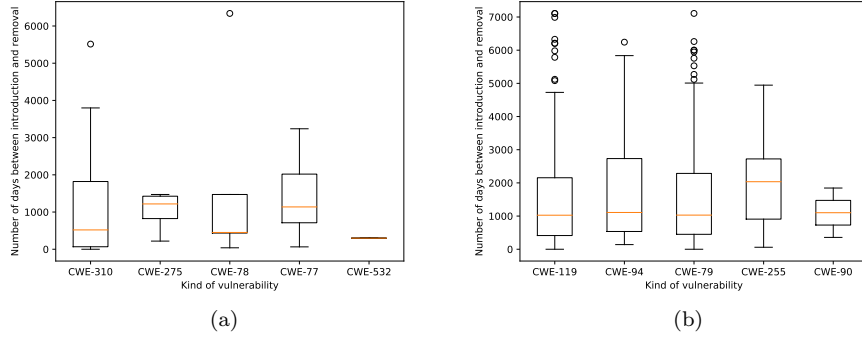


Fig. 4: Survivability of vulnerabilities in terms of days.

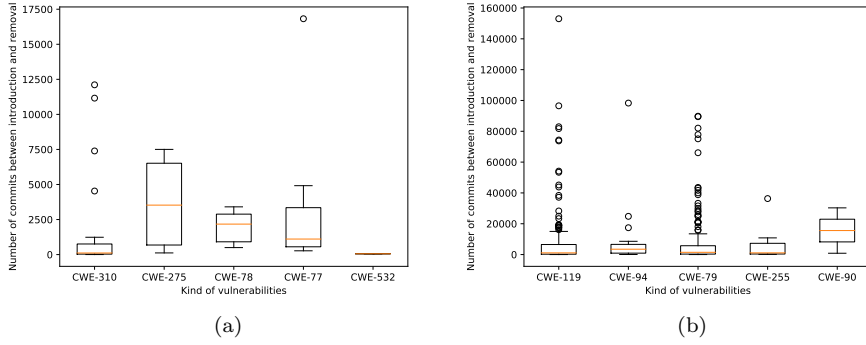


Fig. 5: Survivability of vulnerabilities in terms of commits.

The vulnerability that took more days (7,912) to be fixed was a double free vulnerability on the Kerberos project [6] (CVE-2017-11462). Developers performed 13,415 commits between its introduction and its fix. On the other hand, the vulnerability that took the highest number of commits (351,160) to be fixed was a *out-of-bounds write* in the LibreOffice project [9] (CVE-2017-7870). It took 5,029 days for this vulnerability to be removed. Its turning point was in Mar 2003 and a different developer fixed the vulnerability on Jan 2017. This vulnerability represents the largest difference between the survival days and survival commits of a vulnerability. A *cross-site scripting* of the OAuth2orize project [14] (CVE-2018-11647) took both six days and commits from the last inducing commit to the fixing commit (performed by different developers). Moreover, a *cross-site scripting* vulnerability of the PHPmyadmin project [16] (CVE-2011-3592) took only 148 days to be fixed, even though 20,517 commits were performed on the project before the fix.

Table 6: Methods used to remove vulnerabilities.

CWE	Kind	Removal Method	Total
CWE-79	Improper Neutralization of Input During Web Page Generation (Cross-site Scripting)	Use Escape functions for HTML entities	31%
		Sanitize the text of user inputs	15%
		Manual control of strings to prevent them from containing HTML or PHP elements	1%
CWE-119	Improper Restriction of Operations within the Bounds of a Memory Buffer	Add controls on the buffer size	62%
		Cast between different types of variables containing the buffer size	12%
CWE-264	Permissions, Privileges, and Access Controls	Hide ids within forms	2%
		Check that only authenticated users can access some functions	13%
		Check that authenticated users cannot access sensitive project data	13%
CWE-89	Improper Neutralization of Special Elements used in an SQL Command (SQL Injection)	Control on query parameters to prevent them from containing SQL code	21%
		Queries are parameterized	14%
CWE-352	Cross-Site Request Forgery (CSRF)	Checks to verify if the request is valid	33%
		Added token to identify users	25%

