



A survey on machine learning techniques applied to source code^{☆,☆☆}

Tushar Sharma ^{a,*}, Maria Kechagia ^b, Stefanos Georgiou ^c, Rohit Tiwari ^d, Indira Vats ^e, Hadi Moazen ^f, Federica Sarro ^b

^a Dalhousie University, Canada

^b University College London, United Kingdom

^c Queens University, Canada

^d DevOn, India

^e J.S.S. Academy of Technical Education, India

^f Sharif University of Technology, Iran



ARTICLE INFO

Dataset link: <https://github.com/tushartushar/ML4SCA>

Keywords:

Machine learning for software engineering

Source code analysis

Deep learning

Datasets

Tools

ABSTRACT

The advancements in machine learning techniques have encouraged researchers to apply these techniques to a myriad of software engineering tasks that use source code analysis, such as testing and vulnerability detection. Such a large number of studies hinders the community from understanding the current research landscape. This paper aims to summarize the current knowledge in applied machine learning for source code analysis. We review studies belonging to twelve categories of software engineering tasks and corresponding machine learning techniques, tools, and datasets that have been applied to solve them. To do so, we conducted an extensive literature search and identified 494 studies. We summarize our observations and findings with the help of the identified studies. Our findings suggest that the use of machine learning techniques for source code analysis tasks is consistently increasing. We synthesize commonly used steps and the overall workflow for each task and summarize machine learning techniques employed. We identify a comprehensive list of available datasets and tools useable in this context. Finally, the paper discusses perceived challenges in this area, including the availability of standard datasets, reproducibility and replicability, and hardware resources.

Editor's note: Open Science material was validated by the Journal of Systems and Software Open Science Board.

1. Introduction

In the last two decades, we have witnessed significant advancements in Machine Learning (ML), including Deep Learning (DL) techniques, specifically in the domain of image (Krizhevsky et al., 2012; Szegedy et al., 2015), text (Lee et al., 2017; Abdeljaber et al., 2017), and speech (Sainath et al., 2015; Greff et al., 2017; Graves et al., 2013) processing. These advancements, coupled with a large amount of open-source code and associated artifacts, as well as the availability of accelerated hardware, have encouraged researchers and practitioners to use ML techniques to address software engineering problems (Wan et al., 2019; Zhang and Tsai, 2003; Allamanis et al., 2018a; Le et al., 2020; Alsolai and Roper, 2020).

The software engineering community has employed ML and DL techniques for a variety of applications such as software testing (Lima et al., 2020; Omri and Sinz, 2020; Zhang et al., 2020), source code representation (Allamanis et al., 2018a; Hellendoorn and Devanbu,

2017), source code quality analysis (Alsolai and Roper, 2020; Azeem et al., 2019), program synthesis (Le et al., 2020; Yahav, 2018), code completion (Liu et al., 2020d), refactoring (Aniche et al., 2020), code summarization (Liu et al., 2018; LeClair et al., 2019; Allamanis et al., 2015a), and vulnerability analysis (Shen and Chen, 2020; Shabtai et al., 2009; Ucci et al., 2019) that involve source code analysis. As the field of *Machine Learning for Software Engineering* (ML4SE) is expanding, the number of available resources, methods, and techniques as well as tools and datasets, is also increasing. This poses a challenge, to both researchers and practitioners, to fully comprehend the landscape of the available resources and infer the potential directions that the field is taking. In this context, literature surveys play an important role in understanding existing research, finding gaps in research or practice, and exploring opportunities to improve the state of the art. By systematically examining existing literature, surveys may uncover hidden patterns, recurring themes, and promising research directions. Surveys

[☆] Editor: Dr Alexander Chatzigeorgiou.

^{☆☆} Funding: Maria Kechagia and Federica Sarro are supported by the ERC grant no. 741278 (EPIC).

* Corresponding author.

E-mail address: tushar@dal.ca (T. Sharma).

Table 1

Comparison among surveys. The “Category” column refers to the software engineering task the survey covers. The “Scope” column indicates the focus of the study; TML refers to traditional machine learning and DL refers to deep learning techniques. The “Data&Tools” column indicates if a survey reviews available datasets and tools for ML-based applications, the “Challenges” column shows whether the study identifies challenges in the field studied, the “Type” column refers to the type of literature survey, and the “#Studies” column refers to the number of studies included in a given survey. We use “–” to indicate that a field is not applicable to a certain study and NA for the number of studies column, where the study does not explicitly mention selection criteria and the number of selected studies.

Category	Article	Scope	Data & Tools	Challenges	Type	#Studies
Program comprehension	Nazar et al. (2016)	TML	Tools	No	Lit. survey	59
	Zhang et al. (2022)	DL	Data	No	Lit. survey	NA
	Song et al. (2019b)	TML & DL	No	Yes	Lit. survey	NA
Testing	Omri and Sinz (2020)	DL	No	No	Lit. survey	NA
	Durelli et al. (2019)	TML & DL	No	Yes	Mapping study	48
	Hall and Bowes (2012)	TML	Yes	Yes	Meta-analysis	21
	Zhang et al. (2020)	TML & DL	No	Yes	Lit. survey	46
	Pandey et al. (2021)	TML	No	Yes	Lit. survey	154
Vulnerability analysis	Singh et al. (2018)	TML	No	No	Lit. survey	13
	Li et al. (2019c)	DL	Yes	Yes	Meta-analysis	–
	Shen and Chen (2020)	DL	No	Yes	Meta-analysis	–
	Ucci et al. (2019)	TML	No	Yes	Lit. survey	64
	Jie et al. (2016)	TML	No	No	Lit. survey	19
Quality assessment	Hanif et al. (2021)	TML & DL	No	Yes	Lit. survey	90
	Alsolai and Roper (2020)	TML	No	No	Lit. survey	56
	Tsintzira et al. (2020)	TML	Yes	Yes	Lit. survey	90
	Azeem et al. (2019)	TML	Yes	No	Lit. survey	15
	Caram et al. (2019a)	TML	No	No	Mapping study	25
Prog. synthesis	Lewowski and Madeyski (2022)	TML	Yes	No	Lit. survey	45
	Goues et al. (2019)	TML & DL	No	Yes	Lit. survey	NA
	Le et al. (2020)	DL	Yes	Yes	Lit. survey	NA
Prog. synthesis & code representation	Allamanis et al. (2018a)	TML & DL	Yes	Yes	Lit. survey	39 + 48
	Yang et al. (2022)	DL	Data	Yes	Lit. survey	250
Source-code analysis	Our study	TML & DL	Yes	Yes	Lit. survey	494

also identify untapped opportunities and formulation of new hypotheses. A survey also serves as an educational tool, offering comprehensive coverage of the field to a newcomer.

In fact, there have been numerous recent attempts to summarize the application-specific knowledge in the form of surveys. For example, Allamanis et al. (2018a) present key methods to model source code using ML techniques. Shen and Chen (2020) provide a summary of research methods associated with software vulnerability detection, software program repair, and software defect prediction. Durelli et al. (2019) collect 48 primary studies focusing on software testing using machine learning. Alsolai and Roper (2020) present a systematic review of 56 studies related to maintainability prediction using ML techniques. Recent surveys (Tsintzira et al., 2020; AL-Shaaby et al., 2020; Azeem et al., 2019) summarize application of ML techniques on software code smells and technical debt identification. Similarly, literature reviews on program synthesis (Le et al., 2020) and code summarization (Nazar et al., 2016) have been attempted. We compare in Table 1 the aspects investigated in our survey with respect to existing surveys that review ML techniques for topics such as testing, vulnerabilities, and program comprehension with our survey. Existing studies, in general, kept their focus on only one category; due to that readers could not grasp existing literature belonging to various software engineering categories in a consistent form. In addition, existing surveys do not always provide datasets and tools in the field. Our survey, covers a wide range of software engineering activities; it summarizes a significantly large number of studies; it systematically examines available tools and datasets for ML that would support researchers in their studies in this field; it identifies perceived challenges in the field to encourage the community to explore ways to overcome them.

In this paper, we focus on the usage of ML, including DL, techniques for source code analysis. Source code analysis involves tasks that take the source code as input, process it, and/or produce source code as output. Source code representation, code quality analysis, testing, code summarization, and program synthesis are applications that involve

source code analysis. To the best of our knowledge, the software engineering literature lacks a survey covering a wide range of source code analysis applications using machine learning; this work is an attempt to fill this research gap.

In this survey, we aim to give a comprehensive, yet concise, overview of current knowledge on applied machine learning for source code analysis. We also aim to collate and consolidate available resources (in the form of datasets and tools) that researchers have used in previous studies on this topic. Additionally, we aim to identify and present challenges in this domain. We believe that our efforts to consolidate and summarize the techniques, resources, and challenges will help the community to not only understand the state-of-the-art better, but also to focus their efforts on tackling the identified challenges.

This survey makes the following contributions to the field:

- It presents a summary of the applied machine learning studies attempted in the source code analysis domain.
- It consolidates resources (such as datasets and tools) relevant for future studies in this domain.
- It provides a consolidated summary of the open challenges that require the attention of the researchers.

The rest of the paper is organized as follows. We present the followed methodology, including the literature search protocol and research questions, in Section 2. Sections 3, 4, and 5 provide the detailed results of our findings. We present threats to validity in Section 6, and conclude the paper in Section 7.

2. Methodology

First, we present the objectives of this study and the research questions derived from such objectives. Second, we describe the search protocol we followed to identify relevant studies. The protocol identifies detailed steps to collect the initial set of articles as well as the inclusion and exclusion criteria to obtain a filtered set of studies.

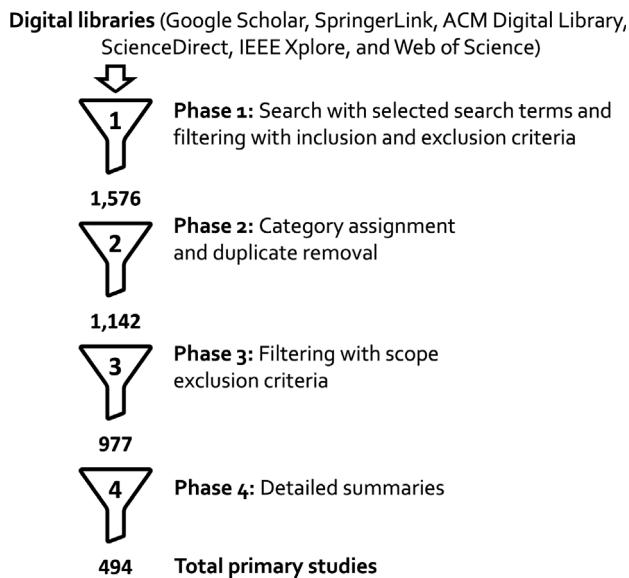


Fig. 1. Overview of the search process.

2.1. Research objectives

This study aims to achieve the following research objectives (ROs).

- RO1. *Identifying specific software engineering tasks involving source code that have been attempted using machine learning.*
Our objective is to explore the extent to which machine learning has been applied to analyze and process source code for SE tasks. We aim to summarize how ML can help engineers tackle specific SE tasks.
- RO2. *Summarizing the machine learning techniques used for these tasks.*
This objective explores the ML techniques commonly applied to source code for performing the software engineering tasks identified above. We attempt to synthesize a mapping of tasks (along with related sub-tasks) and corresponding ML techniques.
- RO3. *Providing a list of available datasets and tools.*
With this goal, we aim to provide a consolidated summary of publicly available datasets and tools along with their purpose.
- RO4. *Identifying the challenges and perceived deficiencies in ML-enabled source code analysis and manipulation for software engineering.*
With this objective, we aim to identify challenges, and opportunities arising when applying ML techniques to source code for SE tasks, as well as to understand the extent to which they have been addressed in the articles surveyed.

2.2. Literature search protocol

We identified 494 relevant studies through a four step literature search. Fig. 1 summarizes the search process. We elaborate on each of these phases in the rest of this section.

2.2.1. Literature search—Phase 1

We split the phase 1 literature search into two rounds. In the first round, we carried out an extensive initial search on six well-known digital libraries—Google Scholar, SpringerLink, ACM Digital Library, ScienceDirect, IEEE Xplore, and Web of Science during Feb–Mar 2021. We formulated a set of search terms based on common tasks and software engineering activities related to source code analysis. Specifically, we used the following terms for the search: *machine learning code*, *machine learning code representation*, *machine learning testing*, *machine learning code synthesis*, *machine learning smell identification*,

machine learning security, *source code analysis*, *machine learning software quality assessment*, *machine learning code summarization*, *machine learning program repair*, *machine learning code completion*, and *machine learning refactoring*. We searched minimum seven pages of search results for each search term manually; beyond seven pages, we continued the search unless we get two continuous search pages without any new and relevant articles. We adopted this mechanism to avoid missing any relevant articles in the context of our study.

In the second round of phase 1, we identified a set of frequently occurring keywords in the articles obtained from the first round for each category individually. To do that, we manually scanned the keywords mentioned in the articles belonging to each category, and noted the keywords that appeared at least three times. If the selected keywords are too generic, we first check whether adding *machine learning* would improve the search results. For example, *machine learning* and *program generation* occurred multiple times in the *program synthesis* category; we combined both of these terms to make one search string i.e., *program generation using machine learning*. In other cases, we tried to reduce the scope of the search term by adding qualifying terms. Consider *feature learning* as an example: it is so generic that would result in many unrelated results. We reduced the search scope by adding *source code* in the search i.e., searching using *feature learning in source code*. We carried out this additional round of literature search to augment our initial search terms and reduce the risk of missing relevant articles. Table 2 summarizes the search terms along with the number of studies found in the second round of phase 1; the full list of used search terms can be found in our replication package (Sharma et al., 2022). Next, we defined inclusion and exclusion criteria to filter out irrelevant articles.

Inclusion criteria:

- Studies and surveys that discuss the application of machine learning (including DL) to source code to perform a software engineering task.
- Resources revealing the deficiencies or challenges in the current set of methods, tools, and practices.

Exclusion criteria:

- Studies focusing on techniques other than ML applied on source code to address software engineering tasks e.g., code smell detection using metrics.
- Articles that are not peer-reviewed (such as articles available only on arXiv.org).
- Articles constituting a keynote, extended abstract, editorial, tutorial, poster, or panel discussion (due to insufficient details and limited length).
- Studies whose full text is not available, or is written in any other language than English.

We considered whether to include studies that do not directly analyze source code. Often, source code is analyzed to extract features, and machine learning techniques are applied to the extracted features. Furthermore, researchers in the field either create their own dataset (in that case, analyze/process source code) or use existing datasets. Removing studies that use a dataset will make this survey incomplete; hence, we decided to include such studies.

During the search, we documented studies that satisfy our search protocol in a spreadsheet including the required meta-data (such as title, bibtex record, and link of the source). The spreadsheet with all the articles from each phase can be found in our online replication package (Sharma et al., 2022). Each selected article went through a manual inspection of title, keywords, and abstract. The inspection applied the inclusion and exclusion criteria leading to inclusion or exclusion of the articles. In the end, we obtained 1576 articles after completing Phase 1 of the search process.

Table 2
Search terms and corresponding relevant studies found in the second round of phase 1.

Category	Search terms	#Studies
Vulnerability analysis	Feature learning in source code	9
	Vulnerability prediction in source code using machine learning	70
	Deep learning-based vulnerability detection	8
	Malicious code detection with machine learning	45
Testing	Word embedding in software testing	2
	Automated Software Testing with machine learning	12
	Optimal machine learning based random test generation	1
Refactoring	Source code refactoring prediction with machine learning	39
	Automatic clone recommendation with machine learning	14
	Machine learning based refactoring detection tools	16
	Search-based refactoring with machine learning	6
Quality assessment	Web service anti-pattern detection with machine learning	25
	Code smell prediction models	34
	Machine learning-based approach for code smells detection	17
	Software design flaw prediction	37
	Linguistic smell detection with machine learning	2
	Software defect prediction with machine learning	66
	Machine learning based software fault prediction	35
Program synthesis	Automated program repair methods with machine learning	45
	Program generation with machine learning	2
	Object-oriented program repair with machine learning	15
	Predicting patch correctness with machine learning	3
	Multihunk program repair with machine learning	9
Program comprehension	Autogenerated code with machine learning	6
	Commits analysis with machine learning	34
	Supplementary bug fixes with machine learning	9
Code summarization	Automatic source code summarization with machine learning	43
	Automatic commit message generation with machine learning	19
	Comments generation with machine learning	11
Code review	Security flaws detection in source code with machine learning	20
	Intelligent source code security review with machine learning	2
Code representation	Design pattern detection with machine learning	10
	Human-machine-comprehensible software representation	1
	Feature learning in source code	6
Code completion	Missing software architectural tactics prediction with machine learning	1
	Software system quality analysis with machine learning	6
	Package-level tactic recommendation generation in source code	3
	Identifier prediction in source code	13
	Token prediction in source code	29

2.2.2. Literature search—Phase 2

We first identified a set of categories and sub-categories for common software engineering tasks. These tasks are commonly referred in recent publications (Ferreira et al., 2021; Allamanis et al., 2018a; Shen and Chen, 2020; Azeem et al., 2019). These categories and sub-categories of common software engineering tasks can be found in Fig. 3. Then, we manually assigned a category and sub-category, if applicable, to each selected article based on the (sub-)category to which an article contributes the most. The assignment was carried out by one of the authors and verified by two other authors. We computed Cohen's Kappa (McHugh, 2012) to measure the initial disagreement; we found a strong agreement among the authors with $\kappa = 0.87$. In case of disagreement, each author specified a key goal, operation, or experiment in the article, indicating the rationale of the category assignment for the article. This exercise resolved the majority of the disagreements. In the rest of the cases, we discussed the rationale identified by individual authors and voted to decide a category or sub-category to which the article contributes the most. In this phase, we also discarded duplicates or irrelevant studies not meeting our inclusion criteria after reading their title and abstract. After this phase, we were left with 1098 studies.

2.2.3. Literature search—Phase 3

In the last decade, the use of ML has increased significantly. The research landscape involving source code and ML, which includes methods, applications, and required resources, has changed significantly in the last decade. To keep the survey focused on recent methods and applications, we focused on studies published after 2011. Also,

we discarded papers that had not received enough attention from the community by filtering out all those having a ‘citation count < (2021 – publication year)’. We chose 2021 as the base year to not penalize studies that came out recently; hence, the studies that are published in 2021 do not need to have any citation to be included in this search. We obtain the citation count from digital libraries manually during Mar-May 2022. After applying this filter, we obtained 977 studies.

2.2.4. Literature search—Phase 4

In this phase, we discarded those studies that do not satisfy our inclusion criteria (such as when the article is too short or do not apply any ML technique to source code for SE tasks) after reading the whole article. The remaining 494 articles are the selected studies that we examine in detail. For each study, we extracted the core idea and contribution, the ML techniques, datasets and tools used as well as challenges and findings unveiled. Next, we present our observations corresponding to each research goal we pose.

2.3. Assigning articles to software engineering task categories

Towards achieving RO1, we tagged each selected article with one of the task categories based on the primary focus of the study. The categories represent common software engineering tasks that involve source code analysis. These categories are *code completion*, *code representation*, *code review*, *code search*, *dataset mining*, *program comprehension*, *program synthesis*, *quality assessment*, *refactoring*, *testing*, and

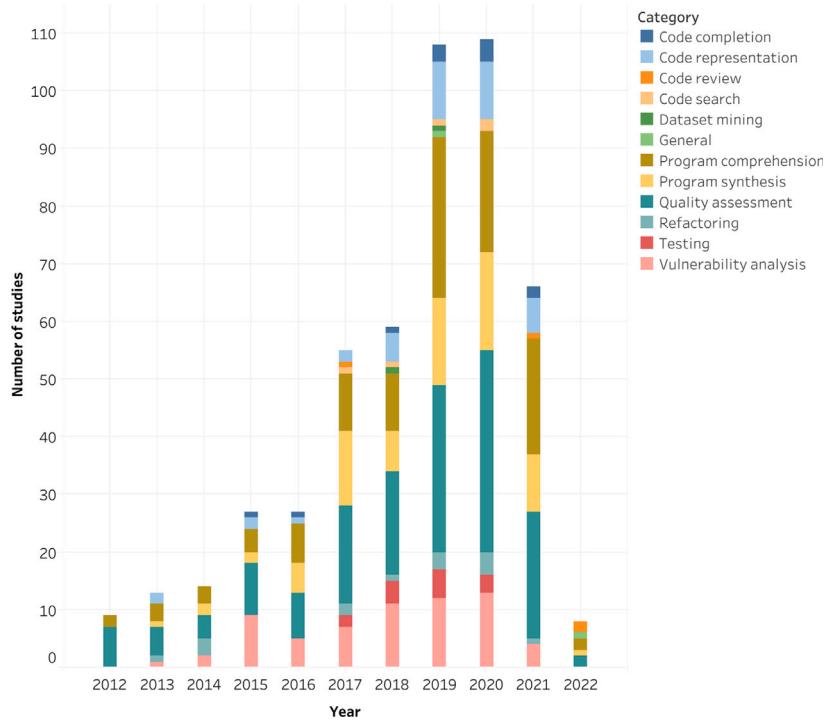


Fig. 2. Category-wise distribution of studies.

