

# Data Preparation for Software Vulnerability Prediction: A Systematic Literature Review

Roland Croft<sup>✉</sup>, Yongzheng Xie, and Muhammad Ali Babar

**Abstract**—Software Vulnerability Prediction (SVP) is a data-driven technique for software quality assurance that has recently gained considerable attention in the Software Engineering research community. However, the difficulties of preparing Software Vulnerability (SV) related data is considered as the main barrier to industrial adoption of SVP approaches. Given the increasing, but dispersed, literature on this topic, it is needed and timely to systematically select, review, and synthesize the relevant peer-reviewed papers reporting the existing SV data preparation techniques and challenges. We have carried out a Systematic Literature Review (SLR) of SVP research in order to develop a systematized body of knowledge of the data preparation challenges, solutions, and the needed research. Our review of the 61 relevant papers has enabled us to develop a taxonomy of data preparation for SVP related challenges. We have analyzed the identified challenges and available solutions using the proposed taxonomy. Our analysis of the state of the art has enabled us identify the opportunities for future research. This review also provides a set of recommendations for researchers and practitioners of SVP approaches.

**Index Terms**—Data preparation, data quality, software vulnerability prediction, systematic literature review

## 1 INTRODUCTION

SOFTWARE security is a paramount concern due to the continued increase in cybersecurity attacks and exploits that are affecting organizations [1]. Software security techniques often focus on detecting and preventing Software Vulnerabilities (SVs) that make their way into a software product or deployment pipeline before release [2]. These SVs are a unique class of software defects that introduce security weaknesses to software and allow for malicious use of products [3]. Due to the high importance of removing these security defects, considerable research efforts have been conducted towards their mitigation [4].

Software Vulnerability Prediction (SVP) is a data-driven process for software quality assurance that aims to leverage historical SV knowledge to classify code modules as vulnerable or not. The granularity of the modules can be set as needed, such as file, function, or code snippet. This area of research has recently surged in popularity within the research community [4], due to its importance and value. SVP can ensure software security early in development and solve the incapacibilities of manually assessing large-scale

software systems for potential SVs, which holds inherent value to an organisation.

Like any data-driven process, data preparation serves as one of the most pivotal components for SVP [5]; *garbage in, garbage out*. Consequently, significant efforts need to be expended for data collection and processing [6]. For SVP, we require examples of both vulnerable and non-vulnerable code for training models. Unfortunately, SV data preparation is not a trivial task [7]. High-quality SV data is notoriously difficult to obtain due to its natural infrequency [8], inconsistent reporting [9], and the unwillingness of organisations to make their sensitive data public [10]. It is widely recognized that data noise can severely impact the quality of an SVP model and eventually negatively impact the validity of the research outcomes [11], [12]. That is why datasets are commonly listed as one of the key challenges for this research area [4], [13]. These data quality issues, in combination with the extreme data collection effort requirements, have led many to view data as the major barrier to industrial adoption of SVP [14], [15].

However, despite the importance and difficulties of SV data preparation for both industry and academia, there has been relatively little effort allocated to systematically understand the known challenges of data preparation for SVP models and how to address them. Whilst there are several secondary studies that have analyzed SVP research [4], [13], [16], [17] and acknowledged the existence of problems with the data preparation, these studies have not focused on thoroughly investigating the data preparation related challenges in SVP research.

Motivated by a lack of an integrated and comprehensive body of knowledge on this important topic, we aim to highlight the state of the practice of data preparation for SVP and consequently identify the associated SV data preparation challenges and solutions. This knowledge is expected to assist practitioners and researchers in gaining better understanding

- Roland Croft and Muhammad Ali Babar are with the Centre for Research on Engineering Software Technologies (CREST), School of Computer Science, University of Adelaide, Adelaide 5005, Australia, and also with the Cyber Security Cooperative Research Centre, Joondalup, WA 6027, Australia. E-mail: {roland.croft, ali.babar}@adelaide.edu.au.
- Yongzheng Xie is with the Centre for Research on Engineering Software Technologies (CREST), School of Computer Science, University of Adelaide, Adelaide 5005, Australia. E-mail: yongzheng.xie@adelaide.edu.au.

Manuscript received 6 Sept. 2021; revised 16 Mar. 2022; accepted 17 Apr. 2022. Date of publication 28 Apr. 2022; date of current version 15 Mar. 2023.

This work was supported by the Cyber Security Cooperative Research Centre Limited whose activities are partially funded by the Australian Government's Cooperative Research Centre Programme.

(Corresponding author: Roland Croft.)

Recommended for acceptance by A. M. Moreno.

Digital Object Identifier no. 10.1109/TSE.2022.3171202

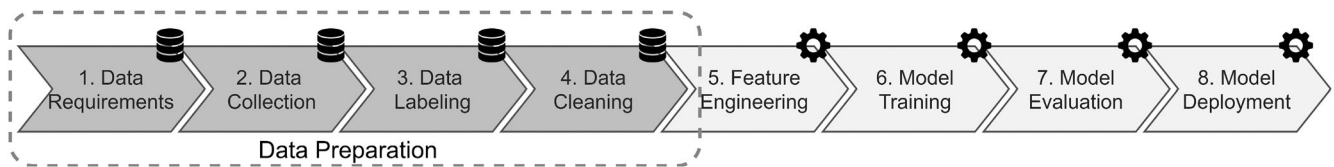


Fig. 1. The machine learning workflow. Adapted from Amershi *et al.* [18].

of the data preparation challenges in SVP and available solutions, in order to support the development and application of more reliable and trustworthy SVP models. In this paper, we empirically examine and synthesize the current practices, challenges and solutions for SVP data preparation through a Systematic Literature Review (SLR). We systematically select 61 peer-reviewed papers on SVP research. We communicate the state of the practice for SVP data preparation techniques, the reported data challenges of these primary studies, and the existing solutions that researchers have used to combat these issues. The main contributions of this research are:

- The first, to the best of our knowledge, systematic review aimed at systematically developing an integrated and comprehensive source of information regarding SVP data preparation practices, challenges and solutions,
- A taxonomy of 16 SV data challenges across six themes. This taxonomy can be used to classify data challenges for future SVP research and practice,
- A mapping of the identified solutions onto the data preparation challenges as per the developed taxonomy,
- A set of recommendations on how to overcome the identified data preparation challenges.

The key contributions of this study are expected to help improve the state of the art and the state of the practice of data preparation for SVP models. The findings can raise awareness and understanding about the important challenges of SV data preparation; such understanding will likely assist to avoid the challenges and improve the reliability of SVP models. Furthermore, we provide recommendations to guide the future research on data preparation and data quality assessment for SVP models. Such future research outcomes are expected to ultimately result in more reliable and trustworthy SVP models. The findings from this study can also be leveraged for enhancing the existing or developing new tools for supporting the construction and application of SVP models in general, and the data preparation for SVP models in particular.

The rest of this paper is organized as follows. Section 2 describes the related work and existing SVP reviews. Section 3 presents the methodology we use to conduct our SLR. The findings of our study are presented in Sections 4, 5 and 6. In Section 7, we provide recommendations for future SV data considerations and research. Finally, in Section 8, we state the applicable threats to validity of our findings, and conclude our study in Section 9.

## 2 BACKGROUND AND RELATED WORK

SVP is a data-driven process that uses learning-based methods to make predictions, and hence follows the standard Machine Learning (ML) workflow.

### 2.1 Data Preparation

Fig. 1 displays the steps involved in a learning-based workflow [18]. For the purposes of SVP and this study, we consider labeling to occur before cleaning.

In our analysis of the reviewed papers, we focus on the data-oriented steps (data collection, data labeling and data cleaning) of the ML workflow. The first step of the ML workflow is the model requirements phase, which identifies the necessary requirements and applications of a model. For our study, we consider this step as *data requirements*, as it is necessary to identify the requirements of the data used to build a model, e.g., what kind of data will be used and from where it will be collected. Hence, these data requirements form a necessary preliminary component of data preparation. We collectively define the first four steps of the ML workflow as *data preparation*.

Practitioners have agreed that the majority of the time taken to construct an ML pipeline is consumed by data preparation [6]. In the 2019 Appen State of AI survey [19], it was reported that a majority of practitioners spend upwards of 25% of their time gathering, cleaning or labeling data. Despite their importance, data preparation processes have rarely been discussed or investigated exclusively [5].

### 2.2 Software Vulnerability Prediction

Software Vulnerability Prediction (SVP) models aim to automatically learn SV knowledge and patterns from historical data. This knowledge can be used to make predictions on the presence of SVs. This process was first notably conceptualized in 2007 by Neuhaus *et al.* [20], and has seen continual technical advancement through research efforts [4].

SVP can be considered as an early form of software security quality assurance, as a trained model can make predictions quickly on static code artefacts, without the need for compilation. In this sense, SVP has been compared to static application security testing methods [21]. Ghaffarian and Shahriari [22] categorized SVP methods into two main approaches: models that do not analyze program syntax and semantics, and models that do. The former utilizes software metrics to describe the code modules of interest, whereas the latter perform directly on source code tokens to perform vulnerable code pattern recognition. Due to the rising popularity of Deep Learning (DL) methods [4], researchers have focused more heavily on approaches that analyze program syntax and semantics, through the use of text-based, sequence-based or graph-based source code feature representations [13].

Data preparation for SVP follows the standard workflow for ML data preparation. SV labels are assigned to the extracted code modules to obtain a labeled SV dataset [10]. The process is heavily dependent on the data sources selected for the codebase and SV labels. Fig. 2 displays the SVP data preparation pipeline.

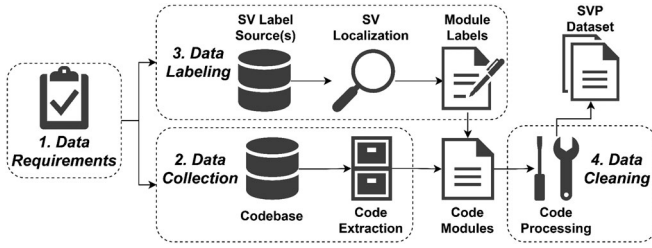


Fig. 2. The SVP data preparation pipeline.

### 2.3 Existing SVP Reviews

With the increasing popularity of data-driven approaches for software vulnerability analysis and discovery, several researchers have reviewed the published SVP approaches and techniques. We briefly describe the key focus areas of the relevant review studies below.

Three papers by different groups of researchers, Li and Shao [23], Coulter *et al.* [10], and Ghaffarian and Shahriari [22], reported separate reviews of the literature on the use of machine learning and data mining for software vulnerability discovery and analysis. Coulter *et al.* [10] provided a more general framework for data-driven cybersecurity tasks, including SVP, whereas Li and Shao [23] and Ghaffarian and Shahriari [22] focused on the specific features and approaches. Le *et al.* [24] also conducted a survey of data-driven methods for SV assessment and prioritization, but they did not consider SV discovery. With the success of DL in fields such as image processing, speech recognition, and natural language processing, researchers have been increasingly motivated to apply DL for the SVP domain. Lin *et al.* [13], Singh and Chaturvedi [25], and Zeng *et al.* [16] all conducted an analysis of the deep learning techniques used by researchers for SVP.

To the best of our knowledge, only two studies have been published that focus on *systematically* reviewing SVP research and knowledge: Semasaba *et al.* [17], and Hasif *et al.*

[4]. The former exclusively investigated Deep Learning techniques, whereas the latter provided a wider view of SV detection, including non-learning based techniques. Similar to the previous secondary studies, the analysis of these systematic reviews focused on the models, techniques and features. Furthermore, all existing secondary studies for SVP focused on the model-oriented steps of the ML workflow, particularly features and techniques (steps 5-6 of Fig. 1).

To this extent, there has been little focus on the data-oriented processes. The SV data used to train a model is the most imperative component of this data-driven process. Although most studies have reported data preparation and data quality as significant issues for this research area [4], [10], [13], [16], [17], [23], they have not performed in-depth analysis of the data quality in SVP research to determine the encountered issues or potential solutions. This knowledge gap fails to provide practitioners and researchers with the specific insights needed to remediate data quality issues. It is vital to gain a better understanding of the quality of data utilized for SVP research; such comprehension is also expected to improve our abilities to better understand how well the SVP approaches work in practice.

Hence, an effort like ours can be of great importance as it not only highlights the critical research gap, but also

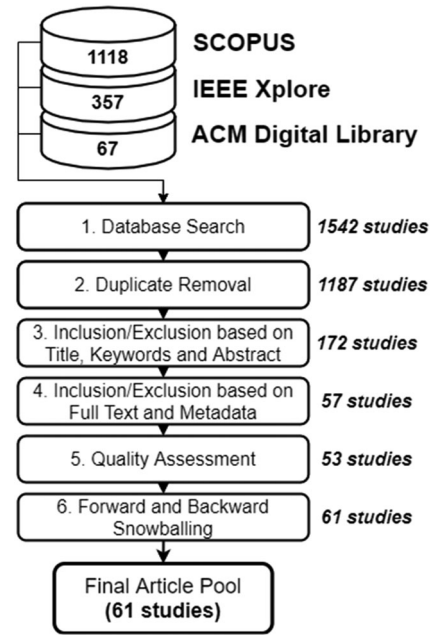


Fig. 3. Stages of the study selection process.

contributes to the evidence-based body of knowledge of data preparation for SVP models. Whilst existing reviews have yielded important insights into this research domain, our systematic review has been motivated by several unique research questions whose answers have enabled us to provide novel findings and potentially useful insights. The knowledge produced by our SLR can be an important complementary piece to the existing secondary studies for providing a consolidated picture of the published literature on different components of the SVP pipeline.

## 3 RESEARCH METHODOLOGY

To obtain insights into the SVP data preparation processes, challenges and solutions, we conducted an SLR of SVP literature. Our findings will potentially be useful to both researchers and practitioners for providing guidance for future SV data preparation and in assessing the validity of existing SV datasets.