We also found that 591 (39.9%) vulnerabilities were introduced (considering the last inducing commits) and removed by the same developer. Among these vulnerabilities, 36 occurred in the ImageMagick project [5] and a unique developer is responsible for both introducing and removing 34 of them from Sep 2009 until Aug 2018. The same developer introduced 40 of the 45 vulnerabilities belonging to this project and has the highest rate of inserting vulnerabilities in the code. A developer of the TCPdump project [17] project has the highest rate of vulnerability removal (48).

#### Summary for RQ<sub>3</sub>

In general, vulnerabilities survive in code for a large number of days and commits. A significant amount of vulnerabilities (39%) is both introduced and removed by the same developer.

#### RQ<sub>4</sub>: How are vulnerabilities removed?

To address our last research question, we randomly analyzed 307 vulnerabilities of 43 CWE belonging to 213 projects. We analyzed both their fixing commit messages and the lines of code that they changed. Even though we found a large variability in how developers removed each vulnerability (also depending, e.g., on the programming language), in our analysis we could identify and classify recurrent methods adopted when removing five specific kinds of vulnerabilities. Table 6 presents an overview of our results.

Most fixes (57%) involve multiple files. The median number of files involved is 2, with 409 being the fix with the most changed files. Developers added and deleted a median of 11 and 4 lines, respectively. No fix is only made of added lines: At least a line of code is always modified/removed, while three fixes do not add any line in code. The highest number of added and deleted lines was 27,911 and 4,686, respectively.

Considering the results of the manual classification, the five kinds of vulnerabilities that developers removed with recurrent methods are (i) *improper neutralization of input during web page generation* (CWE-79) with 65 occurrences; (ii) *improper restriction of operations within the bounds of a memory buffer* (CWE-119) with 34 occurrences; (iii) *permissions, privileges, and access controls* (CWE-264) with 15 occurrences; (iv) *improper neutralization of special elements used in an SQL command* (CWE-89) with 14 occurrences; and (v) *cross-site request forgery* (CWE-352) with 12 occurrences. These kinds are part of the top-10 most frequent vulnerabilities in our study (see Table 1) and, perhaps more importantly, three of them appear to be in the top-25 most dangerous software errors according to CWE.<sup>1</sup>

Table 6 shows that, most of the time, developers apply a few, specific operations to remove the considered vulnerabilities. For example, the *improper neutralization of input during web page generation* is removed by escaping functions in 31% of the cases, followed by the sanitization of the input text (15%) and by manually checking strings (10%). In general, the fixes often involve simple program transformations: a clear example is represented by the *add control on the buffer size* operation, which is a widely used method to remove out of memory buffers (it is used in 62% of the cases). This program transformation just requires the addition of an `if` statement that restricts the bounds of an input. Similarly, most of the other fixing operations relate to simple actions that could, potentially, be automated and used in combination with vulnerability detection approaches to cover the entire identification/refactoring pipeline, thus substantially supporting developers when dealing with security-related issues.

#### Summary for RQ4

57% of the fixes involve at least 2 files, with developers that add and remove a median of 11 and 4 lines, respectively. The most harmful vulnerability types can be removed through simple program transformations that involve the addition of functions or checks in the code, even when large number of files and lines of code are involved.

<sup>1</sup> <https://cwe.mitre.org/top25/index.html>

## 4 Discussion and Implications

The results of our four research questions provided a number of insights that need to be further discussed as well as several implications for the software engineering research community. Specifically:

**The first phase effect.** According to our findings, vulnerabilities are often introduced shortly after the creation of new files (on median, after 2 commits) and while introducing new features in a software project. On the one hand, our findings contradict the results of previous studies where researchers have showed that software reliability issues, and in particular vulnerabilities, are mainly related to *evolutionary* activities [81, 86], thus stimulating further research on the origin of software issues [78]. On the other hand, our findings have implications for future research on automated solutions that would be aimed at finding vulnerabilities in source code: indeed, we argue that novel techniques able to be *time- and context-dependent*, i.e., analyzing the timeline as well as the intended actions of developers of newly committed source code files, could be devised and experimented.