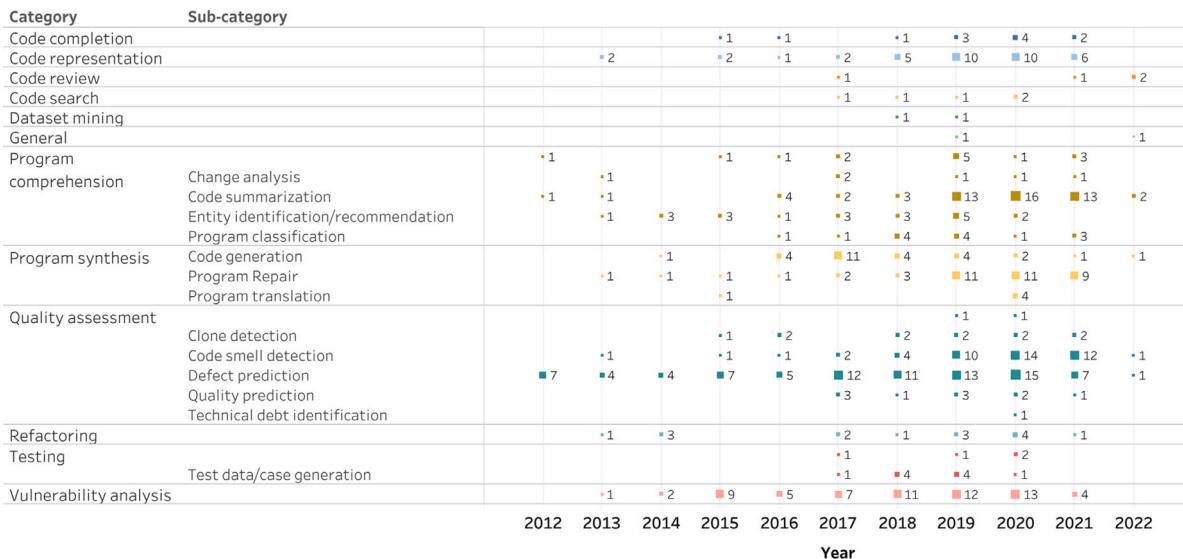


Fig. 3. Category- and sub-categories-wise distribution of studies.

vulnerability analysis. If a given article does not fall in any of these categories but is still relevant to our discussion as it offers overarching discussion on the topic; we put the study in the *general* category. Fig. 2 presents a category-wise distribution of studies per year. It is evident that the topic is engaging the research community more and more and we observe, in general, a healthy upward trend. Interestingly, the number of studies in the scope dropped significantly in the year 2021.

Some of the categories are quite generic and hence further categorization is possible based on specific tasks. For each category, we identified sub-categories by grouping related studies together and assigning an intuitive name representing the set of the studies. For

example, the *testing* category is further divided into *defect prediction*, and *test data/case generation*. We attempted to assign a sub-category to each study; if none of the sub-categories was appropriate for a study, we did not assign any sub-category to the study. One author of this paper assigned a sub-category to each study based on the topic to which that study contributed the most. The initial assignment was verified by two other authors of this paper, where disagreements were discussed and resolved to reach a consensus. Fig. 3 presents the distribution of studies per year w.r.t. each category and corresponding sub-categories.

To quantify the growth of each category, we compute the average increase in the number of articles from the last year for each category

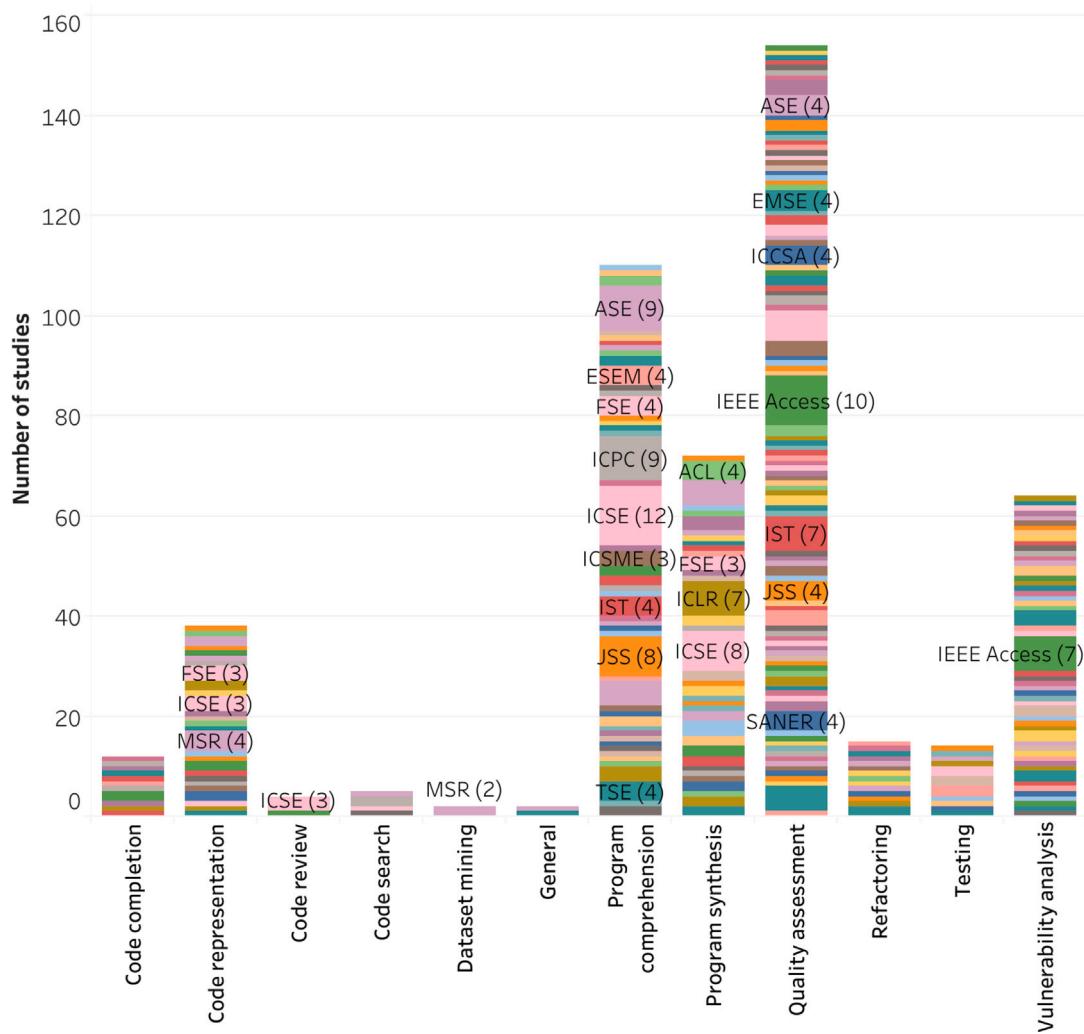


Fig. 5. Top venues for each considered category.

Table 5
Popular models proposed in the selected studies.

Model	#Citations	Model	#Citations
Transfer Naive Bayes (Ma et al., 2012)	513	Code Generation Model (Yin and Neubig, 2017)	651
Path-based code representation (Alon et al., 2018)	230	Multi-headed pointer network (Vasic et al., 2019)	128
Inst2Vec (Ben-Nun et al., 2018)	234	Code-NN (Iyer et al., 2016)	681
DeepCoder (Balog et al., 2016)	612	ASTNN (Zhang et al., 2019)	498
Code2Seq (Alon et al., 2019a)	643	Code2Vec (Alon et al., 2019b)	1093
TBCNN (Mou et al., 2016)	695	Program as graph model (Brockschmidt et al., 2019)	159
SLAMC (Nguyen et al., 2013)	130	Coding criterion (Peng et al., 2015)	128
TransCoder (Roziere et al., 2020)	115	TreeGen (Sun et al., 2020)	124
Codex (Chen et al., 2021b)	897	AlphaCode (Li et al., 2022)	317

Popular models: As part of collecting metadata and summarizing studies, we identified the proposed model, if any, for each selected study. We considered novel proposed models only and not the name of the approach or method in this analysis. We also obtained the number of citations for the study. In Table 5, we present the most popular model, in no particular order, by using the number of citations as the metric to decide the popularity. We collected the number of citations at the end of August 2023 and included all the models with corresponding citations over 100.

In the rest of this section, we delve into each category and sub-category at a time, break down the entire workflow of a code analysis task into fine-grained steps, and summarize the method and ML techniques used. It is worth emphasizing that we structure the discussion around the crucial steps for each category (e.g., model generation, data sampling, feature extraction, and model training).

3.1. Code representation

Raw source code cannot be fed directly to a DL model. Code representation is the fundamental activity to make source code compatible

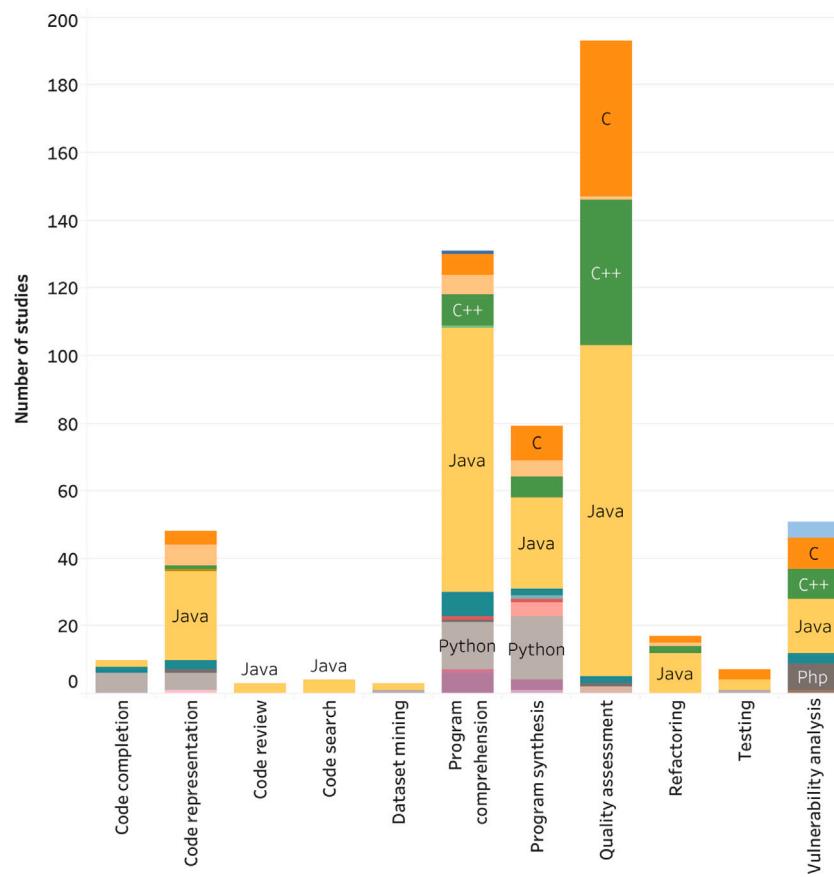


Fig. 6. Target programming languages for each considered category.

with DL models by preparing a numerical representation of the code to further solve a specific software engineering task. Code representation is the process of transforming the textual program source code into a numerical representation i.e., vectors that a DL model can accept and process (Keller et al., 2021). Studies in this category emphasize that source code is a richer construct and hence should not be treated simply as a collection of tokens or text (Nguyen et al., 2018; Allamanis et al., 2018a); the proposed techniques extensively utilize the syntax, structure, and semantics (such as type information from an AST). The activity transforms source code into a numerical representation making it easier to further use the code by ML models to solve specific tasks such as code pattern identification (Mou et al., 2016; Thaller et al., 2019), method name prediction (Alon et al., 2019b), and comment classification (Wang et al., 2020a).

In the training phase, a large number of repositories are processed to train a model which is then used in the inference phase. Source code is pre-processed to extract a source code model (such as an AST or a sequence of tokens) which is fed into a feature extractor responsible to mine the necessary features (for instance, AST paths and tree-based embeddings). Then, an ML model is trained using the extracted features. The model produces a numerical (i.e., a vector) representation that can be used further for specific software engineering applications such as defect prediction, vulnerability detection, and code smells detection.

Dataset preparation: Code representation efforts start with preparing a source code model. The majority of the studies use the AST representation (Nguyen et al., 2018; Alon et al., 2018; Zhang et al., 2019; Allamanis et al., 2015b; Chen et al., 2019e; Alon et al., 2019a,b; Yahav, 2018; Brockschmidt et al., 2019; Wang et al., 2020b; Chakraborty et al., 2022; Peng et al., 2015, 2021). Some studies (Shedko et al., 2020; Allamanis and Sutton, 2013a; Azcona et al., 2019; Chakraborty and Ray, 2021; Zheng et al., 2019; Kanade et al., 2020; Nguyen et al., 2013;

Movshovitz-Attias and Cohen, 2013; Efstathiou and Spinellis, 2019) parsed the source code as tokens and prepared a sequence of tokens in this step. Hoang et al. (2020) generated tokens representing only the code changes. Furthermore, Sui et al. (2020) compiled a program into LLVM-IR. An inter-procedural value-flow graph (ivFG) used was built on top of the intermediate representation. Thaller et al. (2019) used abstract semantic graphs as their code model. Nie et al. (2021) used dataset offered by Jiang et al. (2018) that offers a large number code snippets and comment pairs. Finally, Brauckmann et al. (2020) and Tufano et al. (2018) generated multiple source code models (AST, CFG, and byte code).

Feature extraction: Relevant features need to be extracted from the prepared source code model for further processing. The first category of studies, based on applied feature extraction mechanism, uses token-based features. Nguyen et al. (2018) prepared vectors of syntactic context (referred to as *syntaxeme*), type context (*sememes*), and lexical tokens. Shedko et al. (2020) generated a stream of tokens corresponding to function calls and control flow expressions. Karampatsis et al. (2020) split tokens as subwords to enable subwords prediction. Path-based abstractions is the basis of the second category where the studies extract a path typically from an AST. Alon et al. (2018) used paths between AST nodes. Kovalenko et al. (2020) extracted path context representing two tokens in code and a structural connection along with paths between AST nodes. Alon et al. (2019a) encoded each AST path with its values as a vector and used the average of all of the k paths as the decoder's initial state where the value of k depends on the number of leaf nodes in the AST. The decoder then generated an output sequence while attending over the k encoded paths. Peng et al. (2015) proposed “coding criterion” to capture similarity among symbols based on their usage using AST structural information. Peng et al. (2021) used open-source parser Tree-Sitter to obtain AST for each method. They split code

tokens into sub-tokens respective to naming conventions and generate path using AST nodes. The authors sets 32 as the maximum path length. Finally, [Alon et al. \(2019b\)](#) also used path-based features along with distributed representation of context where each of the path and leaf-values of a path-context is mapped to its corresponding real-valued vector representation.

Another set of studies belong to the category that used graph-based features. [Chen et al. \(2019e\)](#) created AST node identified by an API name and attached each node to the corresponding AST node belonging to the identifier. [Thaller et al. \(2019\)](#) proposed feature maps; feature maps are human-interpretation, stacked, named subtrees extracted from abstract semantic graph. [Brauckmann et al. \(2020\)](#) created a dataflow-enriched AST graph, where nodes are labeled as declarations, statements, and types as found in the Clang¹ AST. [Cvitkovic et al. \(2019\)](#) augmented AST with semantic information by adding a graph-structured vocabulary cache. Finally, [Zhang et al. \(2019\)](#) extracted small statement trees along with multi-way statement trees to capture the statement-level lexical and syntactical information. The final category of studies used DL ([Hoang et al., 2020](#); [Tufano et al., 2018](#)) to learn features automatically.

ML model training: The majority of the studies rely on the RNN-based DL model. Among them, some of the studies ([Wang et al., 2020a](#); [Hellendoorn and Devanbu, 2017](#); [Wang et al., 2020b](#); [Brauckmann et al., 2020](#); [Alon et al., 2019a](#)) employed LSTM-based models; while others ([Zhang et al., 2019](#); [Hoang et al., 2020](#); [Karampatsis et al., 2020](#); [Yahav, 2018](#); [Brockschmidt et al., 2019](#)) used GRU-based models. Among the other kinds of ML models, studies employed GNN-based ([Cvitkovic et al., 2019](#); [Wang et al., 2020c](#)), DNN ([Nguyen et al., 2018](#)), conditional random fields ([Alon et al., 2018](#)), SVM ([Lim, 2018](#); [Rabin et al., 2020](#)), CNN-based models ([Chen et al., 2019e](#); [Mou et al., 2016](#); [Thaller et al., 2019](#)), and transformer-based models ([Peng et al., 2021](#)). Some of the studies rely on the combination of different DL models. For example, [Tufano et al. \(2018\)](#) employed RNN-based model for learning embedding in the first stage which is given to an autoencoder-based model to encode arbitrarily long streams of embeddings.

A typical output of a code representation technique is the vector representation of the source code. The exact form of the output vector may differ based on the adopted mechanism. Often, the code vectors are application specific depending upon the nature of features extracted and training mechanism. For example, Code2Vec produces code vectors trained for method name prediction; however, the same mechanism can be used for other applications after tuning and selecting appropriate features. [Kang et al. \(2019\)](#) carried out an empirical study to observe whether the embeddings generated by Code2Vec can be used in other contexts. Similarly, [Pour et al. \(2021\)](#) used Code2Vec, Code2Seq, and CodeBERT to explore the robustness of code embedding models by retraining the models using the generated adversarial examples.

The semantics of the produced embeddings depend significantly on the selected features. Studies in this domain identify this aspect and hence swiftly focused to extract features that capture the relevant semantics; for example, path-based features encode the order among the tokens. The chosen ML model plays another important role to generate effective embeddings. Given the success of RNN with text processing tasks, due to its capability to identify sequence and pattern, RNN-based models dominate this category.

3.2. Testing

In this section, we point out the state-of-the-art regarding ML techniques applied to software testing. Testing is the process of identifying functional or non-functional bugs to improve the accuracy and reliability of a software. In this section, we offer a discussion on test cases generation by employing ML techniques.

¹ <https://clang.llvm.org/>.

3.2.1. Test data and test cases generation

A usual approach to have a ML model for generating test oracles involves capturing data from an application under test, pre-processing the captured data, extracting relevant features, using an ML algorithm, and evaluating the model.

Dataset preparation: Researchers developed a number of ways for capturing data from applications under test and pre-process them before feeding them to an ML model. [Braga et al. \(2018\)](#) recorded traces for applications to capture usage data. They sanitized any irrelevant information collected from the programs recording components. AppFlow ([Hu et al., 2018c](#)) captures human-event sequences from a smart-phone screen in order to identify tests. Similarly, [Nguyen et al. \(2019\)](#) suggested Shinobi, a framework that uses a fast R-CNN model to identify input data fields from multiple web-sites. [Utting et al. \(2020\)](#) captured user and system execution traces to help generating missing API tests. To automatically identify metamorphic relations, [Nair et al. \(2019\)](#) suggested an approach that leverages ML techniques and test mutants. By using a variety of code transformation techniques, the authors' approach can generate a synthetic dataset for training models to predict metamorphic relations.

Feature extraction: Some authors ([Braga et al., 2018](#); [Utting et al., 2020](#)) used execution traces as features. [Kim et al. \(2018\)](#) suggested an approach that replaces SBST's meta-heuristic algorithms with deep reinforcement learning to generate test cases based on branch coverage information. [Grano et al. \(2018\)](#) used code quality metrics such as coupling, DIT, and NOF to generate test data; they use the test data generated to predict the code coverage in a continuous integration pipeline.

ML model training: Researchers used supervised and unsupervised ML algorithms to generate test data and cases. In some of the studies, the authors utilized more than one ML algorithm to achieve their goal. Specifically, several studies ([Braga et al., 2018](#); [Kim et al., 2018](#); [Utting et al., 2020](#); [Nair et al., 2019](#)) used traditional ML algorithms, such as Support Vector Machine, Naive Bayes, Decision Tree, Multilayer Perceptron, Random Forest, AdaBoost, Linear Regression. [Nguyen et al. \(2019\)](#) used the DL algorithm Fast R-CNN. Similarly, [Godefroid et al. \(2017\)](#) used LSTM to automate generating the input grammar data for fuzzing.

3.3. Program synthesis

This section summarizes the ML techniques used by automated program synthesis tools and techniques in the examined software engineering literature. Apart from a major sub-category *program repair*, we also discuss state-of-the-art corresponds to *code generation* and *program translation* sub-categories in this section.

3.3.1. Program repair

Automated Program Repair (APR) refers to techniques that attempt to automatically identify patches for a given bug (i.e., programming mistakes that can cause an unintended run-time behavior), which can be applied to software with a little or without human intervention ([Goues et al., 2019](#)). Program repair typically consists of two phases. Initially, the repair tool uses fault localization to detect a bug in the software under examination, then, it generates patches using techniques such as search-based software engineering and logic rules that can possibly fix a given bug. To validate the generated patch, the (usually manual) evaluation of the semantic correctness² of that patch follows.

According to [Goues et al. \(2019\)](#), the techniques for constructing repair patches can be divided into three categories (heuristic repair, constraint-based repair, and learning-aided repair) if we consider the following two criteria: what types of patches are constructed and how

² The term semantic correctness is a criterion for evaluating whether a generated patch is similar to the human fix for a given bug ([Liu et al., 2020](#)).

the search is conducted. Here, we are interested in learning-aided repair, which leverages the availability of previously generated patches and bug fixes to generate patches. In particular, learning-aided-based repair tools use ML to learn patterns for patch generation.

Typically, at the pre-processing step, such methods take source code of the buggy revision as an input, and those revisions that fixes the buggy revision. The revision with the fixes includes a patch carried out manually that corrects the buggy revision and a test case that checks whether the bug has been fixed. Learning-aided-based repair is mainly based on the hypothesis that similar bugs will have similar fixes. Therefore, during the training phase, such techniques can use features such as similarity metrics to match bug patterns to similar fixes. Then, the generated patches rely on those learnt patterns. Next, we elaborate upon the individual steps involved in the process of program repair using ML techniques.

Dataset preparation: The majority of the studies extract buggy project revisions and manual fixes from buggy software projects. Most studies leverage source-code naturalness. For instance, Tufano et al. (2019b) extracted millions of bug-fixing pairs from GitHub, Amorim et al. (2018) leveraged the naturalness obtained from a corpus of known fixes, and Chen et al. (2016) used natural language structures from source code. Furthermore, many studies develop their own large-scale bug benchmarks. Ahmed et al. (2018) leveraged 4500 erroneous C programs, Gopinath et al. (2014) used a suite of programs and datasets stemmed from real-world applications, Long and Rinard (2016) used a set of successful manual patches from open-source software repositories, and Mashhadi and Hemmati (2021) used the *ManySSuBs4J* dataset containing natural language description and code snippets to automatically generate code fixes. Le et al. (2015) created an oracle for predicting which bugs should be delegated to developers for fixing and which should be fixed by repair tools. Jiang et al. (2021) used a dataset containing more than 4 million methods extracted. White et al. (2019) used Spoon, an open-source library for analyzing and transforming Java source code, to build a model for each buggy program revision. Pinconschi et al. (2021) constructed a dataset containing vulnerability-fix pairs by aggregating five existing dataset (Mozilla Foundation Security Advisories, SecretPatch, NVD, Secbench, and Big-Vul). The dataset i.e., *PatchBundle* is publicly available on GitHub. Cambronero and Rinard (2019) proposed a method to generate new supervised machine learning pipelines. To achieve the goal, the study trained using a collection of 500 supervised learning programs and their associated target datasets from Kaggle. Liu et al. (2013) prepared their dataset by selecting 636 closed bug reports from the Linux kernel and Mozilla databases. Svyatkovskiy et al. (2021) constructed their experimental dataset from the 2700 top-starred Python source code repositories on GitHub. CODIT (Chakraborty et al., 2020) collects a new dataset—*Code-ChangeData*, consisting of 32,473 patches from 48 open-source GitHub projects collected from Travis Torrent.