To conduct this SLR, we followed the methodological guidance provided by Kitchenham *et al.* [26] and Zhang *et al.* [27] to ensure that our assessment of the existing literature was unbiased and repeatable. The research method was conducted in close collaboration by the first two authors, with guidance from the third author.

Fig. 3 presents the complete search and study selection workflow, and the number of retrieved papers at each stage. The search and study selection process was conducted in February 2021. We obtained a total of 61 studies from our study selection process, which are presented in the Appendix, which can be found on the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TSE.2022.3171202>.

To guide our analysis, we aimed to address the three Research Questions (RQs) presented in Table 1.

### 3.1 Search Strategy

We began with a search strategy to extract all potentially relevant research papers from academic digital libraries.

TABLE 1  
Research Questions Addressed in This Study

Research Question (RQ)	Motivation
<b>RQ1.</b> What considerations do researchers make for each of the data preparation processes when constructing SVP datasets?	We investigate the SVP data preparation practices and choices reported by researchers to help inform the state of the practice and reasoning. The embodiment of this knowledge helps provide workflow guidance for researchers and practitioners, and assists them to identify the important decisions and reasoning when selecting or building an SV dataset.
<b>RQ2.</b> What are the considered challenges and issues for SV data preparation and datasets?	We aim to analyze the reported data challenges to provide an overview of the issues that researchers face when performing SV data preparation or using SV datasets. The results of this RQ will help guide the decisions made by researchers when selecting their data preparation steps, and inform practitioners of the current issues and limitations of the reported empirical SV analysis and prediction.
<b>RQ3.</b> How do researchers address dataset issues and preparation challenges?	This RQ builds upon the findings of RQ2 to analyze the remediation techniques that researchers have used to help overcome the aforementioned challenges. Hence we not only provide a categorization of data challenges for SV-related research, but we also map the solutions that researchers have used to address these issues. These findings can help researchers and practitioners overcome data challenges in the future.

TABLE 2  
Formulation of the Search String

Category	Subject	Search Terms
<i>Population</i>	Software	"software" OR "code"
<i>Intervention</i>	Machine Learning Static Application	"learn" OR "neural network" OR "artificial intelligence" OR "AI-based" OR "predict" NOT ("fuzz" OR "test" OR "attack" OR "adversarial" OR "malware" OR "description")
<i>Comparison</i>	-	-
<i>Outcomes</i>	Software Vulnerability Prediction	"vulnerability" AND ("predict" OR "detect" OR "classify" OR "identify" OR "discover" OR "uncover" OR "locate")

To design the search string, we utilized the PICO (Population, Intervention, Comparison, Outcomes) framework [28]. The *Comparison* component was not applicable to our review because our goal was not to conduct comparison of software with different interventions. Table 2 presents the key terms for each PICO component; we formed our search strings through the union (AND) of the PICO components. We altered the search string suitably to match the differences in the search capabilities of each database. When applicable, we matched the relevant keywords in the title, abstract and keywords of the papers, except for exclusion keywords (prefaced with NOT) which were only matched in the title. Wildcard matching was performed to capture different word variants when available; otherwise we defined the term variants manually, e.g., *predict* and *prediction*. When available, we applied additional search filters to match the exclusion criteria defined in Section 3.2.1 (i.e., limiting to English articles or research papers.) Table 2 only defines the base strings. To find these strings, we consulted papers included in the previous reviews. The full search strings are available in our online appendix, available in the online supplemental material.<sup>1</sup>

We applied this search string to the two most frequently used academic digital libraries for software engineering, as identified by Zhang *et al.* [13]: IEEE Xplore, and ACM digital library. We then additionally included SCOPUS as it is the largest academic literature database available [25], which indexes several other smaller academic databases.

We did not use other search engines such as Google Scholar due to the amount of noise in the search results and need for subjective stopping conditions. We initially retrieved 1542 studies: 1118 studies from SCOPUS, 357 studies from IEEE Xplore, and 67 studies from ACM Digital Library. We downloaded all retrieved studies and then manually removed duplicates, which reduced the total number of studies to 1187.

### 3.2 Study Selection

We sought to select any paper on the topic of Software Vulnerability Prediction (SVP), which we have defined as any model utilising supervised learning-based techniques (ML or DL) for prediction, detection or discovery of an SV in a static code module. We included any study that contributes an SVP model, process or evaluation based on our definition of SVP.

Our definition of SVP hinges on three major principles: *learning-based*, *vulnerability discovery*, and *static code artefacts*. First, we have defined learning-based as the use of a supervised ML or DL algorithm that can learn from training data to make predictions on a dataset [31]. To this extent, we did not include anomaly detection, unsupervised methods, or studies that focus on pure statistical or correlation analysis. Second, the study must have utilized a model that aims to *discover* unknown SVs within a code artefact. This excluded methods that used code clones or similarity detection, as these methods are only able to detect a pre-defined set of SVs and are unable to discover new types. We also did not consider malicious code as SVs. Third, we only included studies that used static code artefacts to make predictions;

<sup>1</sup> [https://github.com/RolandCroit/SVP\\_Data\\_SLR\\_Appendix/blob/main/Search\\_Strings.md](https://github.com/RolandCroit/SVP_Data_SLR_Appendix/blob/main/Search_Strings.md)

TABLE 3  
The Inclusion/Exclusion Criteria

Inclusion Criteria
I1. The study relates to the field of SVP, and informs the practice of Software Engineering.
I2. The study presents a unique SVP process or evaluation.
I3. The study is a full paper longer than six pages.
Exclusion Criteria
E1. Solely a literature review or survey article.
E2. Non peer-reviewed academic literature.
E3. Academic articles other than conference or journal papers, such as book chapters or dissertations.
E4. Studies not written in English.
E5. Studies whose full-text is unavailable.
E6. Studies published to a venue unrelated to the discipline of Computer Science.
E7. Studies that are published to a journal or conference with a CORE ranking of less than A and H-index less than 40, and that have a citation count of less than 20.

either source code, code binaries or an intermediary representation. Hence, we excluded any study that requires run-time analysis of the code (e.g., dynamic testing or attack detection).

### 3.2.1 Inclusion/Exclusion Criteria

The inclusion/exclusion criteria we adopted, displayed in Table 3, were inspired by similar studies [29], [30].

To ensure that we obtained a set of high-quality papers, we adopted a venue assessment approach (E7) used by Sabir *et al.* [31]. We removed the studies published in low quality venues: venue ranking below A using the CORE ranking system,<sup>2,3</sup> and h-index below 40 as recorded in the Scimago database.<sup>4</sup> However, the original influential papers of this domain may have been published in low quality venues. Hence, we only excluded a paper based on venue if it had also not been cited frequently (<20 citations). The citation count is obtained through Google Scholar.<sup>5</sup> The first two authors collaboratively determined suitable thresholds for this criterion through an initial pilot study of 100 papers, to confirm that any papers excluded through this criterion were indeed of lower quality.

We first excluded 1015 studies using information extracted from the title, abstract and keywords. We then excluded an additional 115 studies after processing the full text and metadata (i.e., venue, article type, citations) to obtain a set of 57 studies.

### 3.2.2 Quality Assessment

For SLRs, it is vital to assess the quality of primary studies to ensure that we form a proper and fair representation of the research works [26]. We conducted the assessment process using a quality checklist, and excluded any study that did not pass the checklist. We adopted the quality checklist

TABLE 4  
The Quality Checklist

Data Criteria
DC1. The data source must be reported. If a publicly available dataset is used, the name must be reported.
DC2. A description of the data, such as its size, programming language and class distribution, must be provided.
DC3. The process in which the independent variables are extracted from the data as input to the model must be clearly stated.
DC4. The method in which the data is labeled as vulnerable and non-vulnerable must be clearly stated.
Prediction Model Details Criteria
MC1. The output of the model must be clearly defined.
MC2. The granularity of the dependent variable(s) must be reported.
MC3. The machine learning method and approach must be clearly reported.
Evaluation Criteria
EC1. The performance measure of the model must be reported.
EC2. The predictive performance values must be clearly presented in terms of raw performance numbers, means or medians.

defined by Hall *et al.* [32], and refined by Hosseini *et al.* [29] in their SLR of defect prediction models, as the defect prediction process shares similarities with SVP. This resulted in three stages of assessment: the data, the prediction model details, and the evaluation criteria, displayed in Table 4. Although our study only considers data preparation for analysis, we assessed all three criteria to determine the overall quality of the paper. We removed a total of four studies that did not pass the quality assessment criteria.

### 3.2.3 Snowballing

It is expected that an initial automated search strategy will be unable to identify all relevant studies, as the search string cannot identify obscurely phrased studies, and the digital libraries selected do not exhaustively include all peer-reviewed literature [33]. Hence, after we conducted initial study selection, we utilized manual search processes, both forward and backward snowballing, to obtain additional relevant studies that were not contained in our selected digital libraries or identified by our automatic search. Forward and backward snowballing identify additional relevant studies from papers that cite or are included in the reference lists of the set of included studies, respectively [33]. These identified papers were similarly assessed using the inclusion/exclusion criteria and the quality assessment criteria. We included an additional eight papers in the final set through the snowballing process.

Our final article pool contained 61 studies; 53 studies which passed the initial selection process and eight additional snowballed papers. The studies are listed in the Appendix, available in the online supplemental material.

## 3.3 Overview of the Primary Studies

Fig. 4 displays the number of selected SVP papers over the years. We have not reported values for 2021 as this data is

2. <http://portal.core.edu.au/conf-ranks/>

3. <http://portal.core.edu.au/jnl-ranks/>

4. <https://www.scimagojr.com/journalrank.php>

5. <https://scholar.google.com.au/>

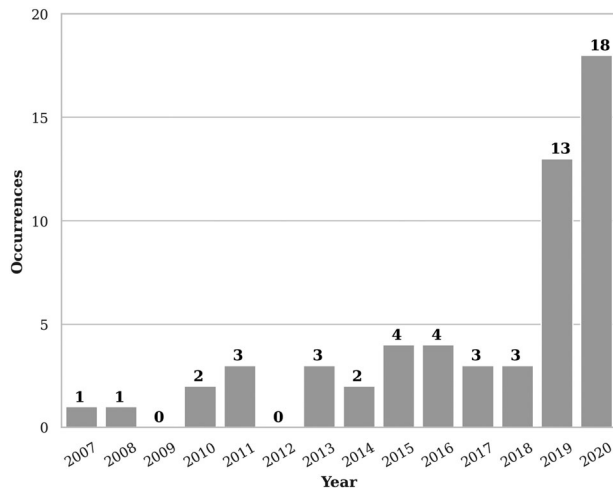


Fig. 4. The number of selected primary studies by year.

incomplete. We observed that this area of research has received exponential popularity within the last two years. This indicates that this is an area undergoing huge growth at the time of this study. Hence, our review contributes important and timely value to this emerging area by synthesizing the current knowledge of the underlying data preparation processes for these empirical studies.

### 3.4 Data Analysis

#### 3.4.1 Data Extraction

We used a data extraction form and the data extraction processes outlined by Garousi and Felderer [34] and Kitchenham *et al.* [26]. Our data extraction form, provided in our online appendix, available in the online supplemental material,<sup>6</sup> consisted of 50 fields describing all data related steps reported by the authors and the details of their dataset(s). These fields consisted of five checkbox questions, 29 multiple choice questions, nine short answer questions, and seven long answer questions. Thirty seven of the 50 fields pertained to the first RQ, and the other 13 collectively related to the latter two RQs. The data extraction form was completed collaboratively using Google Sheets.

An initial pilot study of 10 papers was conducted collaboratively by the first two authors to help design the form and ensure author agreement [34]. The first two authors then performed data extraction individually; the paper set was divided in half randomly for each author to complete. After this process was completed, each author reviewed the data extraction outputs of the other author to ensure consistency. Disagreements were resolved through discussion.

#### 3.4.2 Data Synthesis

The aim of an SLR is to aggregate information from primary studies [26]. For RQ1, we qualitatively examined the outputs of our data extraction form to identify and report the major factors relating to each of the four data preparation steps.

For RQ2 and RQ3, we used thematic analysis to synthesize the data [35]. Specifically, this process was used to identify the reported data challenges and solutions. Any discussion in

a paper that explicitly had mentioned a challenge pertaining to the data, resolved or unresolved, was coded. To ensure that this qualitative coding was grounded by the data, and not affected by any biases of the data extraction form, we imported the full papers into Nvivo [36], a qualitative data analysis tool, and performed coding on the papers directly. We followed the steps for thematic analysis developed by Braun and Clarke [35]:

- 1) *Familiarizing with data*: The initial familiarization was done through the data extraction phase (Section 3.4.1) in which the first two authors read each full paper and filled the data extraction form. This familiarized the first two authors with the relevant factors relating to SV data that were discussed in the papers.
- 2) *Generating initial codes*: To generate initial codes, we used *open coding* of the relevant text in the primary studies using Nvivo. The data was broken down into smaller components and labeled using a code [35], where a code is a word or phrase that acts as a label for a selection of meaningful text in the paper. This process was completed iteratively, with the initial codes being revised and merged in later rounds. Each primary study was usually allocated to more than one code or theme, as each paper can discuss multiple SV data challenges and coding was done on small individual components of the papers.
- 3) *Searching for themes*: We reviewed all the codes and sorted them into themes. As data challenges revolve around data quality, we used existing data quality dimensions [37], [38] to identify potential groupings that the codes might fall under.
- 4) *Reviewing themes, defining and naming themes*: This process involved reviewing all the codes and themes, and revising their allocations.
- 5) *Producing the report*: We present the findings of our thematic analysis in Sections 5 and 6.

## 4 SV DATA PREPARATION CONSIDERATIONS (RQ1)

We first provide an overview of the considerations that researchers have made when performing SV data preparation processes, which we have identified qualitatively through our data extraction process. This documentation of the considerations helps to inform practitioners and researchers of the state of the practice. Furthermore, it can assist these users to better understand how to construct an SVP dataset, and the important aspects to scrutinize. Fig. 5 displays the main decisions that need to be made for each data preparation step.