**More research on vulnerabilities is needed.** Our study constitutes the first, large-scale analysis of the software engineering life cycle of vulnerabilities. Yet, we still have no compelling empirical evidence on the reasons why vulnerabilities remain in a software system for such large periods of time. Similarly, we are still unaware of the motivations why certain actions, e.g., new feature implementation, are more prone to introduce vulnerabilities: for example, this may be due to the lack of proper control mechanisms, a lack of knowledge of basic security principles, or other undocumented reasons. All these aspects represent challenges for the empirical software engineering research community and it is our hope that the results we presented can stimulate a more comprehensive investigation into software vulnerabilities and security practices in general.

**On refactoring effects on software reliability.** As shown in our study, in a non-negligible number of cases developers introduce vulnerabilities when doing refactoring. This is a pretty controversial finding: indeed, this activity is supposed to improve source code quality, while sometimes it leads to the introduction of both software quality and reliability problems [30]. Likely, our results reflect what has been shown in literature, namely that developers do not have proper mechanisms that allow the automation of refactoring activities [89], therefore apply them manually [32, 68]. As such, our findings uncover another potential issue of a wrong application of refactoring, i.e., the introduction of vulnerabilities. It is an important opportunity for the refactoring researchers to further motivate their work and find new applications of study.

**On security testing and more.** Our results indicate that several vulnerabilities are introduced by developers close to deadlines and, perhaps more importantly, remain in a system for a long while. We argue that such results are important for two main aspects. On the one hand, the definition

of advanced testing mechanisms able to assist developers when looking for possible security issues should be devised: this represents a further challenge for the entire testing community and, perhaps, even for automatic test case generation [44, 71, 80]. On the other hand, enabling alternative strategies, e.g., code review practices for vulnerabilities [39], could be an interesting path to allow developers in better spotting errors statically [28, 75].

**Scheduling resources better.** One of the key results of our investigation is that expert developers are those more likely to introduce vulnerabilities. The cause of this phenomenon may be a higher workload or pressure concerning upcoming deadlines. These results call for a more comprehensive overview of software vulnerability research: for instance, novel methods to triage developers' work by optimizing expertise and workload could nicely complement the current research on the topic. Similarly, it is still unknown how more general social issues among developers, e.g., social debt [35, 73, 85], influence the introduction and/or development of security issues.

**Automating vulnerability removal.** Other valuable results from our empirical study concern the study of how software vulnerabilities are removed. Indeed, from the last research question we discovered that most of the actions done by developers to remove vulnerabilities are simple program transformations (e.g., escaping functions) that have the potential to be automated and made available to developers in order to better support the removal of vulnerabilities as well as to implement a comprehensive pipeline covering their entire life-cycle, from discovery to fixing.

## 5 Related Work

Morrison et al. [65] stated that if vulnerabilities behave in the same way as non-security defects in code, then existing defect prevention methods may be used unchanged. Otherwise, understanding how they differ could inform software development process improvement for security. They extended the Orthogonal Defect Classification [36], a scheme software development teams use to collect and analyze defect data to aid software development process improvement, to study process related differences between vulnerabilities and defects, creating ODC + Vulnerabilities (ODC+V). They applied it to classify 583 vulnerabilities and 583 defects across 133 releases of Firefox, phpMyAdmin, and Chrome, and found that they behave differently. Our study complements Morrison et al.'s work [65], and our findings can be used to improve defect preventions methods.