Other studies use existing bug benchmarks, such as DEFECTS4J (Just et al., 2014) and INTROCLASS (Le Goues et al., 2015), which already include buggy revisions and human fixes, to evaluate their approaches. For instance, Saha et al. (2019), Lou et al. (2020), Zhu et al. (2021), Renzullo et al. (2021), Wang et al. (2019), and Chen et al. (2019c) leveraged DEFECTS4J for the evaluations of their approaches. Additionally, Dantas et al. (2019) used the INTROCLASS benchmark and Majd et al. (2020) conducted experiments using 119,989 C/C++ programs within CODE4BENCH. Wu et al. (2020) used the DEEPFIX dataset that contains 46,500 correct C programs and 6975 programs with errors for their graph-based DL approach for syntax error correction.

Some studies examine bugs in different programming languages. For instance, Svyatkovskiy et al. (2020) used 1.2 billion lines of source code in Python, C#, JavaScript, and TypeScript programming languages. Also, Lutellier et al. (2020) used six popular benchmarks of four programming languages (Java, C, Python, and JavaScript).

There are also studies that mostly focus on syntax errors. In particular, Gupta et al. (2019a) used 6975 erroneous C programs with

typographic errors, Santos et al. (2018) used source code files with syntax errors, and Sakkas et al. (2020) used a corpus of 4500 ill-typed OCAML programs that lead to compile-time errors. Bhatia et al. (2018) examined a corpus of syntactically correct submissions for a programming assignment. They used a dataset comprising of over 14,500 student submissions with syntax errors.

Finally, there is a number of studies that use programming assignment from students. For instance, Bhatia et al. (2018), Gupta et al. (2019a), and Sakkas et al. (2020) used a corpus of 4500 ill-typed OCAML student programs.

Feature extraction: The majority of studies utilize similarity metrics to extract similar bug patterns and, respectively, correct bug fixes. These studies mostly employ word embeddings for code representation and abstraction. In particular, Amorim et al. (2018), Svyatkovskiy et al. (2020), Santos et al. (2018), Jiang et al. (2021), and Chen et al. (2016), leveraged source-code naturalness and applied NLP-based metrics. Tian et al. (2020) employed different representation learning approaches for code changes to derive embeddings for similarity computations. Similarly, White et al. (2019) used Word2Vec to learn embeddings for each buggy program revision. Ahmed et al. (2018) used similar metrics for fixing compile-time errors. Additionally, Saha et al. (2019) leveraged a code similarity analysis, which compares both syntactic and semantic features, and the revision history of a software project under examination, from DEFECTS4J, for fixing multi-hunk bugs, i.e., bugs that require applying a substantially similar patch to different locations. Furthermore, Wang et al. (2019) investigated, using similarity metrics, how these machine-generated correct patches can be semantically equivalent to human patches, and how bug characteristics affect patch generation. Sakkas et al. (2020) also applied similarity metrics. Svyatkovskiy et al. (2021) extracted structured representation of code (for example, lexemes, ASTs, and dataflow) and learn directly a task over those representations.

There are several approaches that use logic-based metrics based on the relationships of the features used. Specifically, Van Thuy et al. (2018) extracted twelve relations of statements and blocks for Bi-gram model using Big code to prune the search space, and make the patches generated by PROPHET (Long and Rinard, 2016) more efficient and precise. Alrajeh et al. (2015) identified counterexamples and witness traces using model checking for logic-based learning to perform repair process automatically. Cai et al. (2019) used publicly available examples of faulty models written in the B formal specification language, and proposed B-repair, an approach that supports automated repair of such a formal specification. Cambronero and Rinard (2019) extracted dynamic program traces through identification of relevant APIs of the target library; the extracted traces help the employed machine learning model to generate pipelines for new datasets.

Many studies also extract and consider the context where the bugs are related to. For instance, Tufano et al. (2019b) extracted Bug-Fixing Pairs (BFPS) from millions of bug fixes mined from GitHub (used as meaningful examples of such bug-fixes), where such a pair consists of a buggy code component and the corresponding fixed code. Then, they used those pairs as input to an Encoder–Decoder Natural Machine Translation (NMT) model. For the extraction of the pair, they used the GUMTREE SPOON AST Diff tool (Falleri et al., 2014). Additionally, Soto and Le Goues (2018) constructed a corpus by delimiting debugging regions in a provided dataset. Then, they recursively analyzed the differences between the Simplified Syntax Trees associated with EditEvent's. Mesbah et al. (2019) also generated AST diffs from the textual code changes and transformed them into a domain-specific language called Delta that encodes the changes that must be made to make the code compile. Then, they fed the compiler diagnostic information (as source) and the Delta changes that resolved the diagnostic (as target) into a Neural Machine Translation network for training. Furthermore, Li et al. (2020a) used the prior bug fixes and the surrounding code contexts of the fixes for code transformation learning. Saha et al. (2017) developed

a ML model that relies on four features derived from a program's context, i.e., the source-code surrounding the potential repair location, and the bug report. Similarly, Mashhadi and Hemmati (2021) used a combination of natural language text and corresponding code snippet to generate an aggregated sequence representation for the downstream task. Finally, Bader et al. (2019) utilized a ranking technique that also considers the context of a code change, and selects the most appropriate fix for a given bug. Vasic et al. (2019) used results from localization of variable-misuse bugs. Wu et al. (2020) developed an approach, GGF, for syntax-error correction that treats the code as a mixture of the token sequences and graphs. Lin et al. (2021b) and Zhu et al. (2021) utilized AST paths to generate code embeddings to predict the correctness of a patch. Chakraborty et al. (2020) represent the patches in a parse tree form and extract the necessary information (e.g., grammar rules, tokens, and token-types) from them. They used GumTree,³ a tree-based code differencing tool, to identify the edited AST nodes. To collect the edit context, their proposal, CODIT, converts the ASTs to their parse tree representation and extracts corresponding grammar rules, tokens, and token types.

ML model training: In the following, we present the main categories of ML techniques found in the examined papers.

Neural Machine Translation: This category includes papers that apply neural machine translation (NMT) for enhancing automated program repair. Such approaches can, for instance, include techniques that use examples of bug fixing for one programming language to fix similar bugs for other programming language. Lutellier et al. (2020) developed the repair tool called CoCoNuT that uses ensemble learning on the combination of CNNs and a new context-aware NMT. Additionally, Tufano et al. (2019b) used NMT techniques (Encoder–Decoder model) for learning bug-fixing patches for real defects, and generated repair patches. Mesbah et al. (2019) introduced DEEPDELTA, which used NMT for learning to repair compilation errors. Jiang et al. (2021) proposed CURE, a NMT-based approach to automatically fix bugs. Pinconschi et al. (2021) used SequenceR, a sequence-to-sequence model, to patch security faults in C programs. Zhu et al. (2021) proposed a tool Recoder, a syntax-guided edit decoder that takes encoded information and produces placeholders by selecting non-terminal nodes based on their probabilities. Chakraborty et al. (2020) developed a technique called codit that automates code changes for bug fixing using tree-based neural machine translation. In particular, they proposed a tree-based neural machine translation model, an extension of OpenNMT,⁴ to learn the probability distribution of changes in code.

Natural Language Processing: In this category, we include papers that combine natural language processing (NLP) techniques, embeddings, similarity scores, and ML for automated program repair. Tian et al. (2020) carried out an empirical study to investigate different representation learning approaches for code changes to derive embeddings, which are amendable to similarity computations. This study uses BERT transformer-based embeddings. Furthermore, Amorim et al. (2018) applied, a word embedding model (WORD2VEC), to facilitate the evaluation of repair processes, by considering the naturalness obtained from known bug fixes. Van Thuy et al. (2018) have also applied word representations, and extracted relations of statements and blocks for a Bi-gram model using Big code, to improve the existing learning-aid-based repair tool PROPHET (Long and Rinard, 2016). Gupta et al. (2019a) used word embeddings and reinforcement learning to fix erroneous C student programs with typographic errors. Tian et al. (2020) applied a ML predictor with BERT transformer-based embeddings associated with logistic regression to learn code representations in order to learn deep features that can encode the properties of patch correctness. Saha et al. (2019) used similarity analysis for repairing bugs that may require

applying a substantially similar patch at a number of locations. Additionally, Wang et al. (2019) used also similarity metrics to compare the differences among machine-generated and human patches. Santos et al. (2018) used n-grams and nns to detect and correct syntax errors.

Logic-based rules: Alrajeh et al. (2015) combined model checking and logic-based learning to support automated program repair. Cai et al. (2019) also combined model-checking and ML for program repair. Shim et al. (2020) used inductive program synthesis (DEEPERCODER), by creating a simple Domain Specific Language (DSL), and ML to generate computer programs that satisfies user requirements and specification. Sakkas et al. (2020) combined type rules and ML (i.e., multi-class classification, DNNS, and MLP) for repairing compile errors.

Probabilistic predictions: Here, we list papers that use probabilistic learning and ML approaches such as association rules, Decision Tree, and Support Vector Machine to predict bug locations and fixes for automated program repair. Long and Rinard (2016) introduced a repair tool called PROPHET, which uses a set of successful manual patches from open-source software repositories, to learn a probabilistic model of correct code, and generate patches. Soto and Le Goues (2018) conducted a granular analysis using different statement kinds to identify those statements that are more likely to be modified than others during bug fixing. For this, they used simplified syntax trees and association rules. Gopinath et al. (2014) presented a data-driven approach for fixing of bugs in database statements. For predicting the correct behavior for defect-inducing data, this study uses Support Vector Machine and Decision Tree. Saha et al. (2017) developed the ELIXIR repair approach that uses Logistic Regression models and similarity-score metrics. Bader et al. (2019) developed a repair approach called GETAFIX that uses hierarchical clustering to summarize fix patterns into a hierarchy ranging from general to specific patterns. Xiong et al. (2018) introduced L2S that uses ML to estimate conditional probabilities for the candidates at each search step, and search algorithms to find the best possible solutions. Gopinath et al. (2016) used Support Vector Machine and ID3 with path exploration to repair bugs in complex data structures. Le et al. (2015) conducted an empirical study on the capabilities of program repair tools, and applied Random Forest to predict whether using genetic programming search in APR can lead to a repair within a desired time limit. Aleti and Martinez (2021) used the most significant features as inputs to Random Forest, Support Vector Machine, Decision Tree, and multi-layer perceptron models.

Recurrent neural networks: DL approaches such as RNNs (e.g., LSTM and Transformer) have been used for synthesizing new code statements by learning patterns from a previous list of code statement, i.e., this techniques can be used to mainly predict the next statement. Such approaches often leverage word embeddings. Dantas et al. (2019) combined Doc2Vec and LSTM, to capture dependencies between source code statements, and improve the fault-localization step of program repair. Ahmed et al. (2018) developed a repair approach (TRACER) for fixing compilation errors using RNNs. Recently, Li et al. (2020a) introduced DLFix, which is a context-based code transformation learning for automated program repair. DLFix uses RNNs and treats automated program repair as code transformation learning, by learning patterns from prior bug fixes and the surrounding code contexts of those fixes. Svyatkovskiy et al. (2020) presented INTELLICODE that uses a Transformer model that predicts sequences of code tokens of arbitrary types, and generates entire lines of syntactically correct code. Chen et al. (2016) used the LSTM for synthesizing if-then constructs. Similarly, Vasic et al. (2019) applied the LSTM in multi-headed pointer networks for jointly learning to localize and repair variable misuse bugs. Bhatia et al. (2018) combined neural networks, and in particular RNNs, with constraint-based reasoning to repair syntax errors in buggy programs. Chen et al. (2019c) applied LSTM for sequence-to-sequence learning achieving end-to-end program repair through the SEQUENCEr repair tool they developed. Majd et al. (2020) developed SLDeep, statement-level software defect prediction, which uses LSTM on static code features.

³ <https://github.com/GumTreeDiff/gumtree>.

⁴ <https://opennmt.net/>.

Apart from above-mentioned techniques, White et al. (2019) developed DeepRepair, a recursive unsupervised deep learning-based approach, that automatically creates a representation of source code that accounts for the structure and semantics of lexical elements. The neural network language model is trained from the file-level corpus using embeddings.

3.3.2. Code generation

An automated code generation approach takes specification, typically in the form of natural language prompts, and generates executable code based on the specification (Yin and Neubig, 2017; Rabinovich et al., 2017; Svyatkovskiy et al., 2020). We elaborate on the studies that involve generating source code using ML techniques.

Dataset preparation: Yin and Neubig (2018) proposed a transition-based neural semantic parser, namely TRANX, which generates formal meaning representation from natural language text. They used multiple datasets for their study—dataset proposed by Dong and Lapata (2016) containing 880 geography-related questions, Django dataset (Oda et al., 2015), as well as WikiSQL dataset (Zhong et al., 2017). Similarly, Sun et al. (2020) and Shin et al. (2019a) used the HearthStone dataset (Ling et al., 2016) for Python code generation; in addition, Shin et al. (2019a) used the Spider (Yu et al., 2018) dataset for training. Liang et al. (2017) used the semantic parsing dataset WebQuestionsSP (Yih et al., 2016) consisting 3098 question–answer pairs for training and 1639 for testing. Bielik et al. (2017) used the Linux Kernel dataset (Karpathy et al., 2015), and the Hutter Prize Wikipedia dataset.⁵ Devlin et al. (2017b) evaluated their architecture on 205 real-world Flash-Fill instances (Gulwani et al., 2012). Xiong et al. (2018) used training data stemming from two Defects4J projects and their related JDK packages. Wei et al. (2019) conducted experiments on Java and Python projects collected from GitHub used by previous work (such as by Hu et al., 2018a,b; Wan et al., 2018).

Some studies curated datasets for their experiments. For example, Chen et al. (2021b) created HumanEval, a dataset containing 164 programming problems crafted manually for evaluation. Similarly, Li et al. (2022) first used a curated set of public GitHub repositories implemented in several popular languages such as C++, C#, Java, Go, and Python for pre-training. They created a dataset, CodeContests, for fine-tuning. The dataset includes problems, solutions, and test cases scraped from the Codeforces platform. Furthermore, IntelliCode (Svyatkovskiy et al., 2020) is trained on 1.2 billion lines of source code written in the Python, C#, JavaScript and TypeScript programming languages. Allamanis et al. (2018b) evaluated their models on a large dataset of 2.9 million lines of code. Cai et al. (2017) used a training set that contains 200 traces for addition, 100 traces for bubble sort, 6 traces for topological sort, and 4 traces for quicksort. Devlin et al. (2017a) used programming examples that involve induction, such as I/O examples. Shu and Zhang (2017) used training data to generate programs at various levels of complexity according to 45 predefined tasks (e.g., Split, Join, Select). Murali et al. (2018) used a corpus of about 150,000 API-manipulating Android methods. Shin et al. (2019b) propose a new approach to generate desirable distribution for the target datasets for program induction and synthesis tasks.

Feature extraction: Studies in this category extensively used AST during the feature extraction step. TRANX (Yin and Neubig, 2018) maps natural language text into an AST using a series of tree-construction actions. Similarly, Sun et al. (2020) parsed a program as an AST and decomposed the program into several context-free grammar rules. Also, the study by Yin and Neubig (2017) transformed statements to ASTs. These ASTs are generated for all well-formed programs using parsers provided by the programming language under examination. Furthermore, Rabinovich et al. (2017) developed a model that used a modular decoder, whose

sub-models are composed using natively generated ASTs. Each sub-model is associated with a specific construct in the AST grammar, and, then, it is invoked when that construct is required in the output tree.

Some studies in the category used examples of input and output to learn code generation. Euphony (Lee et al., 2018) learns good representation using easily obtainable solutions for given programs. DeepCoder (Balog et al., 2016) observes inputs and outputs, by leveraging information from interpreters. Then, DeepCoder searches for a program that matches the input–output examples. Similarly, Chen et al. (2019d) developed a neural program synthesis from input–output examples. Shu and Zhang (2017) extracted features from string transformations, i.e., input–output strings, and use the learned features to induce correct programs. Devlin et al. (2017b) used I/O programming examples and developed a DSL for synthesizing related programs.

Finally, the rest of the studies used tokens from source code as their features. For example, Chen et al. (2016) and Li et al. (2022) extracted tokens from source code. Allamanis et al. (2018b) extracted features that refer to program semantics such as variable names. Xiong et al. (2018) extracted several features, including context, variable, expression, and position features, from the source code to train their ML models. Devlin et al. (2017a) focused on extracting features from programs that involve induction. Murali et al. (2018) extracted low-level features (e.g., API calls). Liang et al. (2017) also used tokens and graphs extracted from the data sets used. Shin et al. (2019a) considered idioms (new named operators) from programs in an extended grammar. Bielik et al. (2017) leveraged language features, using datasets of ngrams in their experiments. Maddison and Tarlow (2014) considered features of variables and structural language features. Cummins et al. (2017) used language features to synthesize human-like written programs. Shin et al. (2019b) used different features related to I/O operations e.g., program size, control-flow ratio, and so on. Chen et al. (2018) extracted features from programming-language arguments. Wei et al. (2019) leveraged the power of code summarization and code generation. The input of code summarization is the output of code generation; the approach applies the relations between these tasks and proposes a dual training framework to train these tasks simultaneously using probability and attention weights along with dual constraints.

ML model training: A majority of the studies in this category relies on the RNN-based encoder–decoder architecture. TRANX (Yin and Neubig, 2018) implemented a transition system that generates an AST from a sequence of tree-constructing actions. The system is based on a LSTM-based encoder–decoder model where the encoder encodes the input tokens into its corresponding vector representation and the decoder generates the probabilities of tree-constructing actions. Also, Yin and Neubig (2017) proposed a data-driven syntax-based neural network model for generation of code in general-purpose programming languages such as Python. Cai et al. (2017) implemented recursion in the Neural Programmer-Interpreter framework that uses an LSTM controller on four tasks: grade-school addition, bubble sort, topological sort, and quicksort. Bielik et al. (2017) designed a language TChar for character-level language modeling, and program synthesis using LSTM. Cummins et al. (2017) applied LSTM to synthesize compilable, executable benchmarks. Chen et al. (2018) used reinforcement learning to predict arguments (e.g., CALL, REDUCE). Devlin et al. (2017b) presented a novel variant of the attentional RNN architecture, which allows for encoding of a variable size set of input–output examples. Wei et al. (2019) used Seq2Seq, Bi-LSTM, LSTM-based models to exploit the code summarization and code generation for automatic software development. Furthermore, Rabinovich et al. (2017) introduced Abstract Syntax Networks (ASNs), an extension of the standard encoder–decoder framework.

Some of the studies employed transformer-based models. Sun et al. (2020) proposed TreeGen for code generation. They implemented an AST reader to combine the grammar rules with AST and mitigated the long-dependency problem with the help of the attention mechanism

⁵ <http://prize.hutter1.net/>.

used in Transformers. Similarly, Li et al. (2022) implemented a transformer architecture for *AlphaCode*. Chen et al. (2021b) proposed *Codex* that is a GPT model fine-tuned on publicly available code from GitHub containing up to 12B parameters on code. *IntelliCode* by Svyatkovskiy et al. (2020) is a multilingual code completion tool that predicts sequences of code tokens of arbitrary types. *IntelliCode* is also able to generate entire lines of syntactically correct code. It uses a generative transformer model.

Euphony (Lee et al., 2018) targets a standard formulation, syntax-guided synthesis, by extending the grammar of given programs. To do so, *Euphony* uses a probabilistic model dictating the likelihood of each program. *DeepCoder* (Balog et al., 2016) leverages gradient-based optimization and integrates neural network architectures with search-based techniques. Szydlo et al. (2018) investigated the concept of source code generation of machine learning models as well as the generation algorithms for commonly used ML methods. Chen et al. (2019d) introduced a technique that is based on execution-guided synthesis and uses a synthesizer ensemble. This approach leverages semantic information to ensemble multiple neural program synthesizers. Chen et al. (2016) used latent attention to compute token weights. They found that latent attention performs better in capturing the sentence structure. Allamanis et al. (2018b) used DL models to learn semantics from programs. They used the code's graph structure and learned program representations over the generated graphs. Xiong et al. (2018) applied the gradient boosting tree algorithm to train their models. Devlin et al. (2017a) used the transfer learning and k-shot learning approach for cross-task knowledge transfer to improve program induction in limited-data scenarios. Shu and Zhang (2017) proposed NPBE (Neural Programming by Example) that teaches a DNN to compose a set of predefined atomic operations for string manipulations. Murali et al. (2018) trained a neural generator on program sketches to generate source code in a strongly typed, Java-like programming language. Liang et al. (2017) introduced the Neural Symbolic Machine (NSM), based on a sequence-to-sequence neural network induction, and apply it to semantic parsing. Shin et al. (2019a) employed non-parametric Bayesian inference to mine the code idioms that frequently occur in a given corpus and trained a neural generative model to optionally emit named idioms instead of the original code fragments. Maddison and Tarlow (2014) used models that are based on probabilistic context free grammars (PCFGs) and a neuro-probabilistic language, which are extended to incorporate additional source code-specific structures.

3.3.3. Program translation

In this section, we list studies that use ML that can be used, for instance, for translating source code from one programming language to another by learning source-code patterns. Le et al. (2020) presented a survey on DL techniques including machine translation algorithms and applications. Oda et al. (2015) used statistical machine translation (SMT) and proposed a method to automatically generate pseudo-code from source code for source-code comprehension. To evaluate their approach they conducted experiments, and generated English or Japanese pseudo-code from Python statements using SMT. Then, they found that the generated pseudo-code is mostly accurate, and it can facilitate code understanding. Roziere et al. (2020) applied unsupervised machine translation to create a transcompiler in a fully unsupervised way. TransCoder uses beam search decoding to generate multiple translations. Phan and Jannesari (2020) proposed PREFIXMAP, a code suggestion tool for all types of code tokens in the Java programming language. Their approach uses statistical machine translation that outperforms nMT. They used three corpus for their experiments—a large-scale corpus of English–German translation in NLP (Luong et al., 2017), the Conala corpus (Yin et al., 2018), which contains Python software documentation as 116,000 English sentences, and the MSR 2013 corpus (Allamanis and Sutton, 2013b).