### 4.1 Data Requirements

In the data requirements phase, the requirements for the data to achieve the desired model context and capabilities are specified. There are four main components of the data requirements that need to be specified for SVP:

**Programming Language(s).** Researchers are motivated to explore different approaches to mitigate security risks for different languages. As seen in Fig. 6, C/C++, PHP and Java have been the most commonly investigated languages among the primary studies.

6. [https://github.com/RolandCroit/SVP\\_Data\\_SLR\\_Appendix/blob/main/Data\\_Extraction\\_Form.pdf](https://github.com/RolandCroit/SVP_Data_SLR_Appendix/blob/main/Data_Extraction_Form.pdf)



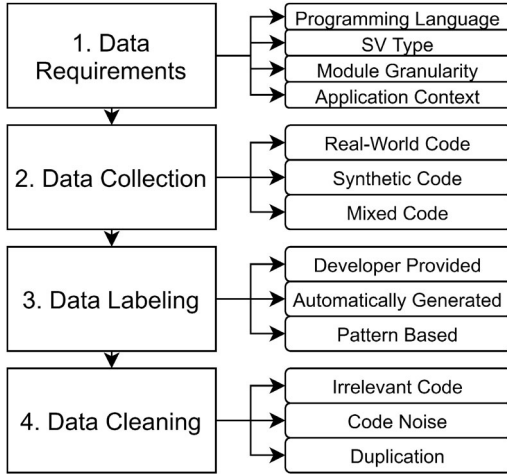


Fig. 5. SVP data preparation step considerations.

C/C++ has been most frequently chosen for analysis as it has a lower level of abstraction and is commonly used to build security critical applications [P4, P16, P54, P58]. PHP has been commonly used for programming web applications, which are highly susceptible to vulnerabilities and exploits [P1, P6, P25, P41, P50, P55], and hence researchers have aimed to ensure its security. Java has also been commonly chosen as it is overall one of the most popular programming languages [P39, P40, P46].

**Vulnerability Type(s).** Similar to the previous consideration, researchers may target detection capabilities towards certain SV types. However, we observed that the majority of methods are capable of detecting a variety of SVs. Over 45% of studies (28 out of 61) did not even report the types of SVs present in the data, and researchers were often limited to the SV types present in their SV label source. However, some studies chose to restrict their analysis to more critical SVs of interest. For example, Fidalgo *et al.* [P6] and Shar and Tan [P25] focused their analysis on SQL injection and cross-site scripting (XSS) as these are common critical web application vulnerability types. Wang *et al.* [P29] and Ghaffarian and Shahriari [P39] limited their studies to just the CWE Top-25 vulnerabilities.<sup>7</sup> Saccente *et al.* [P34] identified that a model trained to predict any SV type produces unreliable predictions in comparison to a model trained to predict just one type. However, the impacts of this data requirement decision on model efficacy are otherwise largely under-explored.

**Code Module Granularity.** The granularity of vulnerability detection has a significant impact on a model and data collection. Depending on the granularity of the inputs used for an SVP model, it can either be used to direct testing efforts by predicting which large scale components are potentially at risk [15], or to explicitly detect fine-grained components that contain vulnerabilities [23]. In our set of primary studies, we identified six levels of granularity (in descending order of granularity): 1) component level, 2) file/class level, 3) function/method level, 4) program slice level, 5) statement level, 6) commit level. In Fig. 7, other than the component-level, the larger granularities are shown to be of more popularity. The file level has been considered as the standard granularity

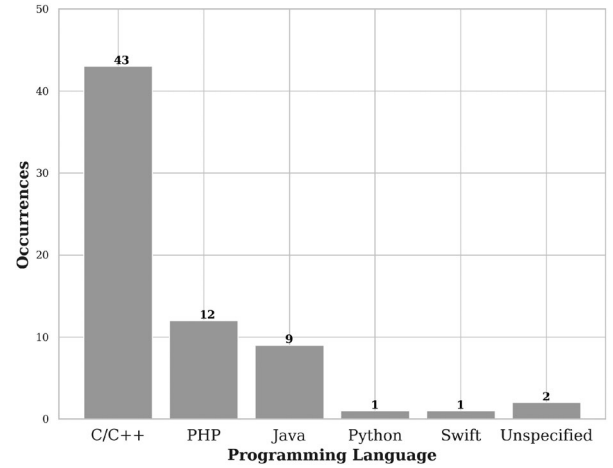


Fig. 6. The number of primary studies for each programming language.

for SVP research [P12, P36, P46] and has been confirmed as actionable by developers [P5]. However, researchers have recently begun to favor finer granularities as they better enclose the scope of vulnerable code snippets, and are more easily inspected [P2, P3, P10, P15, P31, P32].

**Application Contexts.** There are three main vulnerability detection application contexts: within-project, cross-project, and mixed-project, in which each have different requirements. For within-project prediction, both the training data and the testing data come from the same project. In contrast, for cross-project prediction, there is an assumption that there is an insufficient amount of training data available in the target project. Therefore, labeled data from other source projects are used for training. More than one dataset can be collected to compensate for the inadequacy of labeled data in the target project. Mixed-project prediction is a special case of the cross-project setting. The labeled data from multiple projects are combined together to produce sufficient data for both model training and testing. This differs from the cross-project setting in which data from other projects only comprise the training dataset.

Within-project prediction has proven to be the more popular prediction context, with 52% (32 out of 61) of the primary

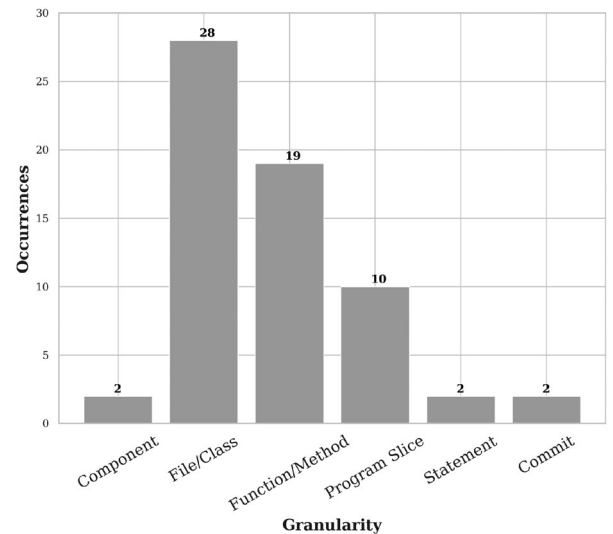


Fig. 7. The number of primary studies for each level of granularity.

<sup>7</sup> [https://cwe.mitre.org/top25/archive/2021/2021\\_cwe\\_top25.html](https://cwe.mitre.org/top25/archive/2021/2021_cwe_top25.html)

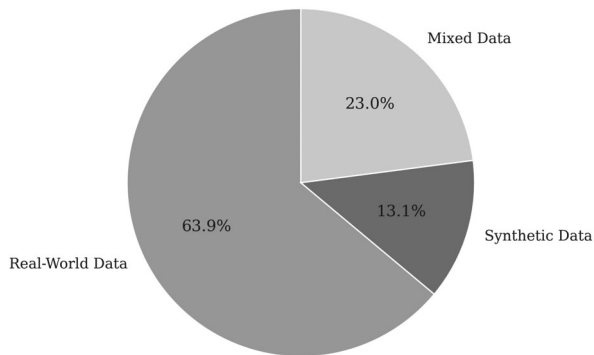


Fig. 8. The proportion of primary studies for each data source type.

studies forming their datasets from a singular project. This is because researchers have considered within-project prediction to be the standard use case of SVP [P5, P36]. Furthermore, cross-project prediction has often performed poorly due to differences in data distribution of the source and target project(s) [P1, P12, P26, P39, P56]. Only 18 of the primary studies considered cross-project prediction, 13 of which also considered within-project prediction. However, due to the various data issues that we will discuss in Section 5, several researchers have considered the reuse of existing datasets from other projects to be a necessity [P1, P10, P18, P21, P22, P26, P31, P37, P41, P43, P56]. Twenty three of the primary studies (38%) performed mixed-project prediction.

## 4.2 Data Collection

SV datasets can be categorized into three main areas based on the type of data sources used to generate the code modules: real-world data, synthetic data, and mixed data. The type of data is the main influence on how the data is collected.

**Real-World Data.** Both the code and the corresponding vulnerability annotations have been derived from real-world repositories. The code has been typically collected for projects hosted on repository hosting sites, such as GitHub [39], or through a different version control system. Experiments conducted using this data are usually considered to be better representative of industrial application because of the reflection of the complexities of the real-world vulnerabilities [40].

**Synthetic Data.** The vulnerable code examples and the labels have been artificially created. The examples from these data sources are synthesized using known vulnerable patterns. Synthetic datasets include SARD [41], OWASP Benchmark [42], and SQLI-Labs [43]. These datasets were originally used for evaluating traditional static and dynamic analysis based vulnerability prediction tools, due to their large test suite size and noise-free information.

**Mixed Data.** Several researchers have opted to create datasets by merging both real-world and synthetic data sources [P32, P44, P48]. This is typically done to achieve a sufficient dataset size whilst maintaining a certain level of real-world representation. Mixed datasets have been primarily constructed for DL-based studies [P32, P44], which are particularly data hungry in comparison to ML.

Fig. 8 displays the proportion of primary studies that use each data type. Real-world datasets have been considered the *de facto* type of data sources in this domain by researchers, primarily as they better represent real-world scenarios

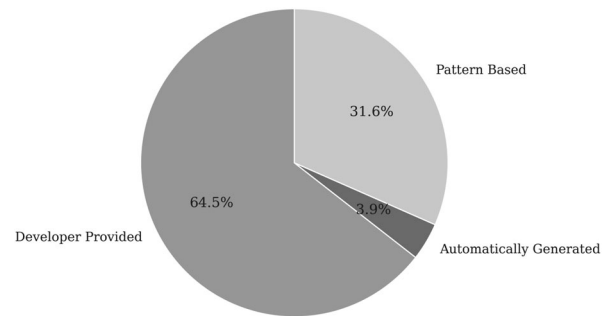


Fig. 9. The ratio of SV data label sources. A study may use more than one label source.

than synthetic examples [40]. The representativeness of a dataset is an important consideration as it helps improve the generalizability and validity of findings [44]. Studies have further claimed that the code patterns of synthetic test cases follow a similar coding format, failing to reflect the characteristics of code patterns in production environments [P2, P3, P9, P29, P31, P39, P42, P43, P61].

However, there are two primary positive traits of synthetic datasets. First, there are a much larger number of labeled synthetic samples that are able to be created in comparison to real-world examples [P8, P23, P34, P38, P39, P44, P54]. SVP is a data-hungry process that requires a sufficient amount of training data [4]. The existence of their labels also significantly reduces the effort of data preparation [P3, P6, P10, P39, P40]. Second, as the code samples are generated with their labels, the labels are cleaner and more reliable than data extracted from noisy real-world repositories [P39, P40], for reasons discussed in the following section.

## 4.3 Data Labeling

For SV data, labels categorize whether a code module is vulnerable or not. Data labeling involves extracting data labels using an external source or tool, to assign to the collected code modules. We observed three SV label sources, that align with the findings of Chakraborty *et al.* [40]: developer-provided, automatically generated, and pattern-based. Fig. 9 displays the ratio of these three SV label types. The choice of labeling approach is often dependent on the data source type collected, as outlined in Section 4.2. Hence, the relative frequency of each method is similar to that of Fig. 8.

Fig. 10 displays the labeling process using developer provided labels. These labels have been extracted from SVs that were identified by developers and reported in security advisories or issue tracking systems, such as NVD [45], Jira [46] or Bugzilla [47]. Whilst these information sources are usually accurate, it is not the same as developers hand-labeling code modules, as patches do not directly equate to labels. Researchers have expended significant effort to trace the label source to the code modules [P2, P18, P29, P43, P61]. This involves identifying the relevant vulnerability reports, extracting code fixes, and localizing these fixes to the relevant granularity and code modules, with each step introducing potential noise to the labels. Labeling of the non-vulnerable class is also quite subjective in this scenario, as there is no associated label source for this class. Consequently, it is a form of weak supervision [48] that can introduce inadequacies into the data.





Fig. 10. The SV labeling process using developer provided SV labels.

Automatically generated labels have not relied on a third-party label source, but instead have used an additional oracle to provide labels directly for the collected code modules. Two of the primary studies used static analysis tools to provide labels to code datasets [P37, P46]. This approach is noticeably the least considered as it heavily relies on the accuracy of the oracle used for labeling.

Pattern-based labels have been obtained for synthetic data sources, that use SV patterns to generate both the code modules and labels. The goal of an SVP model applied to this synthetic data, is to re-learn the patterns that generated the code modules of each class. These labels are largely considered noise-free, unlike the other two approaches, as the modules and labels are inherently connected.

#### 4.4 Data Cleaning

Data cleaning is the fourth step of data preparation. Whilst data cleaning is important, it is the only data preparation step that is non-essential. Collected code can be labeled and used immediately if it is extracted in an appropriate format. Hence, we observed that not all of the primary studies have discussed this step; nor has it been discussed in as much detail as the other steps.

Code modules are the raw data for SVP, so the data cleaning step involves processing collected code modules so that they are in an appropriate format for feature extraction. Table 5 lists the common data cleaning approaches that we observed from the reviewed primary studies. The cleaning approaches fall under three issue types: irrelevant code, code noise, and duplication.

**Irrelevant Code.** First, irrelevant code has been commonly removed from the collected code modules. For all of the primary studies, this involves removing any code that was not of the target programming language, and removing any code that did not fall under the target granularity, such as code not contained in functions for the function-level granularity. Some studies also removed code that was not of relevance to the project or not at risk, such as test cases, third-party code modules, code scripts, and make files [P1, P4, P30, P56].