Shahzad et al. [81] conducted an exploratory measurement study to investigate the life-cycle of 46,310 vulnerabilities. They investigated the phases in the life-cycle, the evolution of vulnerabilities during the years, the requirements to exploit them, their functionality and risk level, and the software vendors and products. They consider that the life-cycle of a vulnerability starts when it is discovered by the vendor, a hacker, or any third-party software analyst. It ends when all users of the software install the vulnerability fixing patch.

In this work, we also investigate the life-cycle of vulnerabilities. However, we consider a life-cycle the period between the vulnerability introduction in the source code and its removal from the code. This way, we investigate aspects that may help developers during the development process by understanding how these problems are introduced and removed in practice.

Smith et al. [84] investigated how developers use a security-focused statistic analysis to resolve security defects. They conducted an exploratory think-aloud study with ten developers who had contributed to a security-critical medical records software system written in Java. During the study, they observed each developer as they assessed potential security vulnerabilities identified by Find Security Bugs [3], a static analysis tool to help developers remove security defects early in the development life-cycle. They found that developers ask questions not only about security vulnerabilities, associated attacks, and fixes, but also questions about the software itself, the social ecosystem that built the software, and related resource and tools. Complementing their work [84], in this paper, we study in which circumstances developers introduce them by analyzing the commits' goals and the actions developers took on the code to remove these vulnerabilities. Together with Smith et al. [84] findings, our analysis results can be used to help creating useful tools to detect security problems, and/or to improve existing ones.

Previous work [27, 53, 58, 59] have evaluated security tools that aim at helping developers find and remove vulnerabilities. Mostly, they evaluate the number of vulnerabilities and false positives the tools report. For instance, Austin and Williams [27] compared the effectiveness of four existing techniques for finding vulnerabilities: systematic manual penetration testing, exploratory manual penetration testing, static analysis, and automated penetration testing. Unlike the previous studies, we studied the complete life-cycle of vulnerabilities and provided a large dataset of vulnerabilities data that allowed a better understanding of when and in which circumstances they occur, how long they stay in code and how developers remove them. Our insights and data may be useful to improve the security finding detection scenario.

Murphy-Hill et al. [69] studied refactoring practice at large. They studied data from a variety of data sets spanning more than 39,000 developers, 240,000 tool-assisted refactorings, 2,5000 developers hours, and 12,000 version control commits. Among their findings, they stated that refactorings are indeed performed and the kind of refactoring performed with a tool differs from the kind performed manually. However, they do not consider vulnerabilities in their study. In this study, we analyzed the life-cycle of 1,195 vulnerabilities, including the manual analysis of how they removed 307 vulnerabilities. In future work, our results can be used to propose refactoring tools to help developers when removing this kind of problems.

## 6 Conclusion

We presented a large-scale empirical study conducted over the life-cycle of 1,195 of 62 kinds belonging to 616 open source projects. We aimed at understanding when and under which circumstances vulnerabilities are introduced, what is their survivability, and how they are removed. These results provide several findings including: (i) vulnerabilities are often introduced shortly after the creation of new files; (ii) most of the times vulnerabilities are introduced after the first year of the project startup; (iii) developers introduce vulnerabilities while introducing new features; (iv) expert developers introduce most vulnerabilities, which is probably due their high workload; and, (v) vulnerabilities have high survivability rates regarding both the number of commits and of days they stay in code.

**Acknowledgements** Bacchelli, Palomba, and Braz gratefully acknowledge the support of the Swiss National Science Foundation through the SNF Projects No. PP00P2\_170529 and PZ00P2\_186090.