3.4. Quality assessment

The *quality assessment* category has sub-categories *code smell detection*, *clone detection*, and *quality assessment/prediction*. In this section, we elaborate upon the state-of-the-art related to each of these categories within our scope.

3.4.1. Code smell detection

Code smells impair the code quality and make the software difficult to extend and maintain (Sharma and Spinellis, 2018). Extensive literature is available on detecting smells automatically (Sharma and Spinellis, 2018); ML techniques have been used to classify smelly snippets from non-smelly code. First, source code is pre-processed to extract individual samples (such as a class, file, or method). These samples are classified into positive and negative samples. Afterwards, relevant features are identified from the source code and those features are then fed into an ML model for training. The trained model classifies a source code sample into a smelly or non-smelly code.

Dataset preparation: The process of identifying code smells requires a dataset as a ground truth for training an ML model. Each sample of the training dataset must be tagged appropriately as smelly sample (along with target smell types) or non-smelly sample. Many authors built their datasets tagged manually with annotations. For example, Fakhouri et al. (2018) developed a manually validated oracle containing 1700 instances of linguistic smells. Pecorelli et al. (2019) created a dataset of 8.5 thousand samples of smells from 13 open-source projects. Some authors (Al-Jamimi and Ahmed, 2013; Mhawish and Gupta, 2020; Cruz et al., 2020a; Jain and Saha, 2021; Hadj-Kacem and Bouassida, 2019) employed existing datasets (Landfill and Qualitas) in their studies. Tummalapalli et al. (2022, 2020b, 2021b) used 226 WSDL files from the tera-PROMISE dataset. Oliveira et al. (2020) relied on historical data and mined smell instances from history where the smells were refactored.

Some efforts such as one by Sharma et al. (2021) used Code-Split (Sharma, 2019b,a) first to split source code files into individual classes and methods. Then, they used existing smell detection tools (Sharma et al., 2016; Sharma, 2018) to identify smells in the subject systems. They used the output of both of these tasks to identify and segregate positive and negative samples. Similarly, Kaur and Kaur (2021) used smells identified by Dr Java, EMMA, and FindBugs as their gold-set. Alazba and Aljamaan (2021) and Dewangan et al. (2021) used the dataset manually labeled instances detected by four code smell detector tools (i.e., iPlasma, PMD, Fluid Tool, Anti-Pattern Scanner, and Marinescu's detection rule). The dataset labeled six code smells collected from 74 software systems. Zhang and Dong (2021) proposed a large dataset BrainCode consisting 270,000 samples from 20 real-world applications. The study used iPlasma to identify smells in the subject systems.

Liu et al. (2019a) adopted an usual mechanism to identify their positive and negative samples. They assumed that popular well-known open-source projects are well-written and hence all of the classes/methods of these projects are by default considered free from smells. To obtain positive samples, they carried out *reverse refactoring* e.g., moving a method from a class to another class to create an instance of feature envy smell.

Feature extraction: The majority of the articles (Barbez et al., 2020; Kaur et al., 2017a; Kumar and Sureka, 2018; Gupta et al., 2019b; Agnihotri and Chug, 2020; Oliveira et al., 2020; Pritam et al., 2019; Fontana et al., 2013; Arcelli Fontana and Zanoni, 2017; Fontana et al., 2015; Thongkum and Mekruksavanich, 2020; Cruz et al., 2020b; Amorim et al., 2015; Cunha et al., 2020; Mhawish and Gupta, 2020; Liu et al., 2019a; Hadj-Kacem and Bouassida, 2018; Tummalapalli et al., 2019; Cruz et al., 2020a; Tummalapalli et al., 2022; Saidani et al., 2020; Tummalapalli et al., 2020b, 2021b; Kaur and Kaur, 2021; Gupta et al., 2021b; Dewangan et al., 2021; Alazba and Aljamaan, 2021; Jain and

Saha, 2021; Zhang and Dong, 2021; Gupta et al., 2021d) in this category use object-oriented metrics as features. These metrics include class-level metrics (such as *lines of code*, *lack of cohesion among methods*, *number of methods*, *fan-in* and *fan-out*) and method-level metrics (such as *parameter count*, *lines of code*, *cyclomatic complexity*, and *depth of nested conditional*). We observed that some of the attempts use a relatively small number of metrics (Thongkum and Mekruksavanich (2020) and Agnihotri and Chug (2020) used 10 and 16 metrics, respectively). However, some of the authors chose to experiment with a large number of metrics. For example, Amorim et al. (2015) employed 62, Mhawish and Gupta (2020) utilized 82, and Arcelli Fontana and Zanoni (2017) used 63 class-level metrics and 84 method-level metrics.

Some efforts diverge from the mainstream usage of using metrics as features and used alternative features. Lujan et al. (2020) used warnings generated from existing static analysis tools as features. Similarly, Ochodek et al. (2019) analyzed individual lines in source code to extract textual properties such as regex and keywords to formulate a set of vocabulary based features (such as bag of words). Tummalapalli et al. (2021a) and Gupta et al. (2021a) used distributed word representation techniques such as Term frequency-inverse Document Frequency (TFIDF), Continuous Bag Of Words (CBW), Global Vectors for Word Representation (GloVe), and Skip Gram. Similarly, Hadj-Kacem and Bouassida (2019) generated AST first and obtain the corresponding vector representation to train a model for smell detection. Furthermore, Sharma et al. (2021) hypothesized that DL methods can infer the features by themselves and hence explicit feature extraction is not required. They did not process the source code to extract features and feed the tokenized code to ML models.

ML model training: The type of ML models usage can be divided into three categories.

Traditional ML models: In the first category, we can put studies that use one or more traditional ML models. These models include *Decision Tree*, *Support Vector Machine*, *Random Forest*, *Naive Bayes*, *Logistic Regression*, *Linear Regression*, *Polynomial Regression*, *Bagging*, and *Multilayer Perceptron*. The majority of studies (Lujan et al., 2020; Kumar and Sureka, 2018; Gupta et al., 2019b; Agnihotri and Chug, 2020; Oliveira et al., 2020; Pritam et al., 2019; Fontana et al., 2013; Fontana et al., 2015; Pecorelli et al., 2019; Thongkum and Mekruksavanich, 2020; Cruz et al., 2020b; Di Nucci et al., 2018; Cunha et al., 2020; Tummalapalli et al., 2019; Cruz et al., 2020a; Tummalapalli et al., 2021a,b; Kaur and Kaur, 2021; Dewangan et al., 2021; Alazba and Aljamaan, 2021; Gupta et al., 2021a; Jain and Saha, 2021; Hadj-Kacem and Bouassida, 2019; Gupta et al., 2021d) in this category compared the performance of various ML models. Some of the authors experimented with individual ML models; for example, Kaur et al. (2017a) and Amorim et al. (2015) used *Support Vector Machine* and *Decision Tree*, respectively, for smell detection.

Ensemble methods: The second category of studies employed ensemble methods to detect smells. Barbez et al. (2020) and Tummalapalli et al. (2020) experimented with ensemble techniques such as *majority training ensemble* and *best training ensemble*. Saidani et al. (2020) used the Ensemble Classifier Chain (ECC) model that transforms multi-label problems into several single-label problems to find the optimal detection rules for each anti-pattern type.

DL-based models: Studies that use DL form the third category. Sharma et al. (2021) used CNN, RNN (LSTM), and autoencoders-based DL models. Hadj-Kacem and Bouassida (2018) employed autoencoder-based DL model to first reduce the dimensionality of data and Artificial Neural Network to classify the samples into smelly and non-smelly instances. Liu et al. (2019a) deployed four different DL models based on CNN and RNN. It is common to use other kinds of layers (such as embeddings, dense, and dropout) along with CNN and RNN. Gupta et al. (2021b) used eight DL models and Zhang and Dong (2021) proposed Metric-Attention-based Residual network (MARS) to detect brain

class/method. MARS used metric-attention mechanism to calculate the weight of code metrics and detect code smells.

Discussion: A typical ML model trained to classify samples into either smelly or non-smelly samples. The majority of the studies focused on a relatively small set of known code smells—*god class* (Barbez et al., 2020; Lujan et al., 2020; Kaur et al., 2017a; Gupta et al., 2019b; Agnihotri and Chug, 2020; Oliveira et al., 2020; Fontana et al., 2013; Grodzicka et al., 2020; Arcelli Fontana and Zanoni, 2017; Cruz et al., 2020b; Caram et al., 2019b; Hadj-Kacem and Bouassida, 2018), *feature envy* (Barbez et al., 2020; Kaur et al., 2017a; Agnihotri and Chug, 2020; Fontana et al., 2013; Arcelli Fontana and Zanoni, 2017; Fontana et al., 2015; Cruz et al., 2020b; Sharma et al., 2021; Hadj-Kacem and Bouassida, 2018), *long method* (Kaur et al., 2017a; Gupta et al., 2019b; Fontana et al., 2013; Grodzicka et al., 2020; Arcelli Fontana and Zanoni, 2017; Fontana et al., 2015; Cruz et al., 2020b; Azeem et al., 2019; Hadj-Kacem and Bouassida, 2018), *data class* (Kaur et al., 2017a; Oliveira et al., 2020; Fontana et al., 2013; Grodzicka et al., 2020; Arcelli Fontana and Zanoni, 2017; Fontana et al., 2015), and *complex class* (Lujan et al., 2020; Gupta et al., 2019b; Oliveira et al., 2020). Results of these efforts vary significantly; F1 score of the ML models vary between 0.3 to 0.99. Among the investigated ML models, authors widely report that *Decision Tree* (Azeem et al., 2019; Fontana et al., 2015; AL-Shaaby et al., 2020; Gupta et al., 2019b) and *Random Forest* (Azeem et al., 2019; Fontana et al., 2015; Kumar and Sureka, 2018; Arcelli Fontana and Zanoni, 2017; Mhawish and Gupta, 2020) perform the best. Other methods that have been reported better than other ML models in their respective studies are *Support Vector Machine* (Tummalapalli et al., 2020), *Boosting* (Luiz et al., 2019), and *autoencoders* (Sharma et al., 2021).

Traditional ML techniques are the prominent choice in this category because these techniques works well with fixed size, fixed column meaning vectors. Code quality metrics capture the features relevant to the identification of smells, and they have fixed size, fixed column meaning vectors. However, such vectors do not capture subjectivity inherent in the context and hence some studies rely on alternative features such as embeddings generated by AST representations to feed DL models such as RNN.

3.4.2. Code clone detection

Code clone detection is the process of identifying duplicate code blocks in a given software system. Software engineering researchers have proposed not only methods to detect code clones automatically, but, also verify whether the reported clones from existing tools are false-positives or not using ML techniques. Studies in this category prepare a dataset containing source code samples classified as clones or non-clones. Then, they apply feature extraction techniques to identify relevant features that are fed into ML models for training and evaluation. The trained models identify clones among the sample pairs.

Dataset preparation: Manual annotation is a common way to prepare a dataset for applying ML to identify code clones (Mostaeen et al., 2019, 2020; White et al., 2016). Mostaeen et al. (2019) used a set of tools (NiCad, Deckard, iClones, CCFinderX and SourcererCC) to first identify a list of code clones; they then manually validated each of the identified clone set. Yang et al. (2014) used existing code clone detection tools to generate their training set. Some authors (such as Bandara and Wijayarathna (2011) and Hammad et al. (2021)) relied on existing code-clone datasets. Zhang and Khoo (2021) used NiCad to detect all clone groups from each version of the software. The study mapped the clones from a consecutive version and used the mapping to predict clone consistency at both the clone-creating and the clone-changing time. Bui et al. (2018) deployed an interesting mechanism to prepare their code-clone dataset. They crawled through GitHub repositories to find different implementations of sorting algorithms; they collected 3500 samples from this process.

Feature extraction: The majority of the studies relied on the textual properties of the source code as features. Bandara and Wijayarathna

(2011) identified features such as the number of characters and words, identifier count, identifier character count, and underscore count using the ANTLR tool. Some studies (Mostaeen et al., 2019, 2020; Mostaeen et al., 2018) utilized line similarity and token similarity. Yang et al. (2014) and Hammad et al. (2021) computed TF-IDF along with other metrics such as position of clones in the file. Cesare et al. (2013) extracted 30 package-level features including the number of files, hashes of the files, and common filenames as they detected code clones at the package level. Zhang and Kho (2021) obtained a set of code attributes (e.g., lines of code and the number of parameters), context attribute set (e.g., method name similarity, and sum of parameter similarity). Similarly, Sheneamer and Kalita (2016) obtained metrics such as the number of constructors, number of field access, and super-constructor invocation from the program AST. They also employed program dependence graph features such as *decl assign* and *control decl*. Along the similar lines, Zhao and Huang (2018) used CFG and DFG (Data Flow Graph) for clone detection. Some of the studies (Bui et al., 2018; White et al., 2016; Fang et al., 2020b) relied on DL methods to encode the required features automatically without specifying an explicit set of features.

ML model training:

Traditional ML models: The majority of studies (Mostaeen et al., 2020; Bandara and Wijayarathna, 2011; Mostaeen et al., 2018; Sheneamer and Kalita, 2016; Zhang and Khoo, 2021) experimented with a number of ML approaches. For example, Mostaeen et al. (2020) used Bayes Network, Logistic Regression, and Decision Tree; Bandara and Wijayarathna (2011) employed Naive Bayes, K Nearest Neighbors, AdaBoost. Similarly, Sheneamer and Kalita (2016) compared the performance of Support Vector Machine, Linear Discriminant Analysis, Instance-Based Learner, Lazy K-means, Decision Tree, Naive Bayes, Multilayer Perceptron, and Logit Boost.

DL-based models: DL models such as ANN (Mostaeen et al., 2019; Mostaeen et al., 2018), DNN (Fang et al., 2020b; Zhao and Huang, 2018), and RNN with Reverse neural network (White et al., 2016) are also employed extensively. Bui et al. (2019) and Bui et al. (2018) combined neural networks for ML models' training. Specifically, Bui et al. (2019) built a Bilateral neural network on top of two underlying sub-networks, each of which encodes syntax and semantics of code in one language. Bui et al. (2018) constructed BiTBCNNs—a combination layer of sub-networks to encode similarities and differences among code structures in different languages. Hammad et al. (2021) proposed a Clone-Advisor, a DNN model trained by fine-tuning GPT-2 over the BigCloneBench code clone dataset, for predicting code tokens and clone methods.

3.4.3. Defect prediction

To pinpoint bugs in software, researchers used various ML approaches. The first step of this process is to identify the positive and negative samples from a dataset where samples could be a type of source code entity such as classes, modules, files, and methods. Next, features are extracted from the source code and fed into an ML model for training. Finally, the trained model can classify different code snippets as buggy or benign based on the encoded knowledge. To this end, we discuss the collected studies based on (1) data labeling, (2) features extract, and (3) ML model training.

Dataset preparation: To train an ML model for predicting defects in source code a labeled dataset is required. For this purpose, researchers have used some well-known and publicly available datasets. For instance, a large number of studies (Cetiner and Sahingoz, 2020; Gondra, 2008; Malhotra, 2014; Singh and Malhotra, 2017; Challagulla et al., 2008; Bhandari and Gupta, 2018; Malhotra and Singh, 2011; Singh and Chug, 2017; Ceylan et al., 2006; Wang and Yao, 2013; Chug and Dhall, 2013; Li et al., 2011; Dhamayanthi and Lavanya, 2019; Prabha and Shivakumar, 2020; Ma et al., 2012; Khan et al., 2020; Chen et al., 2020a; Dam et al., 2019; Wang et al., 2016a; Shi et al.,

2020; Dos Santos et al., 2020; Singh et al., 2020; Zhang and Wu, 2020; Butgereit, 2019; Di Martino et al., 2011; Sarro et al., 2012; Wang et al., 2018; Lin and Lu, 2021; Ren et al., 2014; Li et al., 2017; Kaur and Kaur, 2015; Okutan and Yildiz, 2014; Laradji et al., 2015; Song et al., 2019a; Pan et al., 2019; Malhotra and Kamal, 2019; Qiao et al., 2020; Manjula and Florence, 2019; Suresh et al., 2014; Erturk and Sezer, 2015; Palomba et al., 2017; Yohannese and Li, 2017; Sun et al., 2012; Dejaeger et al., 2012; Al Qasem et al., 2020; Aleem et al., 2015) used the PROMISE dataset (Sayyad Shirabad and Menzies, 2005). Some studies used other datasets in addition to PROMISE dataset. For example, Liang et al. (2019) used Apache projects and Qiao et al. (2020) used MIS dataset (Lyu, 1996). Xiao et al. (2020) utilized a Continuous Integration (CI) dataset and Pradel and Sen (2018) generated a synthetic dataset. Apart from using the existing datasets, some other studies prepared their own datasets by utilizing various GitHub projects (Malhotra and Jangra, 2017; Harman et al., 2014; Singh et al., 2020; Aggarwal, 2019; Malhotra et al., 2017; Pascarella et al., 2018; Tufano et al., 2019a) including Apache (Li et al., 2019b; Bowes et al., 2016; D'Ambros et al., 2012; Fan et al., 2019; Palomba et al., 2016; Sotto-Mayor and Kalech, 2021; Malhotra and Jain, 2012; Choudhary et al., 2018; Rathore and Kumar, 2021), Eclipse (Zimmermann et al., 2007; D'Ambros et al., 2012) and Mozilla (Madhavan and Whitehead, 2007; Knab et al., 2006) projects, or industrial data (Bowes et al., 2016).

Feature extraction: The most common features to train a defect prediction model are the source code metrics introduced by Halstead (1977), Chidamber and Kemerer (1994), and McCabe (1976). Most of the examined studies (Cetiner and Sahingoz, 2020; Gondra, 2008; Malhotra, 2014; Singh and Malhotra, 2017; Challagulla et al., 2008; Malhotra and Singh, 2011; Wang and Yao, 2013; Chug and Dhall, 2013; Malhotra and Jangra, 2017; Malhotra et al., 2017; Ma et al., 2012; Khan et al., 2020; Butgereit, 2019; Chappelly et al., 2017; Knab et al., 2006; Sethi and Gagandeep, 2016; Fan et al., 2019; Kaur and Kaur, 2015; Jing et al., 2014; Okutan and Yildiz, 2014; Laradji et al., 2015; Arar and Ayan, 2015; Aljamaan and Alazba, 2020; Song et al., 2019a; Wang et al., 2016b; Malhotra and Kamal, 2019; Qiao et al., 2020; Manjula and Florence, 2019; Sun et al., 2012; Yohannese and Li, 2017; Suresh et al., 2014; Dejaeger et al., 2012; Choudhary et al., 2018; Erturk and Sezer, 2015; Rathore and Kumar, 2021; Al Qasem et al., 2020; Palomba et al., 2016; Sotto-Mayor and Kalech, 2021; Premalatha and Srikrishna, 2017; Malhotra and Jain, 2012; Aleem et al., 2015; Pascarella et al., 2018; Tsuda et al., 2018) used a large number of metrics such as Lines of Code, Number of Children, Coupling Between Objects, and Cyclomatic Complexity. Some authors (Palomba et al., 2017; Soltanifar et al., 2016) combined detected code smells with code quality metrics. Furthermore, Felix and Lee (2017) used defect metrics such as defect density and defect velocity along with traditional code smells.

In addition to the above, some authors (Ceylan et al., 2006; Dhamayanthi and Lavanya, 2019; Bhandari and Gupta, 2018; Prabha and Shivakumar, 2020) suggested the use of dimensional space reduction techniques—such as Principal Component Analysis (PCA)—to limit the number of features. Pandey and Gupta (2018) used Sequential Forward Search (SFS) to extract relevant source code metrics. Dos Santos et al. (2020) suggested a sampling-based approach to extract source code metrics to train defect prediction models. Kaur et al. (2017b) suggested an approach to fetch entropy of change metrics. Bowes et al. Bowes et al. (2016) introduced a novel set of metrics constructed in terms of mutants and the test cases that cover and detect them.

Other authors (Pradel and Sen, 2018; Zhang and Wu, 2020) used embeddings to train models. Such studies, first generate ASTs (Li et al., 2019b; Fan et al., 2019; Li et al., 2017; Pan et al., 2019; Liang et al., 2019), a variation of ASTs such as simplified ASTs (Lin and Lu, 2021; Chen et al., 2019a), or AST-diff (Wang et al., 2018; Tufano et al., 2019a) for a selected method or file could be considered. Then, embeddings are generated either using the token vector corresponding to each node in the generated tree or extracting a set of paths from an AST.

Singh et al. (2020) proposed a method named *Transfer Learning Code Vectorizer* that generates features from source code by using a pre-trained code representation **DL** model. Another approach for detecting defects is capturing the syntax and multiple levels of semantics in the source code as suggested by **Dam et al. (2019)**. To do so, the authors trained a tree-base **LSTM** model by using source code files as feature vectors. Subsequently, the trained model receives an **AST** as input and predicts if a file is clear from bugs or not.

Wang et al. (2016a) employed the Deep Belief Network algorithm (**DBN**) to learn semantic features from token vectors, which are fetched from applications' **ASTs**. **Shi et al. (2020)** used a **DNN** model to automate the features extraction from the source code. **Xiao et al. (2020)** collected the testing history information of all previous **CI** cycles, within a **CI** environment, to train defect predict models. Likewise to the above study, **Madhavan and Whitehead (2007)** and **Aggarwal (2019)** used the changes among various versions of a software as features to train defect prediction models.

In contrast to the above studies, **Chen et al. (2020a)** suggested the **DTL-DP**, a framework to predict defects without the need of features extraction tools. Specifically, **DTL-DP** visualizes the programs as images and extracts features out of them by using a self-attention mechanism (**Vaswani et al., 2017**). Afterwards, it utilizes transfer learning to reduce the sample distribution differences between the projects by feeding them to a model.