**Code Noise.** Second, some studies have further processed the collected code to remove the potential noise in the data that may impact the created features. These steps are only of value to the studies that extract the features directly from the code tokens, while software-metrics are usually robust to such code noise. Several studies removed comments, blank lines and non-ASCII characters [P1, P10, P20, P28, P32, P34, P44, P56], as these are not relevant to SVs. Some studies also replaced user specific tokens, such as user-defined variables, function names, and string literals, with

TABLE 5  
Common Data Cleaning Approaches for SVP Datasets

Issue	Cleaning Approaches
<b>Irrelevant Code</b>	<ul style="list-style-type: none"> <li>Remove code that is not of the target programming language(s).</li> <li>Remove code that is not of the target granularity.</li> <li>Remove irrelevant code files: test cases, third-party-code, scripts and make files.</li> </ul>
<b>Code Noise</b>	<ul style="list-style-type: none"> <li>Remove blank lines, non-ASCII characters and comments from the code.</li> <li>Ignore code with syntax issues or errors.</li> <li>Replace the user-defined variable and function names with generic tokens.</li> </ul>
<b>Duplication</b>	<ul style="list-style-type: none"> <li>Remove highly correlated items.</li> <li>Remove the duplicated code.</li> </ul>

generic tokens to increase the homogeneity of the features [P7, P23, P28, P29, P32, P50, P56]. Some studies also removed the code that they found to have syntax issues or errors [P9, P50], as this may later impact the feature extraction efforts.

**Duplication.** Third, several studies have attempted to remove duplicated code modules, as duplicate entries may introduce bias into a model [49]. Duplicate code modules can be present due to multiple of the same code file, snippet, or software version being collected [P1, P8, P15].

## 5 DATA CHALLENGES (RQ2)

This section identifies the data challenges that researchers have reported in the primary studies, in relation to the data preparation steps. As discussed in Section 3.4.2, we used thematic analysis to analyze the data quality issues that researchers have explicitly reported. These issues were coded, revised and merged by the first two authors of this study. The themes that we have identified are inspired by existing data quality dimensions [37], [38], as data challenges revolve around data quality. Table 6 details the key data challenges and considerations that we observed through our analysis. The data challenges can be summarized as pertaining to data Generalizability, Accessibility, labeling Effort, Scarcity, and both Label and Data Noise.

### 5.1 Generalizability

Generalizability describes the ability for data to extend to other contexts, both in terms of findings and application [38]. Hence, this largely measures the external validity of the produced analysis, based on the data. This challenge primarily involves the data requirements step, as this is the phase where researchers determine the nature of their dataset.

**Ch1: Real-World Representation.** Most of the challenges described in Table 6 arise when using real-world data. Consequently, several researchers have opted to use synthetically created data to construct their models, as described in Section 4.2. Synthetic datasets are artificially crafted to address data challenges present in the real-world data; these data sources are accessible, large, low-effort, and less noisy. As such, they are an attractive option for researchers.

TABLE 6  
Taxonomy of SV Data Challenges Identified From the Primary Studies

Theme	Challenge	Key Points	Paper ID	#
Generalizability	Ch1: Real-World Representation	<ul style="list-style-type: none"> <li>Synthetic data is not representative of real-world code</li> </ul>	P[2-3, 9, 29, 34, 39, 60]	43
	Ch2: External Generalization	<ul style="list-style-type: none"> <li>Data may be language specific</li> <li>Data may be application or domain specific</li> <li>Data may be specific to vulnerability type</li> </ul>	P[1-2, 11, 15, 20-23, 26, 30, 32, 38, 41-42, 44, 46-47, 55-57, 60-61] P[1, 4-5, 11-13, 19-22, 30, 35-36, 40-42, 46, 52, 55-57, 59] P[15, 24, 26, 32, 38, 61]	
	Ch3: Completeness	<ul style="list-style-type: none"> <li>SVs may span code modules</li> <li>Code representation may have limited scope</li> <li>SVs may be present in non-targeted code files</li> </ul>	P[1, 3, 7, 10, 15, 21, 31, 45] P[1, 3, 8, 13, 32, 44] P[1, 12, 20, 46, 53, 57]	
	Ch4: Cold-Start Problem	<ul style="list-style-type: none"> <li>SVs are required to have originally occurred</li> </ul>	P[4, 10, 21-22, 31, 37, 41, 56]	
Accessibility	Ch5: Data Entry Availability	<ul style="list-style-type: none"> <li>Not all data entries are obtainable</li> </ul>	P[1, 5, 11-13, 20-21, 30, 45, 49, 56]	25
	Ch6: Data Privacy	<ul style="list-style-type: none"> <li>Usage of Version Control Systems is unstable</li> <li>Source code and SV data are required to be available</li> <li>Security advisories can be private or vague</li> </ul>	P[1, 12, 56] P[9-10, 16, 26-27, 32, 37-38, 45, 48, 61] P[1, 13, 56]	
	Ch7: Labor Intensive	<ul style="list-style-type: none"> <li>Manual labeling is highly time-consuming</li> </ul>	P[2, 10, 18, 26, 29, 31, 39, 42-43, 61]	
Effort	Ch8: Expertise Requirements	<ul style="list-style-type: none"> <li>Manual labeling requires high expertise</li> <li>Vulnerabilities are difficult to identify</li> </ul>	P[2, 18, 31, 43] P[20, 30, 33-34, 42-43]	14
	Ch9: Data Imbalance	<ul style="list-style-type: none"> <li>Vulnerable samples are the extreme minority</li> </ul>	P[1, 4-5, 8-9, 11-12, 14-15, 17, 19-20, 22, 24, 30-31, 33, 36, 40, 45, 50, 53-56, 60]	
Data Scarcity	Ch10: Number of Samples	<ul style="list-style-type: none"> <li>Low number of vulnerability samples</li> </ul>	P[1, 4, 10-11, 15, 20, 23, 29-31, 33-34, 36, 45-47, 50]	32
	Ch11: Incomplete Reporting	<ul style="list-style-type: none"> <li>Latent, dormant or unresolved SVs can exist in the dataset</li> <li>SVs can be silently patched</li> </ul>	P[1, 4, 11-13, 16, 19, 21, 28, 30, 36-37, 47, 56, 59-61] P[4-5, 12, 30, 33, 36, 40, 47, 52, 56]	
Label Noise	Ch12: Localization Issues	<ul style="list-style-type: none"> <li>Commit noise causes localization issues</li> <li>Data noise causes localization issues</li> <li>Bug reports do not document code location</li> <li>Version tracking is complex and erroneous</li> </ul>	P[1, 16, 29, 42] P[2, 8, 13, 28, 32, 44] P[28, 30, 33, 42, 56] P[1, 24, 30, 56]	31
	Ch13: Erroneous labeling	<ul style="list-style-type: none"> <li>Manual labeling can be inaccurate or subjective</li> <li>Static analysis tools label modules inaccurately</li> <li>SVs may not actually be exploitable</li> <li>Label quality is unknown</li> </ul>	P[11, 18, 40, 43, 52, 56] P[2, 46] P[1, 52] P[24, 28, 32]	
	Ch14: Code Noise	<ul style="list-style-type: none"> <li>Source code has stylistic differences or syntax issues</li> <li>Binary code is noisy</li> </ul>	P[7, 9-10, 13, 18, 20, 23, 28-29, 32, 34, 38, 41, 50-51, 61] P[27, 48, 53]	
Data Noise	Ch15: Redundancy	<ul style="list-style-type: none"> <li>Some entries are indistinguishable between classes</li> <li>Code versions and localization can add redundancy</li> <li>Vulnerable samples have limited diversity</li> </ul>	P[49] P[1, 8, 12, 15, 19-20, 24, 30, 42, 46] P[14, 28, 40]	35
	Ch16: Heterogeneity	<ul style="list-style-type: none"> <li>Data contains outliers</li> <li>Poor cross-project performance</li> </ul>	P[25, 42] P[1, 12, 26, 39, 56]	

Given that synthetic data may not represent the real-world data, it is considered a big limitation that may render such a dataset unusable unless this limitation is addressed. Synthetic vulnerability examples are considered to be simpler, isolated, less diverse and cleaner than real-world vulnerabilities [P2, P3, P9, P29, P34, P39, P60]. Zheng *et al.* [P29] found that the use of synthetic data sources may significantly inflate the reported model performance in comparison to the models using real-world code. Hence, a model trained using synthetic data is unlikely to be able to detect complex real-world vulnerabilities, which require much deeper semantic understanding and reasoning [40]. Thus, real-world data is the more commonly used data source, as seen in Fig. 8.

**Ch2: External Generalization.** Nearly all studies face external threats to validity of their findings inferred from a specific dataset. In terms of SVP research, this relates to the limited application scope of the selected study datasets.

Datasets may be specific to, and hence have troubles generalising outside of: **programming language** [P1, P11, P15, P20, P21, P22, P23, P25, P26, P30, P32, P38, P41, P42, P44, P46, P47, P55, P56, P57, P60, P61], application or domain type [P1, P4, P5, P11, P12, P13, P19, P20, P21, P22, P30, P35, P36, P40, P41, P42, P46, P52, P55, P56, P57, P59], and SV type [P15, P23, P26, P32, P38, P61].

**Ch3: Completeness.** Completeness is achieved when a dataset has all the relevant parts of an entity's information, which is sufficient to represent every meaningful state of a real-world system [38]. **However, the selected data for analysis can have a limited scope of the overall system, which makes their application context limited.** This data-oriented consideration serves as a challenge for SVP, as it prevents these models from forming a "complete" solution. **First, if the selected granularity of code modules is too fine, researchers are forced to ignore vulnerabilities which span multiple modules** [P1, P3, P7, P10, P31, P45]. For instance, function-level prediction

is unable to predict more complex SVs that span multiple functions. The selected semantic representation of data may also not consider all sources of weaknesses in a software system [P3, P8, P32, P44]. For instance, Tian *et al.* [P3] and Li *et al.* [P8] only consider code snippets of library and API function calls, which would not cover all potential SVs in a system. Static source code is also unable to capture certain necessary dynamic code traits [P1, P13], such as crashes and memory leaks. Similarly, vulnerabilities may be present in code modules which are not of the target programming language of analysis [P1, P12, P20, P46, P53, P57]. In the modern development landscape projects commonly utilize multiple programming languages [50], but SVP models are largely targeted towards a single programming language of choice [P53].

## 5.2 Accessibility

Accessibility describes the ability to retrieve or obtain data from the target data sources [38]. Challenges arise from difficulties in accessing the data, either of the raw code modules during the data collection step or of the data labels during the data labeling step.

*Ch4: Cold-Start Problem.* The cold-start problem is an issue originating from recommender systems in which a system is unable to draw inferences about incoming modules for which it has not yet gathered sufficient information [51]. In terms of SVP, the cold-start problem has been particularly present, as to make future predictions, we require vulnerabilities to have originally occurred and to have been documented [P4, P10, P31]. This makes SVP largely infeasible for new or immature organisations [P10, P21, P22, P31, P37, P41, P56]. Furthermore, the acquisition of initial high-quality training data is a major issue, as seen in the other challenges.

*Ch5: Data Entry Availability.* Since a majority of SV data have been obtained through mining open source repositories, not every part of a system would be necessarily accessible to researchers. For instance, some source code might have been unavailable due to lack of inclusion in public repositories [P1, P11, P12, P20], or unobtainable due to other unstated technical reasons [P5, P21, P45]. Similarly, some vulnerability reports might not have been able to be localized to code modules due to issues in the automatic or manual localization methods [P1, P13, P30, P49, P56].

Some researchers have even pointed out that the reliance on a version control system to track code modules causes issues in itself, as consistent usage of a version control system is unstable [P1, P12, P56]. Version control systems were only widely adopted in 2005 with the introduction of *git*, hence data before this date would be irretrievable [P12, P56]. Furthermore, organisations might switch the version control system they were using, losing previous software history [P1].

*Ch6: Data Privacy.* The potential commercial sensitivity of both software code and SV reports means that organizations are often not willing to share private-source code or data to researchers [52]. This data privacy creates many data accessibility issues. First, several researchers have observed that commercial systems have not provided their source code [P9, P10, P26, P32, P45, P48, P61], making SVP via source code on these systems infeasible. Furthermore, organisations and practitioners might desire to limit the availability of their security advisories by making them private or

vague [P1, P13, P56]. By concealing information about vulnerabilities, it is theoretically harder for an attacker to exploit a system. Similarly, an organisation might not even maintain a public security advisory or document SVs [P16, P27, P37, P38].

To represent real-world data, researchers have often surreptitiously avoided this issue through the use of open-source repositories that have public vulnerability records. Without open sources, data retrieval and reporting can become very difficult due to commercial sensitivity. Whilst this is valid, open-source data is not representative of software engineering practices as a whole; it is unknown whether the derived observations will generalize to private-source code and practices. Only two out of the 61 primary studies used private source data [P5, P11], both of which suffered from data entry availability (Ch5) as a result.

## 5.3 Effort

Effort describes the amount of human-effort required to label a dataset [38]. The standard approach for traditional supervised learning has been to have a subject matter expert hand-label a dataset. However, this is largely infeasible for SV data due to the extreme effort requirements. As such, many researchers have skirted around this challenge by using synthetic labeled data sources or reusing existing datasets. Consequently, this theme was actually the least mentioned by primary studies, as researchers that reused datasets often did not report or discuss the effort.

*Ch7: Labor Intensive.* Labeling code or bug reports as SV-related is a non-trivial task, which makes it highly labor intensive when coupled with the sheer number of modules to examine [P2, P18, P29, P43, P61]. Dowd *et al.* [53] estimated that one hour of security review can cover an average of 500 lines of code. However, most modern software systems contain millions of lines of code, which makes the required man-hours infeasible. Zhou *et al.* [P2] stated that it took 600 man-hours to manually curate their SV dataset. Manual labeling is ultimately the most reliable labeling approach however, due to the large amount of noise for automated labeling (Ch13).