## References

1. Bento4. <https://github.com/axiomatic-systems/Bento4>.
2. Cve search tool. <https://github.com/cve-search/cve-search>.
3. Find security bugs. <https://find-sec-bugs.github.io/>.
4. Github. <https://github.com/>.
5. Imagemagick. <https://github.com/ImageMagick/ImageMagick>.
6. Kerberos. <https://github.com/krb5/krb5>.
7. Libgd. <https://github.com/libgd/libgd>.
8. Libgd - inducing commit. <https://github.com/libgd/libgd/commit/423c9393917a9d1233fe163a45a0791da306b109>.
9. Libreoffice. <https://github.com/LibreOffice/core>.
10. Mediacycenter. <https://www.pandasecurity.com/mediacycenter/security/consequences-not-applying-patches/>.
11. National vulnerability database. <https://nvd.nist.gov/>.
12. ntopng. <https://github.com/ntop/ntopng>.
13. ntopng - inducing commit. <https://github.com/ntop/ntopng/commit/87212ac1c16e17c7b9de3b53eb77c89969cb6b18>.
14. OAuth2orize. <https://github.com/jaredhanson/oauth2orize-fpm>.
15. Paper vulnerabilities database. [TODO:site](https://todo.site).
16. Phpmyadmin. <https://github.com/phpmyadmin/phpmyadmin>.
17. Tcpdump. <https://github.com/the-tcpdump-group/tcpdump>.
18. Teampass. <https://github.com/nilsteampassnet/TeamPass>.
19. Teampass - inducing commit 1. <https://github.com/nilsteampassnet/TeamPass/commit/5a2030d2612fa411c1f5229f1bc6ea9441f3c3b2>.
20. Teampass - inducing commit 2. <https://github.com/nilsteampassnet/TeamPass/commit/de100bcef62b295a575ea2bcd72abfa544a75085>.
21. Teampass - inducing commit 3. <https://github.com/nilsteampassnet/TeamPass/commit/11b8d37601a37ddd9fd64c06f089fa052d533353>.
22. U.s. nist computer security division. <https://www.nist.gov>.
23. Wannacry ransomware attack. [https://en.wikipedia.org/wiki/WannaCry\\_ransomware\\_attack](https://en.wikipedia.org/wiki/WannaCry_ransomware_attack).
24. O. Alhazmi, Y. Malaiya, and I. Ray. Measuring, analyzing and predicting security vulnerabilities in software systems. *Computers & Security*, 26(3):219–228, 2007.