ML model training: In the following, we present the main categories of **ML** techniques found in the examined papers.

Traditional ML models: To train models, most of the studies (**Cetiner and Sahingoz, 2020; Gondra, 2008; Malhotra, 2014; Singh and Malhotra, 2017; Challagulla et al., 2008; Bhandari and Gupta, 2018; Malhotra and Singh, 2011; Singh and Chug, 2017; Ceylan et al., 2006; Chug and Dhall, 2013; Dhamayanthi and Lavanya, 2019; Prabha and Shivakumar, 2020; Malhotra and Jangra, 2017; Malhotra et al., 2017; Hammouri et al., 2018; Pandey and Gupta, 2018; Dos Santos et al., 2020; Singh et al., 2020; Khan et al., 2020; Kaur et al., 2017b; Butgereit, 2019; Wang et al., 2016a; Qiao et al., 2020; Manjula and Florence, 2019; Sun et al., 2012; Yohannese and Li, 2017; Suresh et al., 2014; Dejaeger et al., 2012; Choudhary et al., 2018; Rathore and Kumar, 2021; Palomba et al., 2016; Sotto-Mayor and Kalech, 2021; Soltanifar et al., 2016; Premalatha and Srikrishna, 2017; Malhotra and Jain, 2012; Aleem et al., 2015; Pascarella et al., 2018; Kaur and Kaur, 2015; Okutan and Yildiz, 2014; Laradji et al., 2015; Felix and Lee, 2017; Malhotra and Kamal, 2019; Song et al., 2019a; Aljamaan and Alazba, 2020; Ren et al., 2014**) used traditional **ML** algorithms such as *Decision Tree*, *Random Forest*, *Support Vector Machine*, and *AdaBoost*. Similarly, **Jing et al. (2014)** and **Wang et al. (2016b)** used *Cost Sensitive Discriminative Learning*. In addition, other authors (**Li et al., 2011; Wang and Yao, 2013; Ma et al., 2012**) proposed changes to traditional **ML** algorithms to train their models. Specifically, **Wang and Yao (2013)** suggested a dynamic version of *AdaBoost.NC* that adjusts its parameters automatically during training. Similarly, **Li et al. (2011)** proposed *ACoForest*, an active semi-supervised learning method to sample the most useful modules to train defect prediction models. **Ma et al. (2012)** introduced *Transfer Naive Bayes*, an approach to facilitate transfer learning from cross-company data information and weighting training data.

DL-based models: In contrast to the above studies, researchers (**Chen et al., 2020a; Dam et al., 2019; Pradel and Sen, 2018; Li et al., 2019b; Sethi and Gagandeep, 2016**) used **DL** models such as **CNN** and **RNN**-based models for defect prediction. Specifically, **Chen et al. (2020a)**, **Al Qasem et al. (2020)**, **Li et al. (2017)** and **Pan et al. (2019)** used **CNN**-based models to predict bugs. **RNN**-based methods (**Dam et al., 2019; Tufano et al., 2019a; Chen et al., 2019a; Liang et al., 2019; Fan et al., 2019; Lin and Lu, 2021**) are also frequently used where variations of **LSTM** are used to for defect prediction. Moreover, by using **DL** approaches, authors achieved improved accuracy for defect prediction and they pointed out bugs in real-world applications (**Pradel and Sen, 2018; Li et al., 2019b**).

3.4.4. Quality assessment/prediction

Studies in this category assess or predict issues related to various quality attributes such as reliability, maintainability, and run-time performance. The process starts with dataset pre-processing and labeling to obtain labeled data samples. Feature extraction techniques are applied on the processed samples. The extracted features are then fed into an **ML** model for training. The trained model assesses or predicts the quality issues in the analyzed source code.

Dataset preparation: **Heo et al. (2017)** generated data to train an **ML** model in pursuit to balance soundness and relevance in static analysis by selectively allowing unsoundness only when it is likely to reduce false alarms. Similarly, **Alikhashashneh et al. (2018)** used the **Understand** tool to detect various metrics, and employed them on the Juliet test suite for **C++**. **Reddivari and Raman (2019)** extracted a subset of data belonging to open source projects such as **Ant**, **Tomcat**, and **Jedit** to predict reliability and maintainability using **ML** techniques. **Malhotra and Chug (2012)** also prepared a custom dataset using two proprietary software systems as their subjects to predict maintainability of a class.

Feature extraction: **Heo et al. (2017)** extracted 37 low-level code features for loop (such as number of **Null**, array accesses, and number of exits) and library call constructs (such as parameter count and whether the call is within a loop). Some studies (**Alikhashashneh et al., 2018; Reddivari and Raman, 2019; Malhotra and Chug, 2012**) used source code metrics as features.

ML model training: **Alikhashashneh et al. (2018)** employed *Random Forest*, *Support Vector Machine*, *K Nearest Neighbors*, and *Decision Tree* to classify static code analysis tool warnings as true positives, false positives, or false negatives. **Reddivari and Raman (2019)** predicted reliability and maintainability using the similar set of **ML** techniques. Anomaly-detection techniques such as *One-class Support Vector Machine* have been used by **Heo et al. (2017)**. They applied their method on taint analysis and buffer overflow detection to improve the recall of static analysis. Whereas, some other studies (**Alikhashashneh et al., 2018**) aimed to rank and classify static analysis warnings.

3.5. Code completion

Code auto-completion is a state-of-the-art integral feature of modern source-code editors and IDEs (**Bruch et al., 2009**). The latest generation of auto-completion methods uses **NLP** and advanced **ML** models, trained on publicly available software repositories, to suggest source-code completions, given the current context of the software-projects under examination.

Dataset preparation: The majority of the studies mined a large number of repositories to construct their own datasets. Specifically, **Gopalakrishnan et al. (2017a)** examined 116,000 open-source systems to identify correlations between the latent topics in source code and the usage of architectural developer tactics (such as authentication and load-balancing). **Han et al. (2009, 2011)** trained and tested their system by sampling 4919 source code lines from open-source projects. **Raychev et al. (2016)** used large codebases from **GITHUB** to make predictions for **JavaScript** and **Python** code completion. **Svyatkovskiy et al. (2019)** used 2700 **Python** open-source software **GITHUB** repositories for the evaluation of their novel approach, **Pythia**.

The rest of the approaches employed existing benchmarks and datasets. **Rahman et al. (2020)** trained their proposed model using the data extracted from **Aizu Online Judge (AOJ)** system. **Liu et al. (2020c,d)** performed experiments on three real-world datasets to evaluate the effectiveness of their model when compared with the state-of-the-art approaches. **Li et al. (2018)** conducted experiments on two datasets to demonstrate the effectiveness of their approach consisting of an

attention mechanism and a pointer mixture network on code completion tasks. Schuster et al. (2021) used a public archive of GitHub from 2020 (GitHub, 2020).

Feature extraction: Studies in this category extract source code information in variety of forms. Gopalakrishnan et al. (2017a) extracted relationships between topical concepts in the source code and the use of specific architectural developer tactics in that code. Liu et al. (2020c,d) introduced a self-attentional neural architecture for code completion with multi-task learning. To achieve this, they extracted the hierarchical source code structural information from the programs considered. Also, they captured the long-term dependency in the input programs, and derived knowledge sharing between related tasks. Li et al. (2018) used locally repeated terms in program source code to predict out-of-vocabulary (OoV) words that restrict the code completion. Chen and Wan (2019) proposed a tree-to-sequence (Tree2Seq) model that captures the structure information of source code to generate comments for source code. Raychev et al. (2016) used ASTs and performed prediction of a program element on a dynamically computed context. Svyatkovskiy et al. (2019) introduced a novel approach for code completion called Pythia, which exploits state-of-the-art large-scale DL models trained on code contexts extracted from ASTs.

ML model training: The studies can be classified based on the used ML technique for code completion.

Recurrent Neural Networks: For code completion, researchers mainly try to predict the next token. Therefore, most approaches use RNNs. In particular, Terada and Watanobe (2019) used LSTM for code completion to facilitate programming education. Rahman et al. (2020) also used LSTM. Wang et al. (2019) used an LSTM-based neural network combined with several techniques such as *Word Embedding* models and *Multi-head Attention Mechanism* to complete programming code. Zhong et al. (2019) applied several DL techniques, including LSTM, *Attention Mechanism* (AM), and *Sparse Point Network* (SPN) for JavaScript code suggestions.

Apart from LSTM, researchers have used RNN with different approaches to perform code suggestions. Li et al. (2018) applied neural language models, which involve attention mechanism for RNN, by learning from large codebases to facilitate effective code completion for dynamically-typed programming languages. Hussain et al. (2020) presented CodeGRU that uses GRU for capturing source codes contextual, syntactical, and structural dependencies. Yang et al. (2019a) presented REP to improve language modeling for code completion. Their approach uses learning of general token repetition of source code with optimized memory, and it outperforms LSTM. Schumacher et al. (2020) combined neural and classical ML including RNNs, to improve code recommendations.

Probabilistic Models: Earlier approaches for code completion used statistical learning for recommending code elements. In particular, Gopalakrishnan et al. (2017a) developed a recommender system using prediction models including neural networks for latent topics. Han et al. (2009, 2011) applied *Hidden Markov Models* to improve the efficiency of code-writing by supporting code completion of multiple keywords based on non-predefined abbreviated input. Proksch et al. (2015) used *Bayesian Networks* for intelligent code completion. Raychev et al. (2016) utilized a probabilistic model for code in any programming language with *Decision Tree*. Svyatkovskiy et al. (2019) proposed PYTHIA that employs a *Markov Chain* language model. Their approach can generate ranked lists of methods and API recommendations, which can be used by developers while writing programs.

Other techniques: Recently, new approaches have been developed for code completion based on multi-task learning, code representations, and NMT. For instance, Liu et al. (2020c,d) applied Multi-Task Learning (MTL) for suggesting code elements. Lee et al. (2021) developed MERGELOGGING, a DL-based merged network that uses code representations for automated logging decisions. Chen and Wan (2019) applied TREE2SEQ model with NMT techniques for code comment generation.

3.6. Program comprehension

Program comprehension techniques attempt to understand the theory of comprehension process of developers as well as the tools, techniques, and processes that influence the comprehension activity (Storey, 2005). We summarized, in the rest of the section, program comprehension studies into four sub-categories i.e., code summarization, program classification, change analysis, and entity identification/recommendation.

3.6.1. Code summarization

Code summarization techniques attempt to provide a consolidated summary of the source code entity (typically a method). A variety of attempts has been made in this direction. The majority of the studies (Chen and Zhou, 2018; LeClair et al., 2019; Liu et al., 2019c; Ahmad et al., 2020; Shido et al., 2019; Yao et al., 2019; Hu et al., 2018a; Li et al., 2020b; Wang et al., 2020a; LeClair et al., 2020; Ye et al., 2020; Wang et al., 2020b; Zhang et al., 2020a; Iyer et al., 2016; Li et al., 2021b; Zhou et al., 2022; Haque et al., 2020; Zhou et al., 2021) produces a summary for a small block (such as a method). This category also includes studies that summarize small code fragments (Nazar et al., 2015), code folding within IDEs (Viuginov and Filchenkov, 2019), commit message generation (Jiang et al., 2017a; Liu et al., 2018; Jiang and McMillan, 2017; Jiang, 2019; Chen et al., 2021c; Wang et al., 2020), and title generation for online posts from code (Gao et al., 2020).

Dataset preparation: The majority of the studies (Allamanis et al., 2016; Chen and Zhou, 2018; LeClair et al., 2019; Liu et al., 2019c; Ahmad et al., 2020; Hu et al., 2018a; Chen et al., 2019b; Li et al., 2020b; Wang et al., 2020a; Wan et al., 2018; Wang et al., 2020b; Chen et al., 2021c; Zhou et al., 2021) in this category prepares pairs of code snippets and their corresponding natural language description. Specifically, Chen and Zhou (2018) used more than 66 thousand pairs of C# code and natural language description where source code is tokenized using a modified version of the ANTLR parser. Ahmad et al. (2020) conducted their experiments on a dataset containing Java and Python snippets; sequences of both the code and summary tokens are represented by a sequence of vectors. Hu et al. (2018a) and Li et al. (2020b) prepared a large dataset from 9714 GitHub projects. Similarly, Wang et al. (2020a) mined code snippets and corresponding javadoc comments for their experiment. Chen et al. (2019b) created their dataset from 12 popular open-source Java libraries with more than 10 thousand stars. They considered method bodies as their inputs and method names along with method comments as prediction targets. Psarras et al. (2019) prepared their dataset by using Weka, SystemML, DL4J, Mahout, Neuroph, and Spark as their subject systems. The authors retained names and types of methods, and local and class variables. Choi et al. (2020) collected and refined more than 114 thousand pairs of methods and corresponding code annotations from 100 open-source Java projects. Iyer et al. (2016) mined StackOverflow and extracted title and code snippet from posts that contain exactly one code snippet. Similarly, Gao et al. (2020) used a dump of StackOverflow dataset. They tokenized code snippets with respect to each programming language for pre-processing. The common steps in preprocessing identifiers include making them lower case, splitting the camel-cased and underline identifiers into sub-tokens, and normalizing the code with special tokens such as "VAR" and "NUMBER". Nazar et al. (2015) used human annotators to summarize 127 code fragments retrieved from Eclipse and NetBeans official frequently asked questions. Yang et al. (2021) built a dataset with over 300 K pairs of method and comment to evaluate their approach. Chen et al. (2021c) used dataset provided by Hu et al. (2018a) and manually categorized comments into six intention categories for 20,000 code-comment pairs. Wang et al. (2020) created a Python dataset that contains 128 thousand code-comment pairs. Zhou et al. (2019b) crawled over 6700 Java projects

from Github to extract their methods and the corresponding Javadoc comments to create their dataset.

Jiang (2019) used 18 popular Java projects from GitHub to prepare a dataset with approximately 50 thousand commits to generate commit messages automatically. Liu et al. (2020a) processed 56 popular open-source projects and selected approximately 160 K commits after filtering out the irrelevant commits. Liu et al. (2019d) used RepoRepears to identify Java repositories to process. They collected pull-request meta data by using GitHub APIs. After preprocessing the collected information, they trained a model to generate pull request description automatically. Wang et al. (2021b) prepared a dataset of 107 K commits by mining 10 K open-source repositories to generate context-aware commit messages.

Apart from source code, some of the studies used additional information generated from source code. For example, LeClair et al. (2019) used AST along with code and their corresponding summaries belonging to more than 2 million Java methods. Likewise, Shido et al. (2019) and Zhang et al. (2020a) also generated ASTs of the collected code samples. Liu et al. (2019c) utilized call dependencies along with source code and corresponding comments from more than a thousand GitHub repositories. LeClair et al. (2020) employed AST along with adjacency matrix of AST edges.

Some of the studies used existing datasets such as StaQC (Yao et al., 2018) and the dataset created by Jiang et al. (2017a). Specifically, Liu et al. (2018) and Jiang and McMillan (2017) utilized a dataset of commits provided by Jiang et al. (2017a) that contains two million commits from one thousand popular Java projects. Yao et al. (2019) and Ye et al. (2020) used StaQC dataset (Yao et al., 2018); it contains more than 119 thousand pairs of question title and code snippet related to SQL mined from StackOverflow. Xie et al. (2021) utilized two existing datasets—one each for Java (LeClair and McMillan, 2019) and Python (Barone and Sennrich, 2017). Bansal et al. (2021) evaluated their code summarization technique using a Java dataset of 2.1M Java methods from 28K projects created by LeClair and McMillan (2019). Li et al. (2021b) also used the Java dataset of 2.1M methods LeClair and McMillan (2019) to predict the inconsistent names from the implementation of the methods. Similarly, Haque et al. (2020), LeClair et al. (2021) and Haque et al. (2021) relied on the Java dataset by LeClair and McMillan (2019) for summarizing methods. Zhou et al. (2022) combined multiple datasets for their experiment. The first dataset (Hu et al., 2018a) contains over 87 thousand Java methods. The other datasets contained 2.1M Java methods (LeClair and McMillan, 2019) and 500 thousand Java methods respectively.

Efforts in the direction of automatic code folding also utilize techniques similar to code summarization. Viuginov and Filchenkov (2019) collected projects developed using IntelliJ platform. They identified the `foldable` and `FoldingDescription` elements from `workspace.xml` belonging to 335 JavaScript and 304 Python repositories.

Feature extraction: Studies investigated different techniques for code and feature representations. In the simplest form, Jiang et al. (2017a) tokenized their code and text. Jiang and McMillan (2017) extracted commit messages starting from “verb + object” and computed TFIDF for each word. Haque et al. (2021) extracted top-40 most-common action words from the dataset of 2.1 m Java methods provided by LeClair and McMillan (2019). Psarras et al. (2019) used comments as well as source code elements such as method name, variables, and method definition to prepare bag-of-words representation for each class. Liu et al. (2019c) represented the extracted call dependency features as a sequence of tokens.

Some of the studies extracted explicit features from code or AST. For example, Viuginov and Filchenkov (2019) used 17 languages as independent and 8 languages as dependent features. These features include AST features such as *depth of code blocks' root node*, *number of AST nodes*, and *number of lines in the block*. Hu et al. (2018a) and Li et al.

(2020b) transformed AST into Structure-Based Traversal (SBT). Yang et al. (2021) developed a DL approach, MMTRANS, for code summarization that learns the representation of source code from the two heterogeneous modalities of the AST, i.e., SBT sequences and graphs. Zhou et al. (2022) extracted AST and prepared tokenized code sequences and tokenized AST to feed to semantic and structural encoders respectively. Zhou et al. (2021, 2019b) tokenized source code and parse them into AST. Lin et al. (2021a) proposed block-wise AST splitting method; they split the code of a method based on the blocks in the dominator tree of the Control Flow Graph, and generated a split AST for each block. Liu et al. (2020a) worked with AST diff between commits as input to generate a commit summary. Lu et al. (2017) used Eclipse JDT to parse code snippets at method-level into AST and extracted API sequences and corresponding comments to generate comments for API-based snippets. Huang et al. (2020b) proposed a statement-based AST traversal algorithm to generate the code token sequence preserving the semantic, syntactic and structural information in the code snippet.

The most common way of representing features in this category is to encode the features in the form of embeddings or feature vectors. Specifically, LeClair et al. (2019) used embeddings layer for code, text, as well as for AST. Similarly, Choi et al. (2020) transformed each of the tokenized source code into a vector of fixed length through an embedding layer. Wang et al. (2020a) extracted the functional keyword from the code and perform positional encoding. Yao et al. (2019) used a code retrieval pre-trained model with natural language query and code snippet and annotated each code snippet with the help of a trained model. Ye et al. (2020) utilized two separate embedding layers to convert input sequences, belonging to both text and code, into high-dimensional vectors. Furthermore, some authors encode source code models using various techniques. For instance, Chen et al. (2019b) represented every input code snippet as a series of AST paths where each path is seen as a sequence of embedding vectors associated with all the path nodes. LeClair et al. (2020) used a single embedding layer for both the source code and AST node inputs to exploit a large overlap in vocabulary. Wang et al. (2020b) prepared a large-scale corpus of training data where each code sample is represented by three sequences—code (in text form), AST, and CFG. These sequences are encoded into vector forms using work2vec. Studies also explored other mechanisms to encode features. For example, Liu et al. (2018) extracted commit *diffs* and represented them as bag of words. The corresponding model ignores grammar and word order, but keeps term frequencies. The vector obtained from the model is referred to as *diff vector*. Zhang et al. (2020a) parsed code snippets into ASTs and calculated their similarity using ASTs. Allamanis et al. (2016) and Ahmad et al. (2020) employed attention-based mechanism to encode tokens. Li et al. (2021b) used GloVe, a word embedding technique, to obtain the vector representation of the context; the study included method callers and callee as well as other methods in the enclosing class as the context for a method. Similarly, Li et al. (2021a) calculated edit vectors based on the lexical and semantic differences between input code and the similar code.

ML model training: The ML techniques used by the studies in this category can be divided into the following four categories.

Encoder-decoder models: The majority of the studies used attention-based Encoder-Decoder models to generate code summaries for code snippets. We further classify the studies in three categories based on their ML implementation.

A large portion of the studies use *sequence-to-sequence based approaches*. For instance, Gao et al. (2020) proposed an end-to-end sequence-to-sequence system enhanced with an attention mechanism to perform better content selection. A code snippet is transformed by a source-code encoder into a vector representation; the decoder reads the code embeddings to generate the target question titles. Jiang et al. (2017a) trained an NTM algorithm to “translate” from *diffs* to commit messages. Iyer et al. (2016) used an attention-based neural network

to model the conditional distribution of a natural language summary. Their approach uses an LSTM model guided by attention on the source code snippet to generate a summary of one word at a time. Choi et al. (2020) transformed input source code into a context vector by detecting local structural features with CNNs. Also, attention mechanism is used with encoder CNNs to identify interesting locations within the source code. Similarly, Jiang (2019), Haque et al. (2020), Liu et al. (2019d), Lu et al. (2017) and Takahashi et al. (2019) employed LSTM-based Encoder–Decoder model to generate summaries. Their last module decoder generates source code summary. Ahmad et al. (2020) proposed to use Transformer to generate a natural language summary given a piece of source code. For both encoder and decoder, the Transformer consists of stacked multi-head attention and parameterized linear transformation layers. LeClair et al. (2019) used attention mechanism to not only attend words in the output summary to words in the code word representation but also to attend the summary words to parts of the AST. The concatenated context vector is used to predict the summary of one word at a time. Xie et al. (2021) designed a novel multi-task learning (MLT) approach for code summarization through mining the relationship between method-code summaries and method names. Li et al. (2021b) used RNN-based encoder–decoder model to generate a code representation of a method and check whether the current method name is inconsistent with the predicted name based on the semantic representation. Haque et al. (2021) compared five seq2seq-like approaches (*attendgru*, *ast-attendgru*, *ast-attendgru-fc*, *graph2seq*, and *code2seq*) to explore the role of action word identification in code summarization. Wang et al. (2021b) proposed a new approach, named CoRec, to translate git diffs, using attentional Encoder–Decoder model, that include both code changes and non-code changes into commit messages. Zhou et al. (2019a) presented ContextCC that uses a Seq2Seq Neural Network model with an attention mechanism to generate comments for Java methods.