*Ch8: Expertise Requirements.* Secure code review requires significant security expertise [P2, P18, P31, P43]. To successfully perform secure code review, a practitioner/researcher must have the capability to memorize and recognize thousands of security-related patterns and concepts [54], and this list of required knowledge is continually growing. Furthermore, it has been highly difficult to identify SVs in comparison to regular defects [P20, P30, P33, P34, P42, P43], as they do not necessarily represent functional bugs and are thus hard to verify.

## 5.4 Data Scarcity

Data scarcity refers to the extent to which the quantity or volume of the available data is appropriate for the task at hand [37], [55]. As SVP is a data-driven process, it requires large volumes of data [49]. For SV data, this challenge largely represents the low number of real-world vulnerable examples available. Whilst codebases are often sufficiently large, the number of identifiable vulnerable modules in a codebase is relatively small [8].



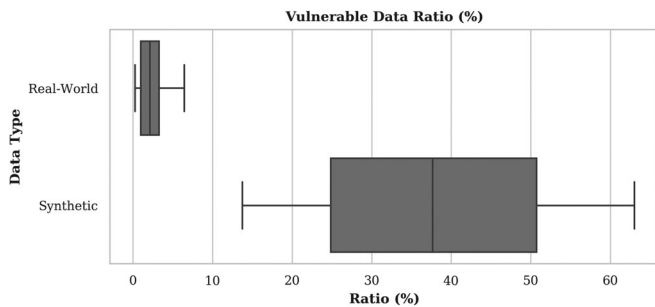


Fig. 11. The percentages of vulnerable files in datasets utilized in the primary studies.

**Ch9: Data Imbalance.** The severe imbalance of vulnerable modules to non-vulnerable modules has been a challenge reported by many researchers using real-world datasets. To quantify this issue, we report the percentages of vulnerable modules in each dataset utilized in the primary studies, displayed in Fig. 11. We display the real-world and synthetic datasets separately, as the synthetic datasets have been artificially constructed to over-represent the vulnerable class. Real-world datasets which were artificially altered to be balanced, or merged datasets consisting of both synthetic and real-world examples are excluded from Fig. 11.

This has been a considerable issue for SVP research, as learning-based models are optimized to perform on balanced classes [56]. The severe class imbalance issue present in real-world SV data can lead to biased classifiers that favor the non-vulnerable class [P15, P19, P24, P30, P31, P33, P50, P54]. This challenge has been referred to as *finding a needle in a haystack* [8] and is notably unique to security defects. Shin and Williams [P33] observed that the number of reported faults was seven times larger than that of vulnerabilities. Nguyen and Tran [P24] observed that class imbalance increased as the module granularity became more fine-grained.

**Ch10: Number of Samples.** Similar to **Ch9**, the severe imbalance of data leads to a very low overall number of samples for the vulnerable class. This has been a major blockade for SVP research as learning-based methods have strict data requirements; they need a large quantity of historical cases to learn from. It is common knowledge in the ML community that *more data beats a cleverer algorithm* [49]. Several studies have particularly exacerbated this issue for DL-based methods, which require even greater data size [P15, P23, P29, P34, P45, P47, P50]. This is an emerging challenge given the rise of DL-based methods for SVP [13].

## 5.5 Label Noise

Label noise relates to whether the labels in a dataset accurately represent the ground truth and are free of error [37], [52], [55]. As SV datasets have been rarely hand-labeled by subject matter experts (excluding synthetic datasets), but instead mined from historical artefacts, a large amount of noise has been introduced into the SV labels. This severely impacts the reliability of SV data.

**Ch11: Incomplete Reporting.** A major problem of using historically reported vulnerabilities to label real world data is that we only have a label source for the vulnerable class, and simply treated the remaining labels as non-vulnerable, creating uncertainty in the labels for the non-vulnerable

class. Modules in the non-vulnerable class are not actually confirmed to be clean, just that no vulnerabilities have been historically reported. However, in reality, there can be dormant, latent or unreported vulnerabilities existing in these modules [P1, P11, P12, P13, P16, P19, p21, P28, P30, P36, P37, P56, P59, P60, P61]. Hence, the non-vulnerable class may be better considered as unlabeled [P21, P47].

Similarly, researchers are limited to the use of *fixed vulnerabilities* as a label source [P52]. Unresolved vulnerabilities have been rarely disclosed as they can be exploited by attackers. Hence, during data labeling, these unresolved SVs would actually be contained in the non-vulnerable class.

Furthermore, the reliance on vulnerabilities to be reported has created a temporal issue for data collection. SVs usually take time to be detected and fixed, and hence a different number of SVs will be documented depending on the time of data collection [P21, P36]. Jimenez *et al.* [P36] observed that the performance of models significantly decreased when only using vulnerabilities observed before the time of model training.

Using bug reports for labeling also has created a reliance on developers or organisations to have thorough reporting practices. However, in reality, some organisations have patched some vulnerabilities “silently” [P36, P40, P56], without any documentation provided for bug reports or security advisories. Furthermore, the difference between SVs and non-security related faults can sometimes be minimal [P5, P12, P30, P33]. Hence, many security defects have not been reported by developers as such [P4, P52]. An organisation may also use multiple bug reporting systems; a reliance on just one data label source (i.e., NVD) would be incomplete [P36, P40].

**Ch12: Localization Issues.** Bug reports and SV records have not always documented the location of SVs, posing a large challenge for the retrieval of SV data [P56]; the vulnerable code modules are often not explicitly listed. Hence, researchers have often relied on the use of patches to trace the location of a vulnerability. This is flawed as not every documented vulnerability has an associated patch [P28, P30, P33, P42, P56], as it may be concealed for privacy or not yet resolved. Furthermore, patches may not always properly disclose the true location of an SV [P1], as patches may instead provide workarounds for separate modules, rather than a fix of the underlying problem. Jimenez *et al.* [P12] found that only 75% of vulnerability reports had an associated fix, which leaves the remaining 25% of reports untraceable. Vulnerability reports are often incomplete and missing references [9].

Furthermore, tracking SV location from a bug fix has been non-trivial. Version information in bug reports is often unreliable [P1, P24, P30, P56], as identifying affected software releases for a vulnerability is highly difficult. Due to the evolving nature of code, the assumption that all previous versions of code were also vulnerable is invalid [P24], and reporting vulnerabilities in prior versions is inaccurate as developers have little benefit from expending efforts to do so [P1]. Furthermore, commits can be noisy [P1, P16, P29, P42] due to tangled code changes. A vulnerability fixing commit may not exclusively patch a vulnerability [57]; other functional changes may be included. A vulnerability fix can also be buried as part of a larger commit including non-security changes. Herzig *et al.* [57] showed that tangled

commits had significant impacts on defect labeling, with an average of 16% of files being mislabeled as a result.

*Ch13: Erroneous Labeling.* Label inaccuracies can additionally come from various other sources. First, with noisy labels stemming from **Ch11**, as well as errors arising from localizing labels to code modules (**Ch12**), manual labeling has been a common approach to help ensure data quality. However, this approach is erroneous in itself. As discussed in **Ch8**, manual labeling is difficult and requires high expertise. Hence, it is inevitably an error-prone task [P18, P40, P43, P52]. Furthermore, the process of labeling modules as vulnerable or not can even be quite subjective [P11, P56]. There is no clear definition of the difference between a vulnerability and a fault, and hence this distinction can be nebulous to a human [P30]. For instance, if a regular function calls a vulnerable function, it is unclear whether this function should also be considered vulnerable [P1].

Some studies have utilized static analysis tools and methods to achieve automatic data labeling without the need for historical vulnerability reports [P12, P37, P46]. However, researchers have observed these methods to be highly inaccurate and hence introduce considerable noise into data labels [P2, P46]. This process can also be flawed from a motivational perspective, as the SVP model is simply relearning the patterns used by the static analysis tool to infer the labels.

Patched vulnerabilities may not actually be exploitable in the real-world either [P1, P52]. Developers may incorrectly hypothesize security weaknesses, or err on the side of caution. This adds unreliability to the accuracy of the labels in the data source, as the SV labels may actually be benign.

Another large and open challenge has been the lack of measures to quantify label quality [P24, P28, P32]. There is no trivial way to measure or quantify the aforementioned label noise issues. Furthermore, with the reliance on historical artefacts for labeling, some sources of label noise are undetectable or unverifiable, e.g., SVs remain dormant (**Ch11**) until they have been detected. This severely impacts the reliability of SVP research, as it is unclear whether the findings have been made using valid data.

## 5.6 Data Noise

Finally, data noise refers to the noise within the raw data entries [38]; code modules for the purposes of SVP. Noise and inaccuracies in these modules may negatively affect the data and any produced features used to train a model, lowering the potential efficacy of that model.

*Ch14: Code Noise.* SVP uses code as the raw data source. However, source code is noisy, which consequently impacts the effectiveness of any produced model. Developers have different coding styles and naming conventions [P28, P34, P38], which adds inconsistency. Li *et al.* [P61] further identified that different projects may have different code quality due to differences in coding practices and guidelines. Real-world code can also contain syntax issues [P9, P50]. These sources of noise can severely impact the versatility of produced SVP models, as they may instead learn specificities of particular coding styles and syntax.

Furthermore, binary code is usually much noisier than regular source code. Binary code snippets can be difficult to trace and identify [P27, P53], or interspersed code and variables can become indistinguishable [P27, P48].

*Ch15: Redundancy.* Redundancy refers to undesirable duplication in a dataset. Too much redundancy in the training data, can lead to bias and overfitting for an SVP model. The major source of redundancy in SV datasets has come from code modules having several different versions and revisions. Datasets that consider versions separately can introduce redundancy into the code entries, as the majority of the code usually stays constant between revisions [P20]. Vulnerable labels can also be duplicated over several versions, as modules can remain vulnerable for an extended period of time [P12, P19, P24, P30, P46]. Code branches are another potential source of redundancy to labels and modules as the majority of code and data is often duplicated across branches [P1, P42]. Automatic extraction of code snippets and program slices can additionally introduce duplication [P8], as duplicate program slices can be created from different entry points.

Vulnerable samples can also not be very distinct from each other [P14, P40], due to consistent SV type or exploit patterns, which limits the learning capacity of an SVP model. This is particularly present for synthetically created samples [P28]. Another issue is that vulnerable and non-vulnerable entries have limited diversity [P49]. Vulnerability patches can only alter a few lines of code, making the majority of the module consistent between its vulnerable and non-vulnerable versions. This is a particularly significant issue as if the model cannot learn these subtle distinctions, it will produce high false positive/negative rates.

*Ch16: Heterogeneity.* Code modules have been observed to be highly heterogeneous, which negatively impacts the diverse application of the produced SVP modules. For example, the coding conventions of one code module often do not match another, due to differences in authorship, functionality or coding style. Learning-based methods operate best when the data, especially the training and test distributions, are homogeneous so that the learnt patterns can be applied uniformly [49]. However, researchers have observed data distributions to be irregular or containing outliers, hence requiring normalization to reduce irregularities [P25, P42]. Similarly, several researchers have observed that SVP models perform poorly in a cross-project setting [P1, P12, P26, P39, P56], due to the heterogeneity of these datasets; the coding conventions and functionality of one project rarely mirror another.

## 6 DATA CHALLENGE SOLUTIONS (RQ3)

In this section, we present the various solutions that researchers have presented in the reviewed studies to solve data challenges or help improve data quality. Fig. 12 displays the main areas of the solutions that we have identified. We again used thematic analysis, described in Section 3.4.2, to identify the categories of the identified solutions. We mapped the solutions to the data challenges based on the data challenge themes that the solutions were connected to when reported in the primary studies.

We provide an overview of the main solutions that have been considered for the identified data preparation challenges. As discussed in **Ch1**, the majority of the data challenges arise from the use of the real-world data; challenges for Accessibility, Effort, Data Scarcity and Label Noise are



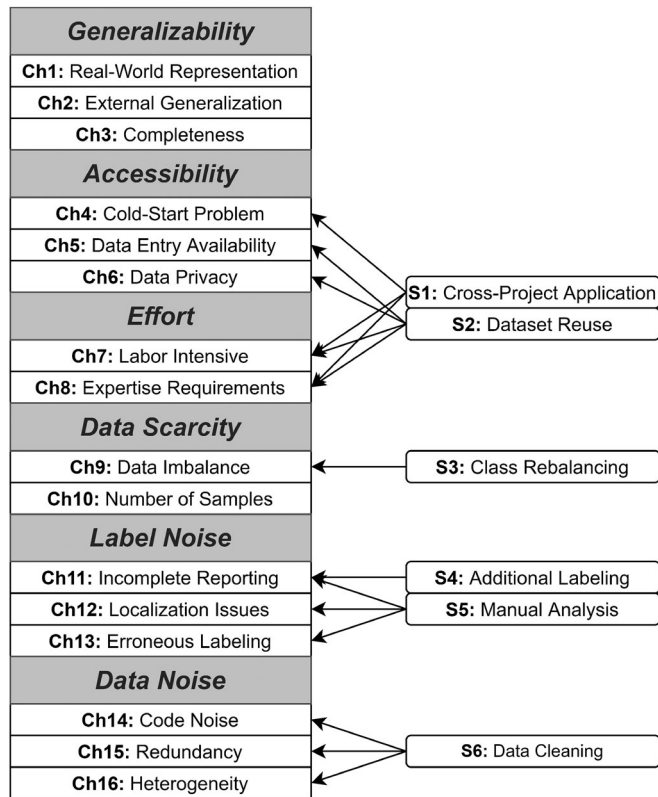


Fig. 12. A mapping of solutions onto challenges.

largely unique to the real-world datasets. Consequently, the majority of the data solutions are similarly aligned with the real-world data challenges. We note that not every primary study provided remediation for the reported data issues.