25. D. Aranha, P. Barbosa, T. Cardoso, C. Araújo, and P. Matias. The return of software vulnerabilities in the brazilian voting machine. *Computers & Security*, 2019.
26. S. Ardi, D. Byers, P. Meland, I. Tondel, and N. Shahmehri. How can the developer benefit from security modeling? In *International Conference on Availability, Reliability and Security*, pages 1017–1025, 2007.
27. A. Austin and L. Williams. One technique is not enough: A comparison of vulnerability discovery techniques. In *International Symposium on Empirical Software Engineering and Measurement*, pages 97–106, 2011.
28. A. Bacchelli and C. Bird. Expectations, outcomes, and challenges of modern code review. In *International Conference on Software Engineering*, pages 712–721, 2013.
29. N. Baddoo and T. Hall. De-motivators for software process improvement: an analysis of practitioners’ views. *Journal of Systems and Software*, 66(1):23–33, 2003.
30. G. Bavota, B. De Carluccio, A. De Lucia, M. Di Penta, R. Oliveto, and O. Strollo. When does a refactoring induce bugs? an empirical study. In *2012 12th International Working Conference on Source Code Analysis and Manipulation*, pages 104–113, 2012.
31. G. Bavota, A. De Lucia, M. Di Penta, R. Oliveto, and F. Palomba. An experimental investigation on the innate relationship between quality and refactoring. *Journal of Systems and Software*, 107:1–14, 2015.
32. G. Bavota, A. De Lucia, A. Marcus, R. Oliveto, and F. Palomba. Supporting extract class refactoring in eclipse: The aries project. In *34th International Conference on Software Engineering*, pages 1419–1422, 2012.
33. G. Canfora and L. Cerulo. Impact analysis by mining software and change request repositories. In *International Software Metrics Symposium*, pages 9–29, 2005.
34. M. Castro, M. Costa, and T. Harris. Securing software by enforcing data-flow integrity. In *Symposium on Operating Systems Design and Implementation*, pages 147–160, 2006.
35. G. Catolino, F. Palomba, D. A Tamburri, A. Serebrenik, and F. Ferrucci. Gender diversity and women in software teams: How do they affect community smells? In *41st International Conference on Software Engineering: Software Engineering in Society*, pages 11–20, 2019.
36. R. Chillarege, I. Bhandari, J. Chaar, M. Halliday, D. Moebus, B. Ray, and M. Wong. Orthogonal defect classification-a concept for in-process measurements. *Transactions on Software Engineering*, 18(11):943–956, 1992.
37. I. Chowdhury and M. Zulkernine. Using complexity, coupling, and cohesion metrics as early indicators of vulnerabilities. *Journal of Systems Architecture*, 57(3):294 – 313, 2011.
38. D. da Costa, S. McIntosh, W. Shang, U. Kulesza, R. Coelho, and A. Hassan. A framework for evaluating the results of the SZZ approach for identifying bug-introducing changes. *Transactions on Software Engineering*, 43(7):641–657, 2017.
39. M. di Biase, M. Bruntink, and A. Bacchelli. A security perspective on code review: The case of chromium. In *International Working Conference on Source Code Analysis and Manipulation*, pages 21–30, 2016.
40. D. Di Nucci, F. Palomba, G. De Rosa, G. Bavota, R. Oliveto, and A. De Lucia. A developer centered bug prediction model. *Transactions on Software Engineering*, 44(1):5–24, 2017.
41. B. Du Bois, S. Demeyer, and J. Verelst. Refactoring-improving coupling and cohesion of existing code. In *Working Conference on Reverse Engineering*, pages 144–151, 2004.
42. M. Finifter, D. Akhawe, and D. Wagner. An empirical study of vulnerability rewards programs. In *USENIX Conference on Security*, pages 273–288, 2013.
43. M. Fowler. *Refactoring: Improving the Design of Existing Code*. Addison-Wesley, 1999.
44. G. Fraser and A. Arcuri. Whole test suite generation. *Transactions on Software Engineering*, 39(2):276–291, 2012.
45. S. Frei, D. Schatzmann, B. Plattner, and B. Trammell. *Modeling the security ecosystem - the dynamics of (In)security*, pages 79–106. Springer US, 2010.
46. Y. Fu, M. Yan, X. Zhang, L. Xu, D. Yang, and J. Kymer. Automated classification of software change messages by semi-supervised latent dirichlet allocation. *Information and Software Technology*, 57:369–377, 2015.
47. P. Heagerty and Y. Zheng. Survival model predictive accuracy and roc curves. *Biometrics*, 61(1):92–105, 2005.