Other studies relied on *tree-based approaches*. For example, Yang et al. (2021) developed a multi-modal transformer-based code summarization approach for smart contracts. Bansal et al. (2021) introduced a project-level encoder DL model for code summarization. Chen et al. (2019b) and Hu et al. (2018a) employed LSTM-based Encoder–Decoder model to generate summaries.

Rest of the studies employed *retrieval-based techniques*. Zhang et al. (2020a) proposed Rencos in which they first trained an attentional Encoder–Decoder model to obtain an encoder for all code samples and a decoder for generating natural language summaries. Second, the approach retrieves the most similar code snippets from the training set for each input code snippet. Rencos uses the trained model to encode the input and retrieves two code snippets as context vectors. It then decodes them simultaneously to adjust the conditional probability of the next word using the similarity values from the retrieved two code snippets. Li et al. (2021a) implemented their retrieve-and-edit approach by using LSTM-based models.

Extended encoder–decoder models: Many studies extended the traditional Encoder–Decoder mechanism in a variety of ways. Among them, *sequence-to-sequence based approaches* include an approach proposed by Liu et al. (2019c); they introduced CallNN that utilizes call dependency information. They employed two encoders, one for the source code and another for the call dependency sequence. The generated output from the two encoders are integrated and used in a decoder for the target natural language summarization. Wang et al. (2020a) implemented a three step approach. In the first step, functional reinforcer extracts the most critical function-indicated tokens from source code which are fed into the second module code encoder along with source code. The output of the code encoder is given to a decoder that generates the target sequence by sequentially predicting the probability of words one by one. LeClair et al. (2020) proposed to use GNN-based encoder to encode AST of each method and RNN-based encoder to model the method as a sequence. They used an attention mechanism to learn important

tokens in the code and corresponding AST. Finally, the decoder generates a sequence of tokens based on the encoder output. Zhou et al. (2022) used two encoders, semantic and structural, to generate summaries for Java methods. Their method combined text features with structure information of code snippets to train encoders with multiple graph attention layers.

Li et al. (2020b) presented a *tree-based approach* Hybrid-DeepCon model containing two encoders for code and AST along with a decoder to generate sequences of natural language annotations. Shido et al. (2019) extended TREE-LSTM and proposed Multi-way TREE-LSTM as their encoder. The rational behind the extension is that the proposed approach not only can handle an arbitrary number of ordered children, but also factor-in interactions among children. Zhou et al. (2021) trained two separate Encoder–Decoder models, one for source code sequence and another for AST via adversarial training, where each model is guided by a well-designed discriminator that learns to evaluate its outputs. Lin et al. (2021a) used a transformer to generate high-quality code summaries. The learned syntax encoding is combined with code encoding, and fed into the transformer.

Rest of the approaches adopted *retrieval-based approaches*. Ye et al. (2020) employed dual learning mechanism by using Bi-LSTM. In one direction, the model is trained for code summarization task that takes code sequence as input and summarized into a sequence of text. On the other hand, the code generation task takes the text sequence and generate code sequence. They reused the outcome of both tasks to improve performance of the other task. Liu et al. (2020a) proposed a new approach ATOM that uses the diff between commits as input. The approach used BiLSTM module to generate a new message by using diff-diff to retrieve the most relevant commit message.

Reinforcement learning models: Some of the studies exploited reinforcement learning techniques for code summary generation. In particular, Yao et al. (2019) proposed code annotation for code retrieval method that generates an natural language annotation for a code snippet so that the generated annotation can be used for code retrieval. They used *Advanced Actor–Critic* model for annotation mechanism and LSTM based model for code retrieval. Wan et al. (2018) and Wang et al. (2020b) used deep reinforcement learning model for training using annotated code samples. The trained model is an *Actor* network that generates comments for input code snippets. The *Critic* module evaluates whether the generated word is a good fit or not. Wang et al. (2020) used a hierarchical attention network for comment generation. The study incorporated multiple code features, including type-augmented abstract syntax trees and program control flows, along with plain code sequences. The extracted features are injected into an actor–critic network. Huang et al. (2020b) proposed a composite learning model, which combines the actor–critic algorithm of reinforcement learning with the encoder–decoder algorithm, to generate block comments.

Other techniques: Jiang and McMillan (2017) used *Naive Bayes* to classify the diff files into the verb groups. For automated code folding, Viuginov and Filchenkov (2019) used *Random Forest* and *Decision Tree* to classify whether a code block needs to be folded. Similarly, Nazar et al. (2015) used *Support Vector Machine* and *Naive Bayes* classifiers to generate summaries from the extracted features. Chen et al. (2021c) compared six ML techniques to demonstrate that comment category prediction can boost code summarization to reach better results. Etemadi and Monperrus (2020) compared NNGen, SimpleNNGen, and EXC-NNGen to explore the origin of nearest diffs selected by the neural network.

3.6.2. Program classification

Studies targeting this category classify software artifacts based on programming language (Ugurel et al., 2002), application domain (Ugurel et al., 2002), and type of commits (such as buggy and adaptive) (Ji et al., 2018; Meqdadi et al., 2019). We summarize these efforts

below from dataset preparation, feature extraction, and ML model training perspective.

Dataset preparation: Ma et al. (2018) identified more than 91 thousand open-source repositories from GitHub as subject systems. They created an oracle by manually classifying software artifacts from 383 sample projects. Shimonaka et al. (2016) conducted experiments on source code generated by four kinds of code generators to evaluate their technique that identify auto-generated code automatically by using ML techniques. Ji et al. (2018) and Meqdadi et al. (2019) analyzed the GitHub commit history. Ugurel et al. (2002) relied on C and C++ projects from Ibiblio and the Sourceforge archives. Levin and Yehudai (2017) used eleven popular open-source projects and annotated 1151 commits manually to train a model that can classify commits into maintenance activities. Similarly, Mariano et al. (2021, 2019) classify commits by maintenance activities; they identify a large number of open-source GitHub repositories. Along the similar lines, Meng et al. (2021) classified commits messages into categories such as bug fix and feature addition and Li et al. (2019a) predicted the impact of single commit on the program. They used popular a small set (specifically, 5 and 10 respectively) of Java projects as their dataset. Furthermore, Sabetta and Bezzi (2018) proposed an approach to classify security-related commits. To achieve the goal, they used 660 such commits from 152 open-source Java projects that are used in SAP software. Gharbi et al. (2019) created a dataset containing 29 K commits from 12 open source projects. Abdalkareem et al. (2020) built a dataset to improve the detection CI skip commits i.e., commits where '[ci skip]' or '[skip ci]' is used to skip continuous integration pipeline to execute on the pushed commit. To build the dataset, the authors used BigQuery GitHub dataset to identify repositories where at least 10% of commits skipped the CI pipeline. Altarawy et al. (2018) used three labeled data sets including one that was created with 103 applications implemented in 19 different languages to find similar applications.

Feature extraction: Features in this category of studies belong to either source code features category or repository features. A subset of studies (Shimonaka et al., 2016; Ma et al., 2018; Ugurel et al., 2002) relies on features extracted from source code token including language specific keywords and other syntactic information. Other studies (Ji et al., 2018; Meqdadi et al., 2019) collect repository metrics (such as number of changed statements, methods, hunks, and files) to classify commits. Ben-Nun et al. (2018) leveraged both the underlying data- and control-flow of a program to learn code semantics performance prediction. Gharbi et al. (2019) used TF-IDF to weight the tokens extracted from change messages. Ghadhab et al. (2021) curated a set of 768 BERT-generated features, a set of 70 code change-based features and a set of 20 keyword-based features for training a model to classify commits. Similarly, Mariano et al. (2021, 2019) extracted a 71 features majorly belonging to source code changes and keyword occurrences categories. Meng et al. (2021) and Li et al. (2019a) computed change metrics (such as number lines added and removed) as well as natural language metrics extracted from commit messages. Abdalkareem et al. (2020) employed 23 commit-level repository metrics. Sabetta and Bezzi (2018) analyzed changes in source code associated with each commit and extracted the terms that the developer used to name entities in the source code (e.g., names of classes). Similarly, LASCAD Altarawy et al. (2018) extracted terms from the source code and preprocessed terms by removing English stop words and programming language keywords.

ML model training: A variety of ML approaches have been applied. Specifically, Ma et al. (2018) used *Support Vector Machine*, *Decision Tree*, and *Bayes Network* for artifact classification. Meqdadi et al. (2019) employed *Naive Bayes*, *Ripper*, as well as *Decision Tree* and Ugurel et al. (2002) used *Support Vector Machine* to classify specific commits. Ben-Nun et al. (2018) proposed an approach based on an RNN architecture and fixed INST2VEC embeddings for code analysis tasks. Levin and Yehudai (2017) and Mariano et al. (2021, 2019) used *Decision Tree* and

Random Forest for commits classification into maintenance activities. Gharbi et al. (2019) applied *Logistic Regression* model to determine the commit classes for each new commit message. Ghadhab et al. (2021) trained a DNN classifier to fine-tune the BERT model on the task of commit classification. Meng et al. (2021) used a CNN-based model to classify code commits. Sabetta and Bezzi (2018) trained *Random Forest*, *Naive Bayes*, and *Support Vector Machine* to identify security-relevant commits. Altarawy et al. (2018) developed LASCAD using *Latent Dirichlet Allocation* and hierarchical clustering to establish similarities among software projects.

3.6.3. Change analysis

Researchers have explored applications of ML techniques to identify or predict relevant code changes (Tollin et al., 2017; Tufano et al., 2019). We briefly describe the efforts in this domain w.r.t. three major steps—dataset preparation, feature extraction, and ML model training.

Dataset preparation: Tollin et al. (2017) performed their study on two industrial projects. Tufano et al. (2019) extracted 236 K pairs of code snippets identified before and after the implementation of the changes provided in the pull requests. Kumar et al. (2017) used eBay web-services as their subject systems. Uchôa et al. (2021) used the data provided by the Code Review Open Platform (CROP), an open-source dataset that links code review data to software changes, to predict impactful changes in code review. Malhotra and Khanna (2013) considered three open-source projects to investigate the relationship between code quality metrics and change proneness.

Feature extraction: Tollin et al. (2017) extracted features related to the code quality from the issues of two industrial projects. Tufano et al. (2019) used features from pull requests to investigate the ability of a NMT model. Abbas et al. (2020) and Malhotra and Khanna (2013) computed well-known C&K metrics to investigate the relationship between change proneness and object-oriented metrics. Similarly, Kumar et al. (2017) computed 21 code quality metrics to predict change-prone web-services. Uchôa et al. (2021) combines metrics from different sources—21 features related to source code, modification history of the files, and the textual description of the change, 20 features that characterize the developer's experience, and 27 code smells detected by DesigniteJava (Sharma, 2018).

ML model training: Tollin et al. (2017) employed *Decision Tree*, *Random Forest*, and *Naive Bayes* ML algorithms for their prediction task. Tufano et al. (2019) used *Encoder-Decoder* architecture of a typical NMT model to learn the changes introduced in pull requests. Malhotra and Khanna (2013) experimented with \square , *Multilayer Perceptron*, and *Random Forest* to observe relationship between code metrics and change proneness. Abbas et al. (2020) compared ten ML models including *Random Forest*, *Decision Tree*, *Multilayer Perceptron*, and *Bayes Network*. Similarly, Kumar et al. (2017) used *Support Vector Machine* to predict change proneness in web-services. Uchôa et al. (2021) used six ML models such as *Support Vector Machine*, *Decision Tree*, and *Random Forest* to investigate whether predicted impactful changes are helpful for code reviewers.

3.6.4. Entity identification/recommendation

This category represents studies that recommend source code entities (such as method and class names) (Allamanis et al., 2015a; Malik et al., 2019; Xu et al., 2019; Jiang et al., 2019; Hellendoorn et al., 2018) or identify entities such as design patterns (Gamma et al., 1994) in code using ML (Uchiyama et al., 2014; Alhusain et al., 2013; Zanoni et al., 2015; Dwivedi et al., 2016; Chaturvedi et al., 2018). Specifically, Linstead et al. (2008) proposed a method to identify functional components in source code and to understand code evolution to analyze emergence of functional topics with time. Huang et al. (2020a) found commenting position in code using ML techniques. Uchiyama et al. (2014) identified design patterns and Abuhamad et al. (2018) recommended code authorship. Similar approaches include recommending

method name (Allamanis et al., 2015a; Jiang et al., 2019; Xu et al., 2019), method signature (Malik et al., 2019), class name (Allamanis et al., 2015a), and type inference (Hellendoorn et al., 2018). We summarize these efforts classified in three steps of applying ML techniques below.

Dataset preparation: The majority of the studies employed GitHub projects for their experiments. Specifically, Linstead et al. (2008) used two large, open source Java projects, Eclipse and ArgoUML in their experiments to apply unsupervised statistical topic models. Similarly, Hellendoorn et al. (2018) downloaded 1000 open-source TypeScript projects and extracted identifiers with corresponding type information. Abuhamad et al. (2018) evaluated their approach over the entire Google Code Jam (GCJ) dataset (from 2008 to 2016) and over real-world code samples (from 1987) extracted from public repositories on GitHub. Allamanis et al. (2015a) mined 20 software projects from GitHub to predict method and class names. Jiang et al. (2019) used the Code2Seq dataset containing 3.8 million methods as their experimental data. Ali et al. (2015) applied information retrieval techniques to automatically create traceability links in three subject systems.

A subset of studies focused on identifying design patterns using ML techniques. Uchiyama et al. (2014) performed experimental evaluations with five programs to evaluate their approach on predicting design patterns. Alhusain et al. (2013) applied a set of design patterns detection tools on 400 open source repositories; they selected all identified instances where at least two tools report a design pattern instance. Zanoni et al. (2015) manually identified 2794 design patterns instances from ten open-source repositories. Dwivedi et al. (2016) analyzed JHotDraw and identified 59 instances of abstract factory and 160 instances of adapter pattern for their experiment. Similarly, Gopalakrishnan et al. (2017b) applied their approach to discover latent topics in source code on 116,000 open-source projects. They recommended architectural tactics based on the discovered topics. Furthermore, Mahmoud and Bradshaw (2017) chose ten open-source projects to validate their topic modeling approach designed for source code.

Feature extraction: Several studies generated embeddings from their feature set. Specifically, Huang et al. (2020a) used embeddings generated from Word2vec capturing code semantics. Similarly, Jiang et al. (2019) employed Code2vec embeddings and Allamanis et al. (2015a) used embeddings that contain semantic information about sub-tokens of a method name to identify similar embeddings utilized in similar contexts. Zhang et al. (2020b) utilized knowledge graph embeddings to extract interrelations of code for bug localization.

Other studies used source code or code metadata as features. Abuhamad et al. (2018) extracted code authorship attributes from samples of code. Malik et al. (2019) used function names, formal parameters, and corresponding comments as features. Ali et al. (2015) extracted source code entity names, such as class, method, and variable names. Bavota et al. (2014) retrieved 618 features from six open-source Java systems to apply Latent Dirichlet Allocation-based feature location technique. Similarly, De Lucia et al. (2014) extracted class name, signature of methods, and attribute names from Java source code. They applied Latent Dirichlet Allocation to label source code artifacts. Gopalakrishnan et al. (2017b) processed tactics in the form of a set of textual descriptions and produced a set of weighted indicator terms. Mahmoud and Bradshaw (2017) extracted code term co-occurrence, pair-wise term similarity, and clusters of terms features and applied their approach Semantic Topic Models (STM) on them.

In addition, Uchiyama et al. (2014), Chaturvedi et al. (2018), Dwivedi et al. (2016) and Alhusain et al. (2013) used several source-code metrics as features to detect design patterns in software programs.

ML model training: The majority of studies in this category use RNN-based DL models. In particular, Huang et al. (2020a) and Hellendoorn et al. (2018) used bidirectional RNN models. Similarly, Abuhamad et al. (2018) and Malik et al. (2019) also employed RNN models to identify

code authorship and function signatures respectively. Zhang et al. (2020b) created a bug-localization tool, KGBUGLOCATOR utilizing knowledge graph embeddings and bi-directional attention models. Xu et al. (2019) employed the GRU-based Encoder–Decoder model for method name prediction. Uchiyama et al. (2014) used a hierarchical neural network as their classifier. Allamanis et al. (2015a) utilized neural language models for predicting method and class names.

Other studies used traditional ML techniques. Specifically, Chaturvedi et al. (2018) compared four ML techniques (*Linear Regression*, *Polynomial Regression*, *support vector regression*, and *neural network*). Dwivedi et al. (2016) used *Decision Tree* and Zanoni et al. (2015) trained *Naive Bayes*, *Decision Tree*, *Random Forest*, and *Support Vector Machine* to detect design patterns using ML. Ali et al. (2015) employed *Latent Dirichlet Allocation* to distinguish domain-level terms from implementation-level terms. Gopalakrishnan et al. (2017b) discovered latent topics using *Latent Dirichlet Allocation* in the large-scale corpus. The study used *Decision Tree*, *Random Forest*, and *Linear Regression* as classifiers to compute the likelihood that a given source file is associated with a given tactic.

3.7. Code review

Code Review is the process of systematically check the code written by a developer performed by one or more different developers. A very small set of studies explore the role of ML in the process of code review that we present in this section.

Dataset preparation: Lal and Pahwa (2017) labeled check-in code samples as *clean* and *buggy*. On code samples, they carried out extensive pre-processing such as normalization and label encoding. Aiming to automate code review process, Tufano et al. (2021) trained two DL architectures one for both contributor and for reviewer. They mined Gerrit and GitHub to prepare their dataset from 8904 projects. Furthermore, Thongtanunam et al. (2022) proposed AutoTransform to better handle new tokens using Byte-Pair Encoding (BPE) approach. They leveraged the dataset proposed by Tufano et al. (2021) consisting 630,858 changed methods to train a Transformer-based NMT model.

Feature extraction: Lal and Pahwa (2017) used TF-IDF to convert the code samples into vectors after applying extensive pre-processing. Tufano et al. (2021) used n-grams extracted from each commit to train their classifiers.

ML model training: Lal and Pahwa (2017) used a *Naive Bayes* model to classify samples into buggy or clean. Tufano et al. (2021) trained two DL architectures one for both contributor and for reviewer. The authors use n-grams extracted from each commit and implement their classifiers using *Decision Tree*, *Naive Bayes*, and *Random Forest*. In their revised work (Tufano et al., 2022), the authors used Text-To-Text Transfer Transformer (T5) model and shown significant improvements in DL code review models.

3.8. Code search

Code search is an activity of searching a code snippet based on individual's need typically in Q&A sites such as StackOverflow (Sachdev et al., 2018; Shuai et al., 2020; Wan et al., 2019). The studies in this category define the following coarse-grained steps. In the first step, the techniques prepare a training set by collecting source code and often corresponding description or query. A feature extraction step then identifies and extracts relevant features from the input code and text. Next, these features are fed into ML models for training which is later used to execute test queries.

Dataset preparation: Shuai et al. (2020) utilized commented code as input. Wan et al. (2019) used source code in the form of tokens, AST, and CFG. Sachdev et al. (2018) employed a simple tokenizer to extract all tokens from source code by removing non-alphanumeric tokens. Ling

et al. (2020) mined software projects from GitHub for the training of their approach. **Jiang et al.** (2017b) used existing McGill corpus and Android corpus.

Feature extraction: Code search studies typically use embeddings representing the input code. **Shuai et al.** (2020) performed embeddings on code, where source code elements (method name, API sequence, and tokens) are processed separately. They generated embeddings for code comments independently. **Wan et al.** (2019) employed a multi-modal code representation, where they learnt the representation of each modality via LSTM, TREE-LSTM and GGNN, respectively. **Sachdev et al.** (2018) identified words from source code and transformed the extracted tokens into a natural language documents. Similarly, **Ling et al.** (2020) used an unsupervised word embedding technique to construct a matching matrix to represent lexical similarities in software projects and used an RNN model to capture latent syntactic patterns for adaptive code search. **Jiang et al.** (2017b) used a fragment parser to parse a tutorial fragment in four steps (API discovery, pronoun and variable resolution, sentence identification, and sentence type identification).

ML model training: **Shuai et al.** (2020) used a CNN-based ML model named CARLCS-CNN. The corresponding model learns interdependent representations for embedded code and query by a co-attention mechanism. Based on the embedded code and query, the co-attention mechanism learns a correlation matrix and leverages row/column-wise max-pooling on the matrix. **Wan et al.** (2019) employed a multi-modal attention fusion. The model learns representations of different modality and assigns weights using an attention layer. Next, the attention vectors are fused into a single vector. **Sachdev et al.** (2018) utilized word and documentation embeddings and performed code search using the learned embeddings. Similarly, **Ling et al.** (2020) used an *autoencoder* network and a metric (believability) to measure the degree to which a sentence is approved or disapproved within a discussion in a issue-tracking system. **Jiang et al.** (2017b) used *Latent Dirichlet Allocation* to segregate all tutorial fragments into relevant clusters and identify relevant tutorial for an API.

Once an ML model is trained, code search can be initiated using a query and a code snippet. **Shuai et al.** (2020) used the given query and code sample to measure the semantic similarity using cosine similarity. **Wan et al.** (2019) ranked all the code snippets by their similarities with the input query. Similarly, **Sachdev et al.** (2018) were able to answer almost 43% of the collected StackOverflow questions directly from code.

3.9. Refactoring

Refactoring transformations are intended to improve code quality (specifically maintainability), while preserving the program behavior (functional requirements) from users' perspective (**Suryanarayana et al.**, 2014). This section summarizes the studies that identify refactoring candidates or predict refactoring commits by analyzing source code and by applying ML techniques on code. A process pipeline typically adopted by the studies in this category can be viewed as a three step process. In the first step, the source code of the projects is used to prepare a dataset for training. Then, individual samples (*i.e.*, either a method, class, or a file) is processed to extract relevant features. The extracted features are then fed to an ML model for training. Once trained, the model is used to predict whether an input sample is a candidate for refactoring or not.