We also note that several of the data solutions were provided to these challenges at a model-level, rather than from the data perspective. For instance, several studies supported the use of ensemble models as they are more robust and hence less susceptible to noise and class imbalance [P42]. Furthermore, feature processing techniques, such as feature selection, were another method used to help remove data noise and increase generalizability [P22, P25, P40, P42, P47, P55]. However, we focus on the data preparation steps, so an analysis of these model-level solutions is out of the scope of this SLR.

**S1: Cross-Project Application.** As discussed in **Ch4**, researchers have been motivated to overcome the cold-start problem and cost of acquiring a dataset. This has led to several efforts to apply cross-project SVP models for which a previous project's dataset is used to train an SVP model that can be applied to a new project. Due to the extreme cost of acquiring data (**Ch7&8**), the cold-start problem is thought to be further exacerbated for SVP. Hence many researchers considered cross-project SVP as an essential solution to this issue [P1, P10, P18, P21, P22, P26, P31, P37, P41, P43, P56].

Whilst this solution does solve the challenges that it seeks to address, it also adds challenges by introducing heterogeneity to the data (**Ch16**), which consequently impacts a model's performance. Furthermore, underlying problems in the original data source can still be present.

**S2: Data Reuse.** With the significant challenges in the data labeling effort (**Ch7&8**), coupled with the challenges in data

accessibility (**Ch5&6**), it is a more efficient choice to reuse or augment the existing datasets. This not only significantly saves time and effort in data preparation, but also enables researchers to conveniently evaluate and benchmark performance on the same datasets. As dataset reuse greatly reduces effort, it is more desirable than self-construction of a dataset. Hence, researchers have often reused the existing datasets when available; 28 out of the 61 reviewed studies reused or augmented the existing datasets in some manner.

Like **S1**, whilst this solution certainly solves the challenges regarding Accessibility and Effort, it does not address the problems in the original dataset like Generalizability, Data Scarcity, Label Noise and Data Noise. However, many researchers have seemingly assumed the validity of the prior datasets, as the studies that utilized former datasets often did not discuss other data challenges.

Data reuse is achieved through data sharing efforts of other researchers. Several researchers made their datasets publicly available for use to assist in construction of SVP models by researchers and practitioners [P1, P8, P10, P14, P16, P26, P31, P32, P36, P38, P44, P52, P56, P61].

However, the actual usability of these datasets can create further issues for this solution. First, the quality and reliability of the information provided in the existing datasets is unknown and unverifiable [P4, P8, P12, P20, P28, P38, P39, P40, P44, P54]. Researchers often prefer to use a dataset about which they have complete knowledge. Researchers also build their SVP models to fit a variety of application contexts, such as specific granularities, programming languages, or SV types. Hence, several researchers found the information provided in the existing datasets insufficient to be applicable to their desired applications' contexts [P2, P4, P12, P16, P23, P27, P31, P39, P42, P43, P48, P56].

Furthermore, although many SV datasets have been created, they are not necessarily available. Researchers have reported that the existing datasets are private [P2, P16, P31] or unavailable [P2, P12, P42]. We observed that five of the shared datasets from the reviewed studies have since become unavailable due to dead links. Similarly, researchers may not share the code used to produce a dataset, which impacts the verifiability and reproducibility. Such problems of availability or usability of the desired datasets can lead researchers to use self-constructed datasets.

**S3: Class Rebalancing.** The severe imbalance of vulnerable to non-vulnerable modules (**Ch9**) is reported by almost half of the reviewed studies (27 out of 61) as a significant data challenge. This imbalance leads to models that bias towards the majority class [49], and do not fairly consider the minority vulnerable class. As such, many studies have employed some form of class rebalancing to help remediate this issue. These rebalancing techniques fall into two main categories: **Undersampling** [P1, P8, P14, P15, P19, P20, P22, P30, P33, P52, P53, P55, P56, P60] and **Oversampling** [P8, P14, P17, P20, P36, P45, P56]. Undersampling is the process of removing samples from the majority class to match the size of the minority class, whereas oversampling duplicates or synthetically adds samples to the minority class until the size matches the majority class. Undersampling was the more popular technique used in the reviewed studies; 14 studies used undersampling in comparison to 7 that used oversampling. Researchers considered undersampling to be more

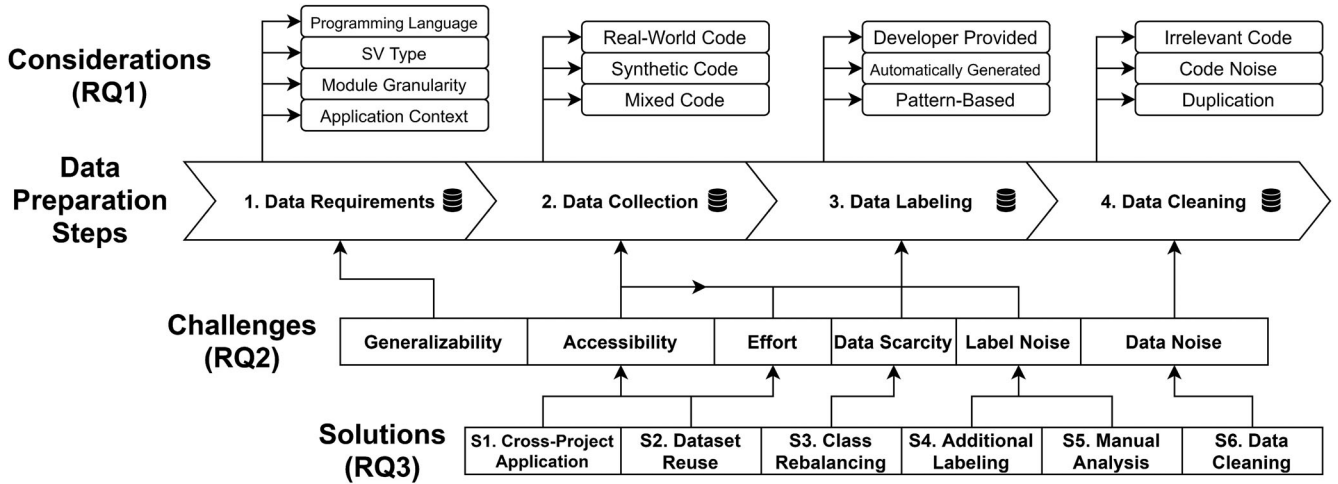


Fig. 13. An overview of the SVP considerations, challenges and solutions for data preparation.

standard [P22, P55], effective [P30, P60] or efficient [P19, P30, P60], in comparison to oversampling. However, we note that undersampling may not always be desirable as it reduces the overall amount of data, which can have significant impacts for SV data due to the low number of samples in the minority class (Ch10). Oversampling does not suffer from information loss, but adds redundancy to training data (Ch15), which can lead to overfitting.

This solution is largely considered a necessity, as the positive impacts of class rebalancing have been widely agreed upon by the community. Furthermore, several studies have explicitly demonstrated the performance increase of models when trained on balanced data in comparison to imbalanced data [P7, P54, P60].

Alternatively, some researchers artificially constructed both their training and test datasets to be balanced [P61]. Balancing the test set like this artificially inflates the model performance, however, as the incoming real-world data will not be balanced [P1, P15]. This is a significant weakness of synthetic datasets that is expected to impact their real-world generalization, as they are constructed to over-represent the vulnerable class.

**S4: Manual Analysis.** The poor labeling provided from bug reports can add a lot of label noise and inaccuracies in the data. To combat this, researchers have often assisted the labeling process through manual analysis. Over 37% of the reviewed studies (23 out of 61) assisted their data collection process with manual inspection. This allowed researchers to manually spot sources of label noise [P1, P2, P4, P7, P8, P10, P11, P14, P25, P26, P28, P29, P30, P31, P32, P42, P44, P35, P46], or made localization to code modules more accurate [P1, P2, P4, P7, P10, P14, P26, P31, P44, P45, P52]. It also helped researchers to obtain an understanding of the quality of their dataset [P13, P16, P32, P42].

Whilst manual labeling assistance can greatly improve the quality of a dataset, manual analysis is problematic in itself. It is highly effort intensive (Ch7), difficult (Ch8) and error-prone (Ch13).

**S5: Additional Labeling.** Manual analysis can only resolve issues in the vulnerable class, as the use of bug reports only provides a source of labels for this class. In reality, however, dormant or latent vulnerabilities can introduce unverifiable

inconsistencies in the non-vulnerable class (Ch11). Hence, several researchers have used automated methods to obtain additional labels. One such approach was to use static analysis tools to label modules, in an attempt to uncover latent vulnerabilities [P37, P46]. However, this process added considerable noise to the labels in itself [P2, P46]. Li *et al.* [P61] used a rule-based approach to obtain the non-vulnerable modules, by filtering out functions which may have security relevance. This decreased the chance of having latent vulnerabilities in the non-vulnerable set. As these additional labeling methods often introduce their own assumptions or inaccuracies, their effectiveness is unclear.

**S6: Data Cleaning.** Data noise problems can be addressed through data cleaning of code modules. Table 5 from Section 4 lists common data cleaning techniques for SV modules. Although the real world code is usually much noisier than the synthetic code due to inconsistent coding styles [P28, P34, P38], cleaning can still be necessary for the synthetic data sources, depending on the method used to construct the vulnerable examples. Fang *et al.* [P50] found a large proportion of coding errors in the synthetic SARD and SQLI-Labs datasets, that they manually remediated.

Data cleaning is one of the four data preparation steps [18], and consequently this solution has been largely considered a necessity. For example, Zheng *et al.* [P28] showed that replacing user defined strings with generic tokens improves a model's performance.

## 7 RECOMMENDATIONS

Based on the findings of this review, we have drawn some actionable recommendations for researchers and practitioners. In order to place our recommendations in the context of the findings from this review, Fig. 13 presents an overview of the findings to our three RQs. Fig. 13 shows the considerations, challenges and current solutions for the data preparation steps of SVP in order to help readers to better understand and interpret our recommendations in this section.

The validity of SVP research largely relies on the quality of the data used to construct an SVP model. However, as identified in Section 5, researchers have faced a plethora of

data preparation challenges that may lead to serious pitfalls that should be avoided while training SVP models.

Due to the current shortcomings of SV datasets and preparation processes, there is an important need for evidence-based research of the existing practices, challenges and solutions in order to advance the state of the art and the state of the practice. Predominantly, there is a need for advancing techniques of real-world data labeling, or creation of more realistic synthetic datasets. Whilst there are several research studies on various related sub-fields, such as fault localization [58], bug classification [59], bug seeding [60], or crowd testing [61], their findings have yet to be leveraged for providing reliable solutions to address the open research problems. Our identification of the existing obstacles to achieving high quality data for SVP is expected to help researchers to develop suitable methods for overcoming the current barriers for improvement.

Several studies have called for a need of a *gold-standard* dataset [4], [13], [16]; one that overcomes the identified data preparation challenges. We have identified several areas of data preparation solutions in Section 6 that researchers have used to help address the data preparation challenges. Whilst we recommend that researchers and practitioners also adopt these solutions as they provide an initial foundation for data quality, these solutions are, by no means, complete to produce a gold-standard dataset. Particularly, many challenges for Generalizability remain unaddressed, and the solutions for Label Noise are insufficient. Manual analysis (S5) is often infeasible or inaccurate, and additional labeling (S6) can add additional noise or inconsistencies into a dataset. Furthermore, SVP research has operated at a variety of different contexts, which makes the use of standardized datasets difficult. Different researchers may favor different levels of granularity or applications domains.

That is why rather than proposing the guidelines for developing a gold-standard dataset, we instead propose a set of recommendations for the future research directions that are expected to help improve the SV data preparation processes by addressing the identified data challenges. Although there are many potential recommendations that can be made, we have formed our discussion from actionable observations obtained from the primary studies and supporting literature, that have the potential of delivering immediate benefit to SVP researchers and practitioners.

**R1. Consideration of Label Noise in the Negative Class.** Fully supervised learning for SVP requires a positive class (vulnerable modules) and a negative class (non-vulnerable modules). The quality of both of these classes is highly important to ensure that a model is able to learn the appropriate patterns in the data. Whilst much efforts and manual inspection have been expended to ensure the quality of the vulnerable labels [P1, P2, P4, P7, P8, P10, P11, P14, P25, P26, P28, P29, P30, P31, P32, P42, P44, P35, P46], equal efforts have not been expended for the non-vulnerable labels. The most commonly used assumption is to simply take all remaining unlabeled files to form the negative class. However, due to incomplete reporting (Ch11), we have observed that there can be considerable inaccuracies in the non-vulnerable class; researchers have only been able to obtain a data source for the vulnerable labels (i.e., the bug reports), but there is no source of labels for non-vulnerable data.

The unavailability of labels for non-vulnerable data adds inherent label noise to the non-vulnerable labels of SV datasets. Such label noise would affect the reliability and trustworthiness of the produced SVP models, as incorrect patterns may be inferred from the mislabeled instances. Hence, there should be more consideration of additional broader label sources to reduce label noise. With proper consideration, additional label sources such as static analysis [P37, P46] can help uncover mislabeled vulnerabilities. Zheng *et al.* [62] used differential analysis to increase the confidence of static analysis labels for vulnerability fixes. Crowd-testing [61] may be another potential label source that can minimize individual effort; crowdsourced labels are already used for a variety of benchmark datasets in the ML domain [63]. Producing more vulnerable labels would also assist in dealing with the challenges in Data Scarcity (Ch9 & 10).

However, ensuring the absence of vulnerabilities from code modules is difficult [64]. As high-quality real world datasets are currently unavailable, we also suggest a relaxation of labeling and fully-supervised learning requirements. The negative class may be more accurately considered as an unlabeled set. In this scenario, semi-supervised methods, such as self-training or Positive and Unlabeled (PU) learning can be applied to help uncover knowledge from the unlabeled set [65].