48. R. Henderson, M. Jones, and J. Stare. Accuracy of point predictions in survival analysis. *Statistics in Medicine*, 20(20):3083–3096, 2001.
49. A. Hindle, D. M German, M. Godfrey, and R. Holt. Automatic classification of large changes into maintenance categories. In *International Conference on Program Comprehension*, pages 30–39, 2009.
50. A. Hindle, D. M German, and R. Holt. What do large commits tell us?: a taxonomical study of large commits. In *International Working Conference on Mining Software Repositories*, pages 99–108, 2008.
51. H. Hsieh and S. Shannon. Three approaches to qualitative content analysis. *Qualitative health research*, 15(9):1277–1288, 2005.
52. S. Huang, H. Tang, M. Zhang, and J. Tian. Text clustering on national vulnerability database. In *International Conference on Computer Engineering and Applications*, volume 2, pages 295–299, 2010.
53. N. Jovanovic, C. Kruegel, and E. Kirda. Pixy: a static analysis tool for detecting web application vulnerabilities. In *Symposium on Security and Privacy*, pages 258–263, 2006.
54. Y. Kamei, E. Shihab, B. Adams, A. Hassan, A. Mockus, A. Sinha, and N. Ubayashi. A large-scale empirical study of just-in-time quality assurance. *Transactions on Software Engineering*, 39(6):757–773, 2013.
55. B. Kitchenham. Evaluating software engineering methods and tool part 1: The evaluation context and evaluation methods. *Software Engineering Notes*, 21(1):11–14, 1996.
56. B. A Kitchenham, T. Dyba, and M. Jorgensen. Evidence-based software engineering. In *International Conference on Software Engineering*, pages 273–281, 2004.
57. F. Li and V. Paxson. A large-scale empirical study of security patches. In *Conference on Computer and Communications Security*, pages 2201–2215, 2017.
58. V. Benjamin Livshits and Monica S. Lam. Finding security vulnerabilities in java applications with static analysis. In *Conference on USENIX Security Symposium*, pages 18–18, 2005.
59. M. Martin, B. Livshits, and M. Lam. Finding application errors and security flaws using pql: A program query language. In *Conference on Object-oriented Programming, Systems, Languages, and Applications*, pages 365–383, 2005.
60. G. McGraw. *Software Security: building security in*. Addison-Wesley, 2006.
61. P. Mell, K. Scarfone, and S. Romanosky. Common vulnerability scoring system. *Security & Privacy*, 4(6):85–89, 2006.
62. T. Menzies, A. Butcher, D. Cok, A. Marcus, L. Layman, F. Shull, B. Turhan, and T. Zimmermann. Local versus global lessons for defect prediction and effort estimation. *Transactions on Software Engineering*, 39(6):822–834, 2013.
63. R. Miller. *Survival Analysis*. John Wiley and Sons, 2011.
64. S. Mirhosseini and C. Parnin. Can automated pull requests encourage software developers to upgrade out-of-date dependencies? In *International Conference on Automated Software Engineering*, pages 84–94, 2017.
65. P. Morrison, R. Pandita, X. Xiao, R. Chillarege, and L. Williams. Are vulnerabilities discovered and resolved like other defects? *Empirical Software Engineering*, 23(3):1383–1421, 2018.
66. P. Morrison, B. Smith, and L. Williams. Surveying security practice adherence in software development. In *Hot Topics in Science of Security: Symposium and Bootcamp*, pages 85–94, 2017.
67. R. Muniz, L. Braz, R. Gheyi, W. Andrade, B. Fonseca, and M. Ribeiro. A qualitative analysis of variability weaknesses in configurable systems with #ifdefs. In *International Workshop on Variability Modelling of Software-Intensive Systems*, pages 51–58, 2018.
68. E. Murphy-Hill, C. Parnin, and A. Black. How we refactor, and how we know it. *Transactions on Software Engineering*, 38(1):5–18, 2011.
69. E. Murphy-Hill, C. Parnin, and A. P. Black. How we refactor, and how we know it. *Transactions on Software Engineering*, 38(1):5–18, 2012.
70. S. Neuhaus, T. Zimmermann, C. Holler, and A. Zeller. Predicting vulnerable software components. In *Conference on Computer and Communications Security*, pages 529–540, 2007.