Dataset preparation: The first set of studies created their own dataset for model training. For instance, **Rodriguez et al.** (2019) and **Amal et al.** (2014) created datasets where each sample is reviewed by a human to identify an applicable refactoring operation; the identified operation is carried out by automated means. **Kosker et al.** (2009) employed four versions of the same repository, computed their complexity metrics, and classified their classes as refactored if their complexity metric values are reduced from the previous version. **Nyamawe et al.** (2019)

analyzed 43 open-source repositories with 13.5 thousand commits to prepare their dataset. Similarly, **Aniche et al.** (2020) created a dataset comprising over two million refactorings from more than 11 thousand open-source repositories. **Sagar et al.** (2021) identified 5004 commits randomly selected from all the commits obtained from 800 open-source repositories where RefactoringMiner (**Tsantalis et al.**, 2020) identified at least one refactoring. Along the similar lines, **Li et al.** (2021b) used RefactoringMiner and RefDiff tools to identify refactoring operations in the selected commits. **Xu et al.** (2017) and **Krasnqi and Cleland-Huang** (2020) used manual analysis and tagging for identifying refactoring operations. **Bavota et al.** (2013) obtained 2329 classes from nine subject systems and applied topic modeling to identify latent topics and move them to an appropriate package. Similarly, **Bavota et al.** (2014) identified all classes from six software systems and applied their proposed technique namely *Methodbook* to identify move method refactoring candidates using relational topic models. Finally, **Kurbatova et al.** (2020) generated synthetic data by moving methods to other classes to prepare a dataset for feature envy smell. The rest of the studies in this category (**Kumar and Sureka**, 2017; **Kumar et al.**, 2019; **Aribandi et al.**, 2019), used the tera-PROMISE dataset containing various metrics for open-source projects where the classes that need refactoring are tagged.

Feature extraction: A variety of features, belonging to product as well as process metrics, has been employed by the studies in this category. Some of the studies rely on code quality metrics. Specifically, **Kosker et al.** (2009) computed cyclomatic complexity along with 25 other code quality metrics. Similarly, **Kumar et al.** (2019) computed 25 different code quality metrics using the SourceMeter tool; these metrics include cyclomatic complexity, class class and clone complexity, loc, outgoing method invocations, and so on. Some of the studies (**Kumar and Sureka**, 2017; **Aribandi et al.**, 2019; **Sidhu et al.**, 2022; **Wang and Godfrey**, 2014) calculated a large number of metrics. Specifically, **Kumar and Sureka** (2017) computed 102 metrics and then applied PCA to reduce the number of features to 31, while **Aribandi et al.** (2019) used 125 metrics. **Sidhu et al.** (2022) used metrics capturing design characteristics of a model including inheritance, coupling and modularity, and size. **Wang and Godfrey** (2014) computed a wide range of metrics related to clones such as number of clone fragments in a class, clone type (type1, type2, or type3), and lines of code in the cloned method.

Some other studies did not limit themselves to only code quality metrics. Particularly, **Yue et al.** (2018) collected 34 features belonging to code, evolution history, diff between commits, and co-change. Similarly, **Aniche et al.** (2020) extracted code quality metrics, process metrics, and code ownership metrics.

In addition, **Nyamawe et al.** (2019, 2020) carried out standard NLP preprocessing and generated TF-IDF embeddings for each sample. Along the similar lines, **Kurbatova et al.** (2020) used *code2vec* to generate embeddings for each method. **Sagar et al.** (2021) extracted keywords from commit messages and used GloVe to obtain the corresponding embedding. **Krasnqi and Cleland-Huang** (2020) tagged each commit message with their parts-of-speech and prepared a language model dependency tree to detect refactoring operations from commit messages. **Bavota et al.** (2013, 2014) extracted identifiers, comments, and string literals from source code. **Bavota et al.** (2013) prepared structural coupling matrix and package decomposition matrix to identify move class candidates. **Bavota et al.** (2014) applied relational topic models to derive semantic relationships between methods and define a probability distribution of topics (topic distribution model) among methods to refactor feature envy code smell.

ML model training: Majority of the studies in this category utilized traditional ML techniques. **Rodriguez et al.** (2019) proposed a method to identify web-service groups for refactoring using *K-means*, COBWEB, and expectation maximization. **Kosker et al.** (2009) trained a *Naive Bayes*-based classifier to identify classes that need refactoring. **Kumar and Sureka** (2017) used *Least Square-Support Vector Machine* (LS-SVM) along

with SMOTE as classifier. They found that LS-SVM with *Radial Basis Function* (RBF) kernel gives the best results. Nyamawe et al. (2019) recommended refactorings based on the history of requested features and applied refactorings. Their approach involves two classification tasks; first, a binary classification that suggests whether refactoring is needed or not and second, a multi-label classification that suggests the type of refactoring. The authors used *Linear Regression*, *Multinomial Naive Bayes* (MNB), *Support Vector Machine*, and *Random Forest* classifiers. Yue et al. (2018) presented CREC—a learning-based approach that automatically extracts refactored and non-refactored clones groups from software repositories, and trains an AdaBoost model to recommend clones for refactoring. Kumar et al. (2019) employed a set of ML models such as *Linear Regression*, *Naive Bayes*, *Bayes Network*, *Random Forest*, *AdaBoost*, and *Logit Boost* to develop a recommendation system to suggest the need of refactoring for a method. Amal et al. (2014) proposed the use of ANN to generate a sequence of refactoring. Aribandi et al. (2019) predicted the classes that are likely to be refactored in the future iterations. To achieve their aim, the authors used various variants of ANN, *Support Vector Machine*, as well as *Best-in-training based Ensemble* (BTE) and *Majority Voting Ensemble* (MVE) as ensemble techniques. Kurbatova et al. (2020) proposed an approach to recommend move method refactoring based on a path-based presentation of code using *Support Vector Machine*. Similarly, Aniche et al. (2020) used *Linear Regression*, *Naive Bayes*, *Support Vector Machine*, *Decision Tree*, *Random Forest*, and *Neural Network* to predict applicable refactoring operations. Sidhu et al. (2022), Xu et al. (2017) and Wang and Godfrey (2014) used DNN, gradient boosting, and *Decision Tree* respectively to identify refactoring candidate. Sagar et al. (2021) and Nyamawe et al. (2020) employed various classifiers such as *Support Vector Machine*, *Linear Regression*, and *Random Forest* to predict commits with refactoring operations.

Bavota et al. (2013, 2014) applied *Latent Dirichlet Allocation* to identify move class and move method refactoring candidates respectively. They model the documents in a given corpus as a probabilistic mixture of latent topics and model the links between document pairs as a binary variable.

3.10. Vulnerability analysis

The studies in this domain analyze source code to identify potential security vulnerabilities. In this section, we point out the state-of-the-art in software vulnerability detection using ML techniques. First, the studies prepare a dataset or identify an existing dataset for ML training. Next, the studies extract relevant features from the identified subject systems. Then, the features are fed into a ML model for training. The trained model is then used to predict vulnerabilities in the source code.

Dataset preparation: Authors used existing labeled datasets as well as created their own datasets to train ML models. Specifically, a set of studies (Pereira et al., 2019; Milosevic et al., 2017; Rahman et al., 2017; Saccente et al., 2019; Kim et al., 2019; Bilgin et al., 2020; Spreitzenbarth et al., 2014; Lin et al., 2020; Yosifova et al., 2021; Sultana et al., 2021; Law and Grépin, 2010; Pang et al., 2017; Abunadi and Alenezi, 2015; Younis and Malaiya, 2014; Vishnu and Jevitha, 2014; Khalid et al., 2019; Kim et al., 2018; Zhang and Li, 2020; Mateless et al., 2020; Du et al., 2019; Shiqi et al., 2018; Du et al., 2015; Yang et al., 2019b; Batur Şahin and Abualigah, 2021; Narayanan et al., 2018; Wang et al., 2020d; Chen et al., 2020b; Li et al., 2015; Ren et al., 2019; Ban et al., 2019) used available labeled datasets for PHP, Java, C, C++, and Android applications to train vulnerability detection models. In other cases, Russell et al. (2018) extended an existing dataset with millions of C and C++ functions and then labeled it based on the output of three static analyzers (i.e., Clang, CppCheck, and Flawfinder).

Many studies (Ma et al., 2019; Ali Alatwi et al., 2016; Cui et al., 2020; Ndichu et al., 2019; Elovici et al., 2007; Medeiros et al., 2014; Ferenc et al., 2019; Piskachev et al., 2019; Kronjee et al., 2018; Pang et al., 2016; Alves et al., 2016; Gupta et al., 2021c; Clemente et al.,

2018; Chernis and Verma, 2018; Moskovich et al., 2009; Hou et al., 2010; Santos et al., 2013; Yang et al., 2018; Zheng et al., 2020; Perl et al., 2015; Shar et al., 2015; Jimenez et al., 2019; Lin et al., 2020, 2018) created their own datasets. Ma et al. (2019), Ali Alatwi et al. (2016), Cui et al. (2020), and Gupta et al. (2021c) created datasets to train vulnerability detectors for Android applications. In particular, Ma et al. (2019) decompiled and generated CFGs of approximately 10 thousand, both benign and vulnerable, Android applications from AndroZoo and Android Malware datasets; Ali Alatwi et al. (2016) collected 5063 Android applications where 1000 of them were marked as benign and the remaining as malware; Cui et al. (2020) selected an open-source dataset comprised of 1179 Android applications that have 4416 different version (of the 1179 applications) and labeled the selected dataset by using the AndroRisk tool; and Gupta et al. (2021c) used two Android applications (Android-universal-image-loader and JHotDraw) which they have manually labeled based on the projects PMD reports (true if a vulnerability was reported in a PMD file and false otherwise). To create datasets of PHP projects, Medeiros et al. (2014) collected 35 open-source PHP projects and intentionally injected 76 vulnerabilities in their dataset. Shar et al. (2015) used phpminder to extract 15 datasets that include SQL injections, cross-site scripting, remote code execution, and file inclusion vulnerabilities, and labeled only 20% of their dataset to point out the precision of their approach. Ndichu et al. (2019) collected 5024 JavaScript code snippets from d3M, JSunpack, and 100 top websites where the half of the code snippets were benign and the other half malicious. In other cases, authors (Yang et al., 2018; Rahman et al., 2017; Perl et al., 2015) collected large number of commit messages and mapped them to known vulnerabilities by using Google's Play Store, National Vulnerability Database (NVD), Synx, Node Security Project, and so on, while in limited cases authors (Piskachev et al., 2019) manually label their dataset. Hou et al. (2010), Moskovich et al. (2009) and Santos et al. (2013) created their datasets by collecting web-page samples from StopBadWare and VxHeavens. Lin et al. (2020) constructed a dataset and manually labeled 1471 vulnerable functions and 1320 vulnerable files from nine open-source applications, named Asterisk, FFmpag, HTTPD, LibPNG, LibTIFF, OpenSSL, Pidgin, VLC Player, and Xen. Lin et al. (2018) have used more than 30,000 non-vulnerable functions and manually labeled 475 vulnerable functions for their experiments.

Feature extraction: Authors used static source code metrics, CFGs, ASTs, source code tokens, and word embeddings as features.

Source code metrics: A set of studies (Medeiros et al., 2014; Ferenc et al., 2019; Alves et al., 2016; Gupta et al., 2021c; Clemente et al., 2018; Rahman et al., 2017; Cui et al., 2020; Piskachev et al., 2019; Ren et al., 2019; Du et al., 2019; Kim et al., 2018; Medeiros et al., 2016; Abunadi and Alenezi, 2015; Law and Grépin, 2010; Sultana et al., 2021) used more than 20 static source code metrics (such as *cyclomatic complexity*, *maximum depth of class in inheritance tree*, *number of statements*, and *number of blank lines*).

Data/control flow and AST: Ma et al. (2012), Kim et al. (2019), Bilgin et al. (2020), Kronjee et al. (2018), Wang et al. (2020d), Du et al. (2015) and Medeiros et al. (2016) used CFGs, ASTs, or data flow analysis as features. More specifically, Ma et al. (2019) extracted the API calls from the CFGs of their dataset and collected information such as the usage of APIs (which APIs the application uses), the API frequencies (how many times the application uses APIs) and API sequence (the order the application uses APIs). Kim et al. (2019) extracted ASTs and GFCs which they tokenized and fed into ML models, while Bilgin et al. (2020) extracted ASTs and translated their representation of source code into a one-dimensional numerical array to feed them to a model. Kronjee et al. (2018) used data-flow analysis to extract features, while Spreitzenbarth et al. (2014) used static, dynamic analysis, and information collected

from `ltrace` to collect features and train a linear vulnerability detection model. Lin et al. (2018) created ASTs and from there they extracted code semantics as features.

Repository and file metrics: Perl et al. (2015) collected GitHub repository meta-data (*i.e.*, programming language, star count, fork count, and number of commits) in addition to source code metrics. Other authors (Pereira et al., 2019; Elovici et al., 2007) used file meta-data such as *files' creation and modification time, machine type, file size, and linker version*.

Code and Text tokens: Chernis and Verma (2018) used simple token features (*character count, character diversity, entropy, maximum nesting depth, arrow count, "if" count, "if" complexity, "while" count, and "for" count*) and complex features (*character n-grams, word n-grams, and suffix trees*). Hou et al. (2010) collected 10 features such as *length of the document, average length of word, word count, word count in a line, and number of NULL characters*. The remaining studies (Russell et al., 2018; Pang et al., 2016; Moskovich et al., 2009; Santos et al., 2013; Yang et al., 2018; Saccente et al., 2019; Zheng et al., 2020; Shar et al., 2015; Chen et al., 2020b; Narayanan et al., 2018; Russell et al., 2018; Mateless et al., 2020; Fang et al., 2020a; Zhang and Li, 2020; Pang et al., 2017; Ban et al., 2019; Yosifova et al., 2021; Lin et al., 2020) tokenized parts of the source code or text-based information with various techniques such as the most frequent occurrences of operational codes, capture the meaning of critical tokens, or applied techniques to reduce the vocabulary size in order to retrieve the most important tokens. In some other cases, authors (Li et al., 2015) used statistical techniques to reduce the feature space to reduce the number of code tokens.

Other features: Ali Alatwi et al. (2016), Ndichu et al. (2019) and Milosevic et al. (2017) extracted permission-related features. In other cases, authors (Yang et al., 2019b) combined software metrics and N-grams as features to train models and others (Shiqi et al., 2018) created text-based images to extract features. Likewise, Sultana (2017) extracted traceable patterns such as CompoundBox, Immutable, Implementor, Overrider, Sink, Stateless, FunctionObject, and LimitSel and used Understand tool to extract various software metrics. Wei et al. (2017) extracted system calls and function call-related information to use as features, while Vishnu and Jevitha (2014) extracted URL-based features like number of chars, duplicated characters, special characters, script tags, cookies, and re-directions. Padmanabhuni and Tan (2015) extracted buffer usage patterns and defensive mechanisms statements constructs by analyzing files.

Model training: To train models, the selected studies used a variety of traditional ML and DL algorithms.

Traditional ML techniques: One set of studies (Ali Alatwi et al., 2016; Ndichu et al., 2019; Pereira et al., 2019; Russell et al., 2018; Pang et al., 2016; Moskovich et al., 2009; Perl et al., 2015; Shar et al., 2015; Yosifova et al., 2021; Sultana et al., 2021; Padmanabhuni and Tan, 2015; Law and Grépin, 2010; Abunadi and Alenezi, 2015; Younis and Malaiya, 2014; Sultana, 2017; Vishnu and Jevitha, 2014; Wei et al., 2017; Du et al., 2019; Fang et al., 2020a; Medeiros et al., 2016; Du et al., 2015; Narayanan et al., 2018; Wang et al., 2020d; Chen et al., 2020b; Ren et al., 2019) used traditional ML algorithms such as *Naive Bayes, Decision Tree, Support Vector Machine, Linear Regression, Decision Tree, and Random Forest* to train their models. Specifically, Ali Alatwi et al. (2016), Russell et al. (2018) and Perl et al. (2015) selected *Support Vector Machine* because it is not affected by over-fitting when having very high dimensional variable spaces. Along the similar lines, Ndichu et al. (2019) used *Support Vector Machine* to train their model with linear kernel. Pereira et al. (2019) used *Decision Tree, Linear Regression, and Lasso* to train their models, while Abunadi and Alenezi (2015) found that *Random Forest* is the best model for predicting cross-project vulnerabilities. Compared to the above studies, Shar et al. (2015) used both supervised (*i.e.*, *Linear Regression* and *Random Forest*) and semi-supervised (*i.e.*, *Co-trained Random Forest*) algorithms to train their

models since most of that datasets were not labeled. Yosifova et al. (2021) used text-based features to train *Naive Bayes, Support Vector Machine, and Random Forest* models. Du et al. (2019) created the LEOPARD framework that does not require prior knowledge about known vulnerabilities and used *Random Forest, Naive Bayes, Support Vector Machine, and Decision Tree* to point them out.

Other studies (Medeiros et al., 2014; Ferenc et al., 2019; Piskachev et al., 2019; Kronjee et al., 2018; Alves et al., 2016; Gupta et al., 2021c; Clemente et al., 2018; Milosevic et al., 2017; Chernis and Verma, 2018; Hou et al., 2010; Santos et al., 2013; Rahman et al., 2017; Cui et al., 2020) used up to 32 different ML algorithms to train models and compared their performance. Specifically, Medeiros et al. (2014) experimented with multiple variants of *Decision Tree, Random Forest, Naive Bayes, K Nearest Neighbors, Linear Regression, Multilayer Perceptron, and Support Vector Machine* models and identified *Support Vector Machine* as the best performing classifier for their experiment. Likewise, Milosevic et al. (2017) and Rahman et al. (2017) employed multiple ML algorithms, respectively, and found that *Support Vector Machine* offers the highest accuracy rate for training vulnerability detectors. In contrast to the above studies, Ferenc et al. (2019) showed that *K Nearest Neighbors* offers the best performance for their dataset after experimenting with DNN, *K Nearest Neighbors, Support Vector Machine, Linear Regression, Decision Tree, Random Forest, and Naive Bayes*. In order to find out which is the best model for the SWAN tool, Piskachev et al. (2019) evaluated the *Support Vector Machine, Naive Bayes, Bayes Network, Decision Tree, Stump, and Ripper*. Their results pointed out the *Support Vector Machine* as the best performing model to detect vulnerabilities. Similarly, Kronjee et al. (2018), Cui et al. (2020), and Gupta et al. (2021c) compared different ML algorithms and found *Decision Tree and Random Forest* as the best performing algorithms.

DL techniques: A large number of studies (Yang et al., 2018; Saccente et al., 2019; Kim et al., 2019; Lin et al., 2020; Ban et al., 2019; Kim et al., 2018; Mateless et al., 2020; Lin et al., 2018; Shiqi et al., 2018; Batur Şahin and Abualigah, 2021) used DL methods such as CNN, RNN, and ANN to train models. In more details, Yang et al. (2018) utilized the BP-ANN algorithm to train vulnerability detectors. For the project Achilles, Saccente et al. (2019) used an array of LSTM models to train on data containing Java code snippets for a specific set of vulnerability types. In another study, Kim et al. (2019) suggested a DL framework that makes use of RNN models to train vulnerability detectors. Specifically, the authors framework first feeds the code embeddings into a Bi-LSTM model to capture the feature semantics, then an attention layer is used to get the vector weights, and, finally, passed into a dense layer to output if a code is safe or vulnerable. Compared to the studies that examined traditional ML or DL algorithms, Zheng et al. (2020) examined both of them. They used *Random Forest, K Nearest Neighbors, Support Vector Machine, Linear Regression* among the traditional ML algorithms along with Bi-LSTM, GRU, and CNN. Their results indicate Bi-LSTM as the best performing model. Lin et al. (2020) developed a benchmarking framework that can use Bi-LSTM, LSTM, Bi-GRU, GRU, DNN and Text-CNN, but can be extended to use more deep learning models. Kim et al. (2018) generating graphical semantics that reflect on code semantic features and use them for Graph Convolutional Network to automatically identify and learn semantic and extract features for vulnerability detection, while Shiqi et al. (2018) created textual images and fed them to Deep Belief Networks to classify malware.

3.11. Summary

In this section, we briefly summarize the usage of ML in a software engineering task involving source code analysis. Fig. 7 presents an overview of the pipeline that is typically used in a software engineering task that uses ML.

Dataset preparation: Preparing a dataset is the first major activity in the pipeline. The activity starts with identifying the source of required

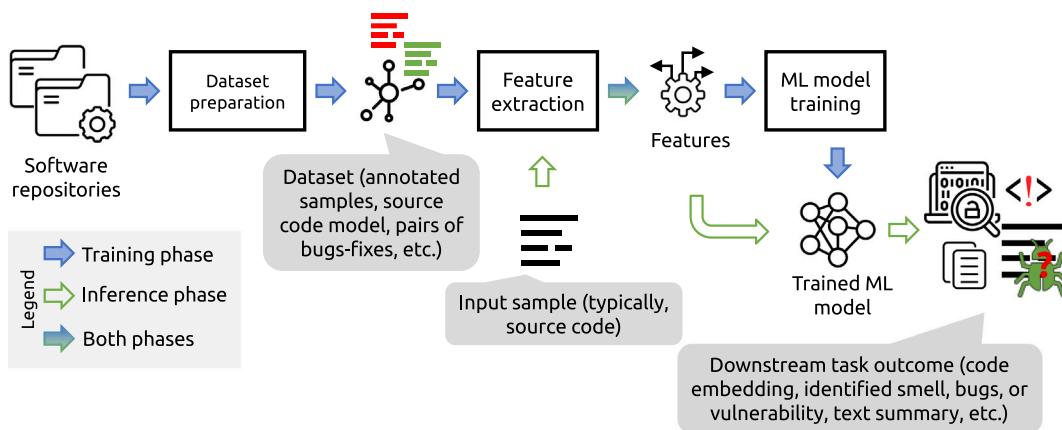


Fig. 7. Overview of the software engineering task implementation pipeline using ML.

data, typically source code repositories. The activity involves selecting and downloading the required repositories, collecting supplementary data (such as GitHub issues), create individual samples sometimes by combining information, and annotate samples. Depending upon the specific software engineering task at hand, these steps are customized and extended.

The outcome of this activity is a dataset. Depending upon the context, the dataset may contain information such as annotated code samples, source code model (e.g., AST), and pairs of buggy code and fixed code.

Feature extraction: Performance of a ML model depends significantly on the provided kind and quality of features. Various techniques are applied on the prepared dataset to extract the required features that help the ML model perform well for the given task. Features may take variety of form and format; for source code analysis applications, typical features include source code metrics, source code tokens, their properties, and representation, changes in the code (code diff), vector representation of code and text, dependency graph, and vector representation of AST, CFG, or AST diff. Obviously, selection of the specific features depends on the downstream task.

ML model training: Selecting a ML model for a given task depends on many factors such as the nature of the problem, the properties of training and input samples, and the expected output. Below, we provide an analysis of employed ML models based on these factors.

- One of the factors that influence the choice of ML models is the chosen features and their properties. Studies in the *quality assessment* category majorly relied on token-based features and code quality metrics. Such features allowed studies in this categories to use traditional ML models. Some authors applied DL models such as DNN when higher-granularity constructs such as CFG and DFG are used as features.
- Similarly, the majority of the studies in *testing* category relied on code quality metrics. Therefore, they have fixed size, fixed meaning (for each column) vectors to feed to a ML model. With such inputs, traditional ML approaches, such as *Random Forest* and *Support Vector Machine*, work well. Other studies used a variation of AST or AST of the changes to generate the embeddings. DL models including DNN and RNN-based models are used to first train a model for embeddings. A typical ML classifier use the embeddings to classify samples in buggy or benign.
- Typical output of a *code representation* study is embeddings representing code in the vector form. The semantics of the produced embeddings significantly depend on the selected features. Studies in this domain identify this aspect and, hence, they are swiftly focused to extract features that capture the relevant semantics; for example, path-based features encode the order among the

tokens. The chosen ML model plays another important role to generate effective embeddings. Given the success of RNN with text processing tasks, due to its capability to identify a sequence or pattern, RNN-based models dominate this category.

- *Program repair* is typically a sequence to sequence transformation i.e., a sequence of buggy code is the input and a sequence of fixed code is the output. Given the nature of the problem, it is not surprising to observe that the majority of the studies in this category used Encoder–Decoder-based models. RNN are considered a popular choice to realize Encoder–Decoder models due to its capability to remember long sequences.

4. Datasets and tools

For RO3, this section provides a consolidated summary of available datasets and tools that are used by the studies considered in the survey. We carefully examined each selected study and noted the resources (i.e., datasets and tools). We define the following criteria to include a resource in our catalog.

- The referenced resource must have been used by at least one primary study.
- The referenced resource must be publicly available at the time of writing this article (Dec 2022).
- The resource provides bare-minimum usage instructions to build and execute (wherever applicable) and to use the artifact.
- The resource is useful either by providing an implementation of a ML technique, helping the user to generate information/data which is further used by a ML technique, or by providing a processed dataset that can be directly employed in a ML study.

Table 6 lists all the tools that we found in this exploration. Each resource is listed with its category, name and link to access the resource, number of citations (as of Dec 2022), and the time when it was first introduced along with the time when the resource was last updated. We collected the metadata about the resources manually by searching the digital libraries, repositories, and authors' websites. The cases where we could not find the required information, are marked as “–”. We also provide a short description of the resource.

The list of datasets found in our exploration is presented in **Table 7**. Similar to the Tools' table, **Table 7** lists each resource with its category, name and link to access the resource, number of citations (as of July 2022), the time when it was first introduced along with the time when the resource was last updated, and a short description of the resource.

Table 6

A list of tools useful for applying machine learning to source code.

Category	Name	#Citation	Introd.	Updated	Description
Code representation	ncc (Ben-Nun et al., 2018)	234	Dec 2018	Aug 2021	Learns representations of code semantics
	Code2vec (Alon et al., 2019b)	487	Jan 2019	Feb 2022	Generates distributed representation of code
	Code2seq (Alon et al., 2019a)	536	May 2019	Jul 2022	Generates sequences from structured representation of code
	Vector representation for coding style (Kovalenko et al., 2020)	3	Sep 2020	Jul 2022	Implements vector representation of individual coding style
	CC2Vec (Hoang et al., 2020)	69	Oct 2020	–	Implements distributed representation of code changes
	AutoenCODE (Tufano et al., 2018)	75	–	–	Encodes source code fragments into vector representations
	Graph-based code modeling (Allamanis et al., 2018b)	544	May 2018	May 2021	Generates code modeling with graphs
	Vocabulary learning on code (Cvirkovic et al., 2019)	34	Jan 2019	–	Generates an augmented AST from Java source code
Code search	User2code2vec (Azcon et al., 2019)	29	Mar 2019	May 2019	Generates embeddings for developers based on distributed representation of code
	Deep Code Search (Gu et al., 2018)	472	May 2018	May 2022	Searches code by using code embeddings
	FRAPt (Jiang et al., 2017b)	43	Jul 2017	–	Searches relevant tutorial fragments for APIs
	Obfuscated-code2vec (Compton et al., 2020)	23	Oct 2022	–	Embeds java classes with Code2vec
Program comprehension	DEEPTYPEr (Hellendoorn et al., 2018)	87	Oct 2018	Feb 2020	Annotates types for JavaScript and TypeScript
	CallINN (Liu et al., 2019c)	9	Oct 2019	–	Implements a code summarization approach by using call dependencies
	NeuralCodeSum (Ahmad et al., 2020)	277	May 2020	Oct 2021	Implements a code summarization method by using transformers
	Summarization_tf (Shido et al., 2019)	30	Jul 2019	–	Summarizes code with Extended TREE-LSTM
	CoaCor (Yao et al., 2019)	36	Jul 2019	May 2020	Explores the role of rich annotation for code retrieval
	DeepCom (Li et al., 2020b)	102	Nov 2020	May 2021	Generates code comments
	Rencos (Zhang et al., 2020a)	79	Oct 2020	–	Generates code summary by using both neural and retrieval-based techniques
	CODES (Panichella et al., 2012)	121	Jul 2012	Jul 2016	Extracts method description from StackOverflow discussions
	CFS	–	–	–	Summarizes code fragments using SVM and NB
	TASSAL	–	–	–	Summarizes code using autofolding
	ChangeScribe (Cortes-Coy et al., 2014)	180	Dec 2014	Dec 2015	Generates commit messages
	CodeInsight (Rahman et al., 2015)	59	Nov 2015	May 2019	Recommends insightful comments for source code
	CodeNN (Iyer et al., 2016)	681	Aug 2016	May 2017	Summarizes code using neural attention model
	Code2Que (Gao et al., 2020)	25	Jul 2020	Aug 2021	Suggests improvements in question titles from mined code in StackOverflow
	Bi-TBCNN (Bui et al., 2018)	34	Mar 2019	May 2019	Implements a Bi-TBCNN model to classify algorithms
	DeepSim (Zhao and Huang, 2018)	139	Oct 2018	–	Implements a DL approach to measure code functional similarity
	FCDetector (Fang et al., 2020b)	48	Jul 2020	–	Proposes a fine-grained granularity of source code for functionality identification
	LASCAD (Altarawy et al., 2018)	12	Aug 2018	–	Categorizes software into relevant categories
	FunCom (LeClair et al., 2019)	46	May 2019	–	Summarizes code

(continued on next page)

5. Challenges and perceived deficiencies

The aim of this section is to focus on RO4 of the study by consolidating the perceived deficiencies, challenges, and opportunities in applying ML techniques to source code observed from the selected studies. We document challenges or deficiencies mentioned in the considered studies while studying and summarizing them. After the summarization phase was over, we consolidated all the documented notes and prepared a summary that we present below.

- **Standard datasets:** ML is by nature data hungry; specifically, supervised learning methods need a considerably large, cleaned, and annotated dataset. Though the size of available open software

engineering artifacts is increasing day by day, the lack of high-quality datasets (*i.e.*, clean and reliably annotated) is one of the biggest challenges in the domain (Ghaffarian and Shahriari, 2017; Ucci et al., 2019; Gondra, 2008; Kumar and Singh, 2012; Durelli et al., 2019; Chen et al., 2020a; Barbez et al., 2020; Alsolai and Roper, 2020; Tsintzira et al., 2020; Soto and Le Goues, 2018; Tian et al., 2020; Svyatkovskiy et al., 2020; Gopinath et al., 2016; Sakkas et al., 2020; Liu et al., 2019a; Wan et al., 2019; Shen and Chen, 2020; Jimenez et al., 2019). Therefore, there is a need for defining standardized datasets. Authors have cited low performance, poor generalizability, and over-fitting due to poor dataset quality as the results of the lack of standard validated high-quality datasets.

Table 6 (continued).

	SONARQUBE	–	–	–	Analyzes code quality
	SVF (Sui and Xue, 2016)	317	Mar 2016	Jul 2022	Enables inter-procedural dependency analysis for LLVM-based languages
	Designite (Sharma et al., 2016)	101	Mar 2016	Jul 2023	Detects code smells and computes quality metrics in Java and C# code
	CloneCognition (Mostaaine et al., 2018)	10	Nov 2018	May 2019	Proposes a ML framework to validate code clones
	SMAD (Barbez et al., 2020)	25	Mar 2020	Feb 2021	Implements smell detection (God class and Feature envy) using ML
Quality assessment	Checkstyle	–	–	–	Checks for coding convention in Java code
	FindBugs	–	–	–	Implements a static analysis tool for Java
	PMD	–	–	–	Finds common programming flaws in Java and six other languages
	py-ccflex (Ochodek et al., 2019)	12	Mar 2017	Oct 2020	Mimics code metrics by using ML
	Deep learning smells (Sharma et al., 2021)	27	Jul 2021	Nov 2020	Implements DL (CNN, RNN, and autoencoder-based models) to identify four smells
	GREC (Yue et al., 2018)	26	Nov 2018	–	Recommends clones for refactoring
	ML for software refactoring (Aniche et al., 2020)	31	Sep 2020	–	Recommends refactoring by using ML
	DTLDP (Chen et al., 2020a)	28	Aug 2019	–	Implements a deep transfer learning framework
	BugDetection (Li et al., 2019b)	66	Oct 2019	May 2021	Trains models for defect prediction
	DeepBugs (Pradel and Sen, 2018)	210	Nov 2018	May 2021	Implements a framework for learning name-based bug detectors
Program synthesis	CoCoNuT (Lutellier et al., 2020)	97	Jul 2020	Sep 2021	Repairs Java programs
	DeepFix (Gupta et al., 2017)	498	Feb 2017	Dec 2017	Fixes common C errors
	TRANX (Yin and Neubig, 2018)	187	Oct 2018	–	Translates natural language text to formal meaning representations
Testing	TreeGen	83	Nov 2019	–	Generates code
	AppFlow (Hu et al., 2018c)	47	Oct 2018	–	Automates UI tests generation
	DeepFuzz (Liu et al., 2019b)	72	Jul 2019	Mar 2020	Grammar fuzzer that generates C programs
	Agilika (Uutting et al., 2020)	7	Aug 2020	Mar 2022	Generates tests from execution traces
	TestDescriber	–	–	–	Implements test case summary generator and evaluator
Vulnerability analysis	Randoop	–	–	Jul 2022	Generates tests automatic for Java code
	WAP (Medeiros et al., 2013)	9	Oct 2013	Nov 2015	Detects and corrects input validation vulnerabilities
	SWAN (Piskachev et al., 2019)	8	Oct 2019	May 2022	Identifies vulnerabilities
General	VCCFinder (Perl et al., 2015)	174	Oct 2015	May 2017	Finds potentially dangerous code in repositories
	BERT (Devlin et al., 2018)	76,767	Oct 2018	Mar 2020	NLP pre-trained models
	BC3 Annotation Framework	–	–	–	Annotates emails/conversations easily
	JGibLDA	–	–	–	Implements Latent Dirichlet Allocation
	Stanford NLP Parser	–	–	–	A statistical NLP parser
	srcML	–	–	May 2022	Generates XML representation of sourcecode
	CallGraph	–	Oct 2017	Oct 2018	Generates static and dynamic call graphs for Java code
	ML for programming	–	–	–	Offers various tools such as JSNice, Nice2Predict, and DEBIN

Mitigation: Although available datasets have increased, given a wide number of software engineering tasks and variations in these tasks as well as the need of application-specific datasets, the community still looks for application-specific, large, and high-quality datasets. To mitigate the issue, the community has focused on developing new datasets and making them publicly available by organizing a dedicated track, for example, the MSR data showcase track. Dataset search engines such as the Google dataset search,⁶ Zenodo,⁷ and Kaggle datasets⁸ could be used to search available datasets. Researchers may also propose generic datasets that can

serve multiple application domains or at least different variations of a software engineering task. In addition, recent advancements in ML techniques such as active learning (Prince, 2004; Settles, 2009; Ren et al., 2020) may reduce the need of large datasets. Besides, the way the data is used for model validation must be improved. For example, Jimenez et al. (2019) showed that previous studies on vulnerability prediction trained predictive models by using perfect labeling information (*i.e.*, including future labels, as yet undiscovered vulnerabilities) and showed that such an unrealistic labeling assumption can profoundly affect the scientific conclusions of a study as the prediction performance worsen dramatically when one fully accounts for realistically available labeling. Such issues can be avoided by proposing standards for datasets laying out the minimum expectations from any public dataset.

⁶ <https://datasetsearch.research.google.com/>.⁷ <https://zenodo.org/>.⁸ <https://www.kaggle.com/datasets>.

Table 7

A list of datasets useful for applying machine learning to source code.

Category	Name	#Citation	Introd.	Updated	Description
Code representation	Code2seq (Alon et al., 2019b)	418	Jan 2019	Feb 2022	Sequences generated from structured representation of code
	GHTorrent (Gousios, 2013)	728	Oct 2013	Sep 2020	Meta-data from GitHub repositories
Code completion	Neural Code Completion	148	Nov 2017	Sep 2019	Dataset and code for code completion with neural attention and pointer networks
Program synthesis	CoNaLACorpus (Yin et al., 2018)	201	Dec 2018	Oct 2021	Python snippets and corresponding natural language description
	IntroClass (Le Goues et al., 2015)	299	Jul 2015	Feb 2016	Program repair dataset of C programs
	Code contest (Li et al., 2022)	84	Dec 2022	–	Code generation dataset for AlphaCode
Program comprehension	Program comprehension dataset (Stapleton et al., 2020)	61	May 2018	Aug 2021	Contains code for a program comprehension user survey
	CommitGen (Jiang et al., 2017a)	116	–	–	Commit messages and the diffs from 1006 Java projects
	StaQC (Yao et al., 2018)	80	Nov 2019	Aug 2021	148K Python and 120K SQL question-code pairs from StackOverflow
	TL-CodeSum (Hu et al., 2018b)	241	Feb 2019	Sep 2020	Dataset for code summarization
	DeepCom (Hu et al., 2018a)	–	May 2018	–	Dataset for code completion
Quality assessment	src-d datasets	–	–	–	Various labeled datasets (commit messages, duplicates, DockerHub, and Nuget)
	BigCloneBench (Svajlenko et al., 2014)	272	Dec 2014	Mar 2021	Known clones in the IJADataset source repository
	Multi-label smells (Guggulothu and Moiz, 2020)	28	May 2020	–	A dataset of 445 instances of two code smells and 82 metrics
	Deep learning smells (Sharma et al., 2021)	27	Jul 2021	Nov 2020	A dataset of four smells in tokenized form from 1072 C# and 100 Java repositories
	ML for software refactoring (Aniche et al., 2020)	31	Nov 2019	–	Dataset for applying ML to recommend refactoring
	QScored (Sharma and Kessentini, 2021)	11	Aug 2021	–	Code smell and metrics dataset for more than 86 thousand open-source repositories
	Landfill (Palomba et al., 2015)	34	May 2015	–	Code smell dataset with public evaluation
	KeepItSimple (Fakhoury et al., 2018)	16	Jul 2018	–	A dataset of linguistic antipatterns of 1753 instances of source code elements
	Code smell dataset (Cruz et al., 2020a)	8	Sept 2018	–	A dataset of four code smells
	Defects4J (Just et al., 2014)	858	Jul 2014	Jul 2022	Java reproducible bugs
	PROMISE (Sayyad Shirabad and Menzies, 2005)	434	–	Jan 2021	Various datasets including defect prediction and cost estimation
Testing	BugDetection (Li et al., 2019b)	59	Oct 2019	May 2021	A bug prediction dataset containing 4.973M methods belonging to 92 different Java project versions
	DEEPBUGS (Pradel and Sen, 2018)	155	Oct 2018	Apr 2021	A JavaScript code corpus with 150K code snippets
	DTLDP (Chen et al., 2020a)	28	Oct 2020	–	Dataset for deep transfer learning for defect prediction
	DAMT (Nair et al., 2019)	15	Aug 2019	Dec 2019	Metamorphic testing dataset
	wpscan	–	–	–	a PHP dataset for WordPress plugin vulnerabilities
	Genome (Zhou and Jiang, 2012)	2898	Jul 2012	Dec 2015	1200 malware samples covering the majority of existing malware families
	Juliet (Boland and Black, 2012)	147	–	–	81K synthetic C/C++ and Java programs with known flaws
Vulnerability analysis	AndroZoo (Allix et al., 2016)	–	–	–	15.7M APKs from Google's Play Store
	TRL (Lin et al., 2018)	108	Apr 2018	Jan 2019	Vulnerabilities in six C programs
	Draper vDISC (Russell et al., 2018)	479	Jul 2018	Nov 2018	1.27 million functions mined from C and C++ applications
	SAMATE (Black, 2007)	–	–	–	A set of known security flaws from NIST for C, C++, and Java programs
	jsVuln (Ferenc et al., 2019)	3	–	–	JavaScript Vulnerability Analysis dataset

(continued on next page)

- Reproducibility and replicability:** Reproducibility and replicability of any ML implementation can be compromised by the factors discussed below.

- Insufficient information:** Aspects such as the ML model, their hyper-parameters, data size and ratio (of benign and faulty samples, for instance) are required to understand and repli-

Table 7 (continued).

General	SWAN (Piskachev et al., 2019)	8	Jul 2019	Jul 2022	A Vulnerability Analysis collection of 12 Java applications
	Project-KB (Ponta et al., 2019)	49	Aug 2019	–	A Manually-Curated dataset of fixes to vulnerabilities of open-source software
	GitHub Java Corpus (Allamanis and Sutton, 2013a)	411	–	–	A large collection of Java repositories
	150k Python dataset (Raychev et al., 2016)	89	–	–	Contains parsed AST for 150K Python files
	uci source code dataset (Lopes et al., 2010)	38	Apr 2010	Nov 2013	Various large scale source code analysis datasets

cate the study. During our exploration, we found numerous studies that do not present even the bare-minimum pieces of information to replicate and reproduce their results. Likewise, Di Nucci et al. (2018) carried out a detailed replication study and reported that the replicated results were lower by up to 90% compared to what was reported in the original study.

- Handling of data imbalance: It is very common to have imbalanced datasets in software engineering applications. Authors use techniques such as under-sampling and over-sampling to overcome the challenge for training. However, test datasets must retain the original sample ratio as found in the real world (Di Nucci et al., 2018); carrying out a performance evaluation based on a balanced dataset is flawed. Obviously, the model will perform significantly inferior when it is put at work in a real-world context. We noted many studies (Agnihotri and Chug, 2020; Oliveira et al., 2020; Guggulothu and Moiz, 2020; Fontana et al., 2013; Fontana et al., 2015; Thongkum and Mekruksavanich, 2020; Cunha et al., 2020) that used balanced samples and often did not provide the size and ratio of the training and testing dataset. Such improper handling of data imbalance contributes to poor reproducibility.

Mitigation: The importance of reproducibility and replicability has been emphasized and understood by the software engineering community (Liu et al., 2020b). It has lead to a concrete artifact evaluation mechanism adopted by leading software engineering conferences. For example, FSE artifact evaluation divides artifacts into five categories—functional, reusable, available, results reproduced, and results replicated.⁹ Such thorough evaluation encouraging software engineering authors to produce high-quality documentation along with easily replicate experiment results using their developed artifacts. In addition, efforts (such as model engineering process Banna et al., 2021) are being made to support ML research reproducible and replicable. Finally, identifying practices (such as assumptions related to hardware or dependencies) that may hinder reproducibility improve reproducibility.

- Maturity in ML development:** Development of ML systems are inherently different from traditional software development (Wan et al., 2019). Phases of ML development are very exploratory in nature and highly domain and problem dependent (Wan et al., 2019). Identifying the most appropriate ML model, their appropriate parameters, and configuration is largely driven by trial and error manner (Wan et al., 2019; Azeem et al., 2019; Shen and Chen, 2020). Such an ad hoc and immature software development environment poses a huge challenge to the community. A related challenge is lack of tools and techniques for various phases and tasks involved in ML software development. It includes effective tools for testing ML programs, ensuring that the dataset are pre-processed adequately, debugging, and effective data management (Wan et al., 2019; Patel et al., 2008; Giray, 2021). In addition, quality aspects such as explainability and

trust-worthiness are new desired quality aspects especially applicable for ML code where current practices and knowledge is inadequate (Giray, 2021).

Mitigation: The ad-hoc trial and error ML development can be addressed by improved tools and techniques. Even though the variety of ML development environments including managed services such as AWS Sagemaker and Google Notebooks attempt to make ML development easier, they essentially do not offer much help in reducing the ad-hoc nature of the development. A significant research push from the community would make ML development relatively systematic and organized.

Recent advancements in the form of available tools not only help a developer to comprehend the process but also let them effectively manage code, data, and experimental results. Examples of such tools and methods include DARVIZ (Sankaran et al., 2017) for DL model visualization, MLFlow¹⁰ for managing the ML lifecycle, and DeepFault (Eniser et al., 2019) for identifying faults in DL programs. Such efforts are expected to address the challenge.

Software Engineering for Machine Learning (SE4ML) brings another perspective to this issue by bringing best practices from software engineering to ML development. Efforts in this direction not only can make ML specific code maintainable and reliable but also can contribute back to reproducibility and replicability.

- Data privacy and bias:** Data hungry ML models are considered as good as the data they are consuming. Data collection and preparation without data diversity leads to bias and unfairness. Although we are witnessing more efforts to understand these sensitive aspects (Zhang and Harman, 2021; Brun and Meliou, 2018), the present set of methods and practices lack the support to deal with data privacy issues at large as well as data diversity and fairness (Brun and Meliou, 2018; Giray, 2021).

Mitigation: Data standards and best practices focusing on data privacy could be considered as an evaluation criterion to mitigate issues concerning data privacy and bias. In addition, mitigation of the issue is also linked with appropriate data pre-