**R2. Consideration of Timeliness.** Timeliness describes the temporal aspects of data [38]. There are two main components for timeliness. The first is the currency of data; the age of data in use compared to when it was collected. Practitioners ought to generally avoid using out-of-date data as it can lose relevancy to contemporary settings. This is particularly an issue for SV data, as vulnerabilities take time to be discovered (Ch11), hence, delayed data collection effort can usually obtain more complete information [11]. For instance, 10 vulnerabilities may have been reported for a codebase after three months, but an additional 10 vulnerabilities may have been reported one year later. Hence, a lack of consideration for data timeliness will lead to mislabelled instances and a lack of completeness. Researchers ought to develop methods for updating SV datasets to ensure they use more current label data and increase the number of positive samples. This will ensure the datasets are more complete upon use, and consequently allow for more reliable SVP models.

The second aspect of timeliness relates to the temporal nature of the data accumulation; historical artefacts are created incrementally over the lifecycle of a project. This poses an issue for SV data due to concept drift; the vulnerability data and patterns change over time with the emergence of new concepts [66]. Failure to account for this temporal nature, that is, preserving data order for model training and validation, has been shown to produce unreliable models and impact the real-world generalizability [67], [68]. However, only 13% of the reviewed studies (8 out of 61) considered a time-based ordering of the data. Hence, there needs to be further investigations into concept drift and use of time-based ordering in order to produce more reliable SVP models.

**R3. Use of Data Visualization.** Data understandability describes the ability to comprehend data [37]. Exploratory data analysis is an integral part of data science, as it helps practitioners to understand data quality issues, imbalances,



and relationships within data. This is particularly important for ML as it allows practitioners to identify the necessary data cleaning practices and to conduct feature engineering [69]. SV data can be complex due to the intricacies of the code and security weaknesses. Despite this, not much effort has been reported by researchers into data understanding and visualization.

Data visualization is one of the most powerful techniques for data understandability [69], but only a few of the reviewed studies did visual exploration of their data. Neuhaus and colleagues [P13] visualized the locations of SVs in a codebase and observed that the distribution of SVs are scarce and irregular. This motivated their search for specific code patterns that can be used to describe the irregular distribution. Chowdhury and Zulkerine [P19] used data visualization to identify that files are unlikely to contain SVs in multiple different versions; hence, they determined that vulnerability history of a module may be a poor indicator. The future efforts in data understandability for SV datasets, using similar techniques to the aforementioned ones, may be able to yield novel impactful insights for this data, and assist in SVP model creation. Better data understandability can also help practitioners who lack data science expertise with setting up SVP models.

**R4. Creation and Use of Diverse Language Datasets.** We observe in Fig. 6 that researchers have primarily constructed datasets for singular languages, as the semantic and syntactic differences of programming languages creates difficulties towards applying models across languages. However, single-language datasets pose inherent coverage problems (Ch3) as modern software development is typically not conducted using a single language [50]. SVs can also be present in cross-language operations that use foreign function interfaces. Furthermore, as dataset creation is expensive (Ch7), the lack of multi-language datasets creates scalability problems if practitioners intend to add or switch programming languages. Existing SV datasets have only been constructed for popular languages such as C/C++, Java or PHP, which makes application of SVP models difficult to more obscure languages. Hence, there is need for techniques that can efficiently create or utilise diverse language datasets. Our review found only one study [P53] investigated a cross-language approach.

Software engineering researchers need to increase their efforts for developing datasets supporting cross-language analysis by leveraging the outcomes of the research in the related areas of software engineering. For example, multilingual source code analysis is an emerging research domain [70], that has seen success predominantly in static analysis. Bug seeding is another relevant research field that intends to import bugs from the existing projects to new ones [60]. We speculate that bug seeding can be used to import bugs from one language into a dataset of a semantically similar language, hence, allowing efficient dataset creation of obscure languages. SVP research will benefit from investigation and adoption of these techniques. More diverse language datasets will increase the completeness and scalability of SVP models for practitioners.

**R5. Use of Data Quality Assessment Criteria.** We observed in Ch13 that the quality of most SVP datasets is unknown and even unconsidered. If researchers and practitioners are

able to determine the quality of their datasets, they are expected to make informed decisions about the validity and reliability of the constructed SVP models. Data quality is an inherent requirement of most ML-based systems [18]. Hence, data quality assessment criteria ought to be developed to assess SV datasets in order to enable quantifiable analysis and verification of data quality issues.

Whilst data quality assessment is a common practice for organisations and data-centric industries [37], [38], general data quality dimensions may not be directly applicable or measurable for SV datasets. These generalised data quality patterns need to be customized for SVP. Specific data quality assessment criteria can be determined through requirements-driven approaches [71]. Our categorisation of SV data challenges in Table 6 is expected to potentially help in identifying the relevant data quality dimensions and requirements.

**R6. Better Data Sharing and Governance.** Although data sharing and reuse have been popular for SVP (S2), we observed that there are many issues regarding the availability of these datasets. Furthermore, the incomplete reporting practices of SV data has led to a lack of trustworthiness. The reporting for the datasets has often been insufficient to allow for proper replication, making the datasets hard to validate. For instance, many of the reviewed studies did not report the version of the data source or specific extraction steps. Furthermore, most of the data preparation processes use extensive manual inspection or labeling (S5), which is hard to replicate due to its subjectivity. Whilst some of the studies have attempted to address this problem by making their datasets available, the code and methods used to create these datasets have been rarely shared; this situation has made the reproduction and adaptation efforts very difficult. For example, Riom and colleagues [72] found the replication of a seminal work [P16] infeasible due to these challenges.

We observe that the underlying problem of the above-mentioned challenges is a lack of proper data reporting and storage efforts. Hence, it can be recommended that researchers make more efforts in the future to specify the exact details and processes of data preparation in order to improve data reuse and augmentation. Gebru and colleagues [73] propose a process called *Datasheets for Datasets*. Such processes of dataset documentation promises to also help practitioners to better understand the capabilities and implications of the produced SVP models.

Another potential solution is open data sharing platforms or repositories for enhancing data availability. Such platforms can also provide vital information regarding data provenance and quality. These efforts will potentially also encourage data maintenance, which is another issue that we have discussed in R2. Considerable efforts towards such a platform have already been made in the Software Engineering domain through the PROMISE repository [74] and the SEACRAFT repository [75]. However, these efforts still currently fail to ensure data provenance and maintenance.

Ultimately there is a need for better industrial and corporate engagement with public datasets. Direct industrial collaboration often leads to higher impact software engineering research [76]. Such involvement will provide more realistic, complete and trustworthy datasets, due to increased data provenance and accessibility. This cooperation would also allow SVP evaluation to extend beyond

open-source repositories to private-source code. As technical advancements continue to be made for SVP models, researchers should move towards industrial demonstrations and case studies that can encourage further corporate participation and application. Garousi *et al.* [76] have outlined several best practices towards ensuring effective industry-academia collaboration.

## 8 THREATS TO VALIDITY

This review has been designed and executed by carefully following the guidelines for SLR provided by Kitchenham *et al.* [26]. Hence, we have identified the potential validity threats to this SLR and taken appropriate steps to minimize the potential impact of the identified threats as per the SLR guidelines. We discuss the validity threats considered for this SLR below.

A standard threat to any SLR is selection bias; some relevant papers may be missed during the selection process of the SLR. Whilst there is a possibility that our search and study selection process may have missed some relevant papers, we systematically drove the paper selection process by following the SLR guidelines and the recommended practices to minimize such possibility of missing the relevant papers. For example, we chose a meta search engine, SCOPUS, that indexes all the well-known computer science and software engineering digital libraries such as IEEE, ACM, Elsevier, and Springer. Furthermore, we iteratively refined our search string using the quasi-gold sensitivity approach defined by Zhang *et al.* [27] until we were confident that our search string had retrieved the majority of the key papers in this research area. Finally, we used backward and forward snowballing to help capture any studies that might have been missed during our automatic search. To avoid the study selection bias by the authors, we initially conducted a collaborative pilot study selection on 100 papers to ensure consistency. Furthermore, any paper that an individual author was not confident about including/excluding was discussed between the first two authors before making a final decision.

Additional validity threats can be introduced through the quality assessment, data extraction and thematic analysis processes of this SLR. Inaccuracies in these processes can be introduced by human-error and researcher-bias. The first two authors jointly carried out the pilot activities for these processes to help ensure consistency. Furthermore, all the data-extraction activities were cross-checked by the authors, with disagreements resolved through discussions.

Finally, our results may be affected by publication bias; research is biased towards the publication of positive results over negative results [26]. Hence, it is expected that researchers would also be reluctant to report the major (data) limitations if it is not the focus of a study. As our findings are grounded in the data extracted from the primary studies, we are only able to report the considerations and challenges explicitly discussed in the papers. Hence, our findings may not be exhaustive. However, we assume that our findings are able to capture the major data challenges and considerations, due to the quality and integrity of our selected primary studies which we ensured through rigorous use of our inclusion/exclusion and data quality assessment criteria from Section 3.2.

## 9 CONCLUSION

Software Vulnerability Prediction (SVP) approaches have gained significant attention of the software engineering community for ensuring software security. Whilst researchers have developed and evaluated a large number of SVP approaches, they have also identified several challenges that need to be addressed for wider adoption of the reported approaches. Given the importance of the topic and the amount of available literature that is largely dispersed, it was needed and timely to invest efforts in a Systemization of Knowledge (SoK) of the available peer-reviewed literature in order to highlight the key challenges of data preparation phases for developing, evaluating, and deploying SVP approaches.

We have carried out a systematic literature review of 61 papers by following the well-known SLR guidelines. The main aim of our study was to identify the considerations, challenges and solutions of data preparation for SVP, by answering three research questions. First, we have identified the major decisions made by researchers for the data preparation processes, to help inform the state of the practice. Second, through our thematic analysis, we have derived a taxonomy of 16 identified data challenges that researchers face for constructing data-driven methods for SVP. These challenges involve the data Generalizability, Accessibility, label collection Effort and Scarcity, and both Label and Data Noise. Third, we also categorized and mapped the data preparation related solutions reported in the reviewed papers in order to highlight their understanding and utility.

We have found that the identified challenges are particularly pertinent to the data labelling process that has consequently attracted the majority of the identified solutions. Due to these significant challenges, data reuse is a common solution to reduce data construction effort and difficulties. Alternatively, the use of synthetic datasets is an attractive option as these datasets artificially solve challenges for Accessibility, Effort, Data Scarcity and Label Noise, but these datasets experience significant challenges with Generalizability, which severely limits their value.

We assert that the findings of our SLR will help researchers and practitioners to understand the key SVP data preparation considerations and challenges. By consolidating the state of the practice into an integrated source of information, SoK, this study is expected to assist practitioners in improving their data preparation practices for building and deploying SVP models. Furthermore, we believe that the reported taxonomy of the data preparation challenges will be used to identify and classify the data preparation for SVP related challenges that practitioners may encounter; and our categorization of the identified solutions is expected to help identify the best practices of data preparation for SVP models. Such improvements are expected to improve the quality of SVP models, as their efficacy hinges on the data quality; *Garbage In, Garbage Out*.

The findings of this review will inform the future research that is expected to address the challenges for which the reviewed studies do not provide appropriate solutions. We have derived six recommendations from the findings of this review for providing actionable suggestions to SVP



researchers and practitioners. Whilst we acknowledge that our study does not provide a complete overview of the SVP process on its own, the data is undoubtedly one of the most important components for any data-driven process; hence, we take one of the first steps to highlight this area where the future research efforts can make significant advances in the data preparation and data quality state-of-the-art. Such advances will enable the creation and use of more reliable and trustworthy SVP approaches supporting automated software security analytics.

## REFERENCES

- [1] R. Sobers, 134 cybersecurity statistics and trends for 2021, 2021. [Online]. Available: <https://www.varonis.com/blog/cybersecurity-statistics/>
- [2] H. Shahriar and M. Zulkernine, "Mitigating program security vulnerabilities: Approaches and challenges," *ACM Comput. Surv.*, vol. 44, no. 3, pp. 1–46, 2012.
- [3] Y. Shin and L. Williams, "Can traditional fault prediction models be used for vulnerability prediction?," *Empir. Softw. Eng.*, vol. 18, no. 1, pp. 25–59, 2013.
- [4] H. Hanif, M. H. N. M. Nasir, M. F. Ab Razak, A. Firdaus, and N. B. Anuar, "The rise of software vulnerability: Taxonomy of software vulnerabilities detection and machine learning approaches," *J. Netw. Comput. Appl.*, vol. 179, 2021, Art. no. 103009.
- [5] D. Pyle, *Data Preparation for Data Mining*. Morgan Kaufmann, Amsterdam, The Netherlands: Elsevier, 1999.
- [6] A. Zheng and A. Casari, *Feature Engineering for Machine Learning: Principles and Techniques for Data Scientists*, Newton, MA, USA: O'Reilly Media, Inc., 2018.
- [7] J. Walden, J. Stuckman, and R. Scandariato, "Predicting vulnerable components: Software metrics vs text mining," in *Proc. IEEE 25th Int. Symp. Softw. Rel. Eng.*, 2014, pp. 23–33.
- [8] T. Zimmermann, N. Nagappan, and L. Williams, "Searching for a needle in a haystack: Predicting security vulnerabilities for windows vista," in *Proc. 3rd Int. Conf. Softw. Testing, Verification Validation*, 2010, pp. 421–428.
- [9] A. Anwar, A. Abusnaina, S. Chen, F. Li, and D. Mohaisen, "Cleaning the NVD: Comprehensive quality assessment, improvements, and analyses," 2020, *arXiv: 2006.15074*.
- [10] R. Coulter, Q.-L. Han, L. Pan, J. Zhang, and Y. Xiang, "Code analysis for intelligent cyber systems: A data-driven approach," *Inf. Sci.*, vol. 524, pp. 46–58, 2020.
- [11] M. Jimenez, R. Rwemalika, M. Papadakis, F. Sarro, Y. Le Traon, and M. Harman, "The importance of accounting for real-world labelling when predicting software vulnerabilities," in *Proc. 27th ACM Joint Meeting Eur. Softw. Eng. Conf. Symp. Found. Softw. Eng.*, 2019, pp. 695–705.
- [12] C. Tantithamthavorn, S. McIntosh, A. E. Hassan, A. Ihara, and K. Matsumoto, "The impact of mislabelling on the performance and interpretation of defect prediction models," in *Proc. IEEE/ACM 37th IEEE Int. Conf. Softw. Eng.*, 2015, pp. 812–823.
- [13] G. Lin, S. Wen, Q.-L. Han, J. Zhang, and Y. Xiang, "Software vulnerability detection using deep neural networks: A survey," *Proc. IEEE Proc.*, vol. 108, no. 10, pp. 1825–1848, Oct. 2020.
- [14] B. Turhan, T. Menzies, A. B. Bener, and J. Di Stefano, "On the relative value of cross-company and within-company data for defect prediction," *Empirical Softw. Eng.*, vol. 14, no. 5, pp. 540–578, 2009.
- [15] P. Morrison, K. Herzig, B. Murphy, and L. Williams, "Challenges with applying vulnerability prediction models," in *Proc. Symp. Bootcamp Sci. Secur.*, 2015, pp. 1–9.
- [16] P. Zeng, G. Lin, L. Pan, Y. Tai, and J. Zhang, "Software vulnerability analysis and discovery using deep learning techniques: A survey," *IEEE Access*, vol. 8, pp. 197158–197172, 2020.
- [17] A. O. A. Semasaba, W. Zheng, X. Wu, and S. A. Agyemang, "Literature survey of deep learning-based vulnerability analysis on source code," *IET Softw.*, vol. 14, no. 6, pp. 654–664, 2020.
- [18] S. Amershi *et al.*, "Software engineering for machine learning: A case study," in *Proc. IEEE/ACM 41st Int. Conf. Softw. Eng.: Softw. Eng. Pract.*, 2019, pp. 291–300.
- [19] "The state of AI and machine learning," Appen, Sydney, Australia, Tech. Rep. WC-2019, 2019.
- [20] S. Neuhaus, T. Zimmermann, C. Holler, and A. Zeller, "Predicting vulnerable software components," in *Proc. 14th ACM Conf. Comput. Commun. Secur.*, 2007, pp. 529–540.
- [21] R. Croft, D. Newlands, Z. Chen, and M. A. Babar, "An empirical study of rule-based and learning-based approaches for static application security testing," 2021, *arXiv:2107.01921*.
- [22] S. M. Ghaffarian and H. R. Shahriari, "Software vulnerability analysis and discovery using machine-learning and data-mining techniques: A survey," *ACM Comput. Surv.*, vol. 50, no. 4, pp. 1–36, 2017.
- [23] Z. Li and Y. Shao, "A survey of feature selection for vulnerability prediction using feature-based machine learning," in *Proc. 11th Int. Conf. Mach. Learn. Comput.*, 2019, pp. 36–42.
- [24] T. H. Le, H. Chen, and M. A. Babar, "A survey on data-driven software vulnerability assessment and prioritization," 2021, *arXiv:2107.08364*.
- [25] S. K. Singh and A. Chaturvedi, "Applying deep learning for discovery and analysis of software vulnerabilities: A brief survey," *Soft Comput.: Theories Appl.*, vol. 8, pp. 649–658, 2020.
- [26] B. Kitchenham, "Procedures for performing systematic reviews," *Keele, U.K., Keele Univ.*, vol. 33, no. 2004, pp. 1–26, 2004.
- [27] H. Zhang, M. A. Babar, and P. Tell, "Identifying relevant studies in software engineering," *Inf. Softw. Technol.*, vol. 53, no. 6, pp. 625–637, 2011.
- [28] C. Schardt, M. B. Adams, T. Owens, S. Keitz, and P. Fontelo, "Utilization of the pico framework to improve searching pubmed for clinical questions," *BMC Med. Informat. Decis. Mak.*, vol. 7, no. 1, pp. 1–6, 2007.
- [29] S. Hosseini, B. Turhan, and D. Gunarathna, "A systematic literature review and meta-analysis on cross project defect prediction," *IEEE Trans. Softw. Eng.*, vol. 45, no. 2, pp. 111–147, Feb. 2019.
- [30] N. Li, M. Shepperd, and Y. Guo, "A systematic review of unsupervised learning techniques for software defect prediction," *Inf. Softw. Technol.*, vol. 122, 2020, Art. no. 106287.
- [31] B. Sabir, F. Ullah, M. A. Babar, and R. Gaire, "Machine learning for detecting data exfiltration: A review," *ACM Comput. Surv.*, vol. 54, no. 3, pp. 1–47, 2021.
- [32] T. Hall, S. Beecham, D. Bowes, D. Gray, and S. Counsell, "A systematic literature review on fault prediction performance in software engineering," *IEEE Trans. Softw. Eng.*, vol. 38, no. 6, pp. 1276–1304, Nov./Dec. 2012.
- [33] C. Wohlin, "Guidelines for snowballing in systematic literature studies and a replication in software engineering," in *Proc. 18th Int. Conf. Eval. Assessment Softw. Eng.*, 2014, pp. 1–10.
- [34] V. Garousi and M. Felderer, "Experience-based guidelines for effective and efficient data extraction in systematic reviews in software engineering," in *Proc. 21st Int. Conf. Eval. Assessment Softw. Eng.*, 2017, pp. 170–179.
- [35] V. Braun and V. Clarke, "Using thematic analysis in psychology," *Qualitative Res. Psychol.*, vol. 3, no. 2, pp. 77–101, 2006.
- [36] QSR International, Nvivo, 2020. [Online]. Available: <https://www.qsrinternational.com/nvivo-qualitative-data-analysis-software/home>
- [37] L. L. Pipino, Y. W. Lee, and R. Y. Wang, "Data quality assessment," *Commun. ACM*, vol. 45, no. 4, pp. 211–218, 2002.
- [38] F. Sidi, P. H. S. Panahy, L. S. Affendey, M. A. Jabar, H. Ibrahim, and A. Mustapha, "Data quality: A survey of data quality dimensions," in *Proc. Int. Conf. Inf. Retrieval Knowl. Manage.*, 2012, pp. 300–304.
- [39] Github, 2008. [Online]. Available: <https://github.com/>
- [40] S. Chakraborty, R. Krishna, Y. Ding, and B. Ray, "Deep learning based vulnerability detection: Are we there yet?," *IEEE Trans. Softw. Eng.*, early access, Jun. 08, 2021, doi: [10.1109/TSE.2021.3087402](https://doi.org/10.1109/TSE.2021.3087402).
- [41] NIST, Software assurance and reference dataset (SARD), 2005. [Online]. Available: <https://samate.nist.gov/SARD/testsuite.php>
- [42] OWASP, Owasp benchmark project, 2015. [Online]. Available: <https://owasp.org/www-project-benchmark/>
- [43] VulnSpy, Sqli labs, 2018. [Online]. Available: <https://www.vulnspy.com/sqli-labs/>
- [44] M. Ivarsson and T. Gorschek, "A method for evaluating rigor and industrial relevance of technology evaluations," *Empirical Softw. Eng.*, vol. 16, no. 3, pp. 365–395, 2011.
- [45] NIST, National vulnerability database, 1999. [Online]. Available: <https://nvd.nist.gov/>
- [46] Atlassian, Jira, 2002. [Online]. Available: <https://www.atlassian.com/software/jira>
- [47] Mozilla, Bugzilla, 1998. [Online]. Available: <https://www.bugzilla.org/>

- [48] J. Hernández-González, I. Inza, and J. A. Lozano, "Weak supervision and other non-standard classification problems: A taxonomy," *Pattern Recognit. Lett.*, vol. 69, pp. 49–55, 2016.
- [49] P. Domingos, "A few useful things to know about machine learning," *Commun. ACM*, vol. 55, no. 10, pp. 78–87, 2012.
- [50] P. S. Kochhar, D. Wijedasa, and D. Lo, "A large scale study of multiple programming languages and code quality," in *Proc. IEEE 23rd Int. Conf. Softw. Anal., Evol., Reengineering*, 2016, pp. 563–573.
- [51] B. Lika, K. Kolomvatsos, and S. Hadjiefthymiades, "Facing the cold start problem in recommender systems," *Expert Syst. Appl.*, vol. 41, no. 4, pp. 2065–2073, 2014.
- [52] M. F. Bosu and S. G. MacDonell, "A taxonomy of data quality challenges in empirical software engineering," in *Proc. 22nd Australian Softw. Eng. Conf.*, 2013, pp. 97–106.
- [53] M. Dowd, J. McDonald, and J. Schuh, *The Art of Software Security Assessment: Identifying and Preventing Software Vulnerabilities*, London, U.K.: Pearson Education, 2006.
- [54] S. Barnum and G. McGraw, "Knowledge for software security," *IEEE Secur. Privacy*, vol. 3, no. 2, pp. 74–78, Mar./Apr. 2005.
- [55] R. Y. Wang and D. M. Strong, "Beyond accuracy: What data quality means to data consumers," *J. Manage. Inf. Syst.*, vol. 12, no. 4, pp. 5–33, 1996.
- [56] H. He and E. A. Garcia, "Learning from imbalanced data," *IEEE Trans. Knowl. Data Eng.*, vol. 21, no. 9, pp. 1263–1284, Sep. 2009.
- [57] K. Herzig, S. Just, and A. Zeller, "The impact of tangled code changes on defect prediction models," *Empirical Softw. Eng.*, vol. 21, no. 2, pp. 303–336, 2016.
- [58] J. Zhou, H. Zhang, and D. Lo, "Where should the bugs be fixed? more accurate information retrieval-based bug localization based on bug reports," in *Proc. 34th Int. Conf. Softw. Eng.*, 2012, pp. 14–24.
- [59] Y. Jiang, P. Lu, X. Su, and T. Wang, "Ltrwes: A new framework for security bug report detection," *Inf. Softw. Technol.*, vol. 124, 2020, Art. no. 106314.
- [60] Y. Jia and M. Harman, "An analysis and survey of the development of mutation testing," *IEEE Trans. Softw. Eng.*, vol. 37, no. 5, pp. 649–678, Sep./Oct. 2010.
- [61] N. Leicht, I. Blohm, and J. M. Leimeister, "Leveraging the power of the crowd for software testing," *IEEE Softw.*, vol. 34, no. 2, pp. 62–69, Mar./Apr. 2017.
- [62] Y. Zheng *et al.*, "D2a: A dataset built for ai-based vulnerability detection methods using differential analysis," in *Proc. IEEE/ACM 43rd Int. Conf. Softw. Eng.: Softw. Eng. Pract.*, 2021, pp. 111–120.
- [63] D. Tsipras, S. Santurkar, L. Engstrom, A. Ilyas, and A. Madry, "From imagenet to image classification: Contextualizing progress on benchmarks," in *Proc. Int. Conf. Mach. Learn.*, 2020, pp. 9625–9635.
- [64] G. M. Weinberg, *Perfect Software and Other Illusions About Testing*. New York, NY, USA: Dorset House Pub., 2008.
- [65] T. H. M. Le, D. Hin, R. Croft, and M. A. Babar, "Puminer: Mining security posts from developer question and answer websites with PU learning," in *Proc. 17th Int. Conf. Mining Softw. Repositories*, 2020, pp. 350–361.
- [66] T. H. M. Le, B. Sabir, and M. A. Babar, "Automated software vulnerability assessment with concept drift," in *Proc. IEEE/ACM 16th Int. Conf. Mining Softw. Repositories*, 2019, pp. 371–382.
- [67] S. McIntosh and Y. Kamei, "Are fix-inducing changes a moving target? a longitudinal case study of just-in-time defect prediction," *IEEE Trans. Softw. Eng.*, vol. 44, no. 5, pp. 412–428, May 2018.
- [68] D. Falessi, J. Huang, L. Narayana, J. F. Thai, and B. Turhan, "On the need of preserving order of data when validating within-project defect classifiers," *Empirical Softw. Eng.*, vol. 25, no. 6, pp. 4805–4830, 2020.
- [69] T. M. Mitchell, "Machine learning and data mining," *Commun. ACM*, vol. 42, no. 11, pp. 30–36, 1999.
- [70] Z. Mushtaq, G. Rasool, and B. Shehzad, "Multilingual source code analysis: A systematic literature review," *IEEE Access*, vol. 5, pp. 11307–11336, 2017.
- [71] R. Zhang, M. Indulska, and S. Sadiq, "Discovering data quality problems," *Bus. Inf. Syst. Eng.*, vol. 61, no. 5, pp. 575–593, 2019.
- [72] T. Riom, A. Sawadogo, K. Allix, T. F. Bissyandé, N. Moha, and J. Klein, "Revisiting the VCCFinder approach for the identification of vulnerability-contributing commits," *Empirical Softw. Eng.*, vol. 26, no. 3, pp. 1–30, 2021.
- [73] T. Gebru *et al.*, "Datasheets for datasets," *Commun. ACM*, vol. 64, no. 12, pp. 86–92, Dec. 2021.
- [74] J. S. Shirabad and T. Menzies, "The PROMISE repository of software engineering databases," School of Information Technology and Engineering, University of Ottawa, Canada, 2005. [Online]. Available: <http://promise.site.uottawa.ca/SERepository>
- [75] T. Menzies, R. Krishna, and D. Pryor, "The seacraft repository of empirical software engineering data," 2017. [Online]. Available: <https://zenodo.org/communities/seacraft>
- [76] V. Garousi, K. Petersen, and B. Ozkan, "Challenges and best practices in industry-academia collaborations in software engineering: A systematic literature review," *Inf. Softw. Technol.*, vol. 79, pp. 106–127, 2016.

► For more information on this or any other computing topic, please visit our Digital Library at [www.computer.org/csdl](http://www.computer.org/csdl).