71. F. Palomba, A. Panichella, A. Zaidman, R. Oliveto, and A. De Lucia. Automatic test case generation: What if test code quality matters? In *25th International Symposium on Software Testing and Analysis*, pages 130–141, 2016.
72. F. Palomba, A. Panichella, A. Zaidman, R. Oliveto, and A. De Lucia. The scent of a smell: An extensive comparison between textual and structural smells. *Transactions on Software Engineering*, 44(10):977–1000, 2018.
73. F. Palomba, D. Tamburri, F. Fontana, R. Oliveto, A. Zaidman, and A. Serebrenik. Beyond technical aspects: How do community smells influence the intensity of code smells? *Transactions on Software Engineering*, 2018.
74. L. Pascarella, F. Palomba, and A. Bacchelli. Fine-grained just-in-time defect prediction. *Journal of Systems and Software*, to appear, 2019.
75. L. Pascarella, D. Spadini, F. Palomba, M. Bruntink, and A. Bacchelli. Information needs in contemporary code review. *Human-Computer Interaction*, 2:135, 2018.
76. C. Pfleeger and S. Pfleeger. *Security in computing*. Prentice Hall Professional Technical Reference, 2002.
77. G. Rodríguez-Pérez, G. Robles, and J. González-Barahona. Reproducibility and credibility in empirical software engineering: A case study based on a systematic literature review of the use of the szz algorithm. *Information and Software Technology*, 99:164–176, 2018.
78. G. Rodríguez-Pérez, A. Zaidman, A. Serebrenik, G. Robles, and J. González-Barahona. What if a bug has a different origin?: Making sense of bugs without an explicit bug introducing change. In *12th International Symposium on Empirical Software Engineering and Measurement*, page 52, 2018.
79. N. Sae-Lim, S. Hayashi, and M. Saeki. Context-based code smells prioritization for prefactoring. In *International Conference on Program Comprehension*, pages 1–10, 2016.
80. H. Shahriar and M. Zulkernine. Automatic testing of program security vulnerabilities. In *33rd Annual International Computer Software and Applications Conference*, volume 2, pages 550–555, 2009.
81. M. Shahzad, M. Z. Shafiq, and A. X. Liu. A large scale exploratory analysis of software vulnerability life cycles. In *International Conference on Software Engineering*, pages 771–781, 2012.
82. R. Shirey. Internet security glossary. In *RFC*, 2000.
83. J. Sliwinski, T. Zimmermann, and A. Zeller. When do changes induce fixes? In *International Workshop on Mining Software Repositories*, pages 1–5, 2005.
84. J. Smith, B. Johnson, E. Murphy-Hill, B. Chu, and H. Richter. How developers diagnose potential security vulnerabilities with a static analysis tool. *Transactions on Software Engineering*, XX(X):XX–XX, 2018.
85. D. Tamburri, F. Palomba, A. Serebrenik, and A. Zaidman. Discovering community patterns in open-source: A systematic approach and its evaluation. *Empirical Software Engineering*, 24(3):1369–1417, 2019.
86. M. Tufano, F. Palomba, G. Bavota, M. Di Penta, R. Oliveto, A. De Lucia, and D. Poshyvanyk. An empirical investigation into the nature of test smells. In *2016 31st International Conference on Automated Software Engineering*, pages 4–15, 2016.
87. M. Tufano, F. Palomba, G. Bavota, R. Oliveto, M. Di Penta, A. De Lucia, and D. Poshyvanyk. When and why your code starts to smell bad (and whether the smells go away). *Transactions on Software Engineering*, 43(11):1063–1088, 2017.
88. B. Vasilescu, D. Posnett, B. Ray, M. GJ van den Brand, A. Serebrenik, P. Devanbu, and V. Filkov. Gender and tenure diversity in github teams. In *Conference on Human Factors in Computing Systems*, pages 3789–3798, 2015.
89. C. Vassallo, G. Grano, F. Palomba, H. Gall, and A. Bacchelli. A large-scale empirical exploration on refactoring activities in open source software projects. *Science of Computer Programming*, 180:1–15, 2019.
90. C. Vassallo, F. Palomba, and H. Gall. Continuous refactoring in ci: A preliminary study on the perceived advantages and barriers. In *34th International Conference on Software Maintenance and Evolution*, pages 564–568, 2018.
91. C. Vassallo, S. Panichella, F. Palomba, S. Proksch, A. Zaidman, and H. Gall. Context is king: The developer perspective on the usage of static analysis tools. In *International Conference on Software Analysis, Evolution and Reengineering*, pages 38–49, 2018.

92. A. Zeller. *Why Programs Fail: A Guide to Systematic Debugging*. Morgan Kaufmann Publishers Inc., 2005.
93. S. Zhang, D. Caragea, and X. Ou. An empirical study on using the national vulnerability database to predict software vulnerabilities. In *International Conference on Database and Expert Systems Applications*, pages 217–231, 2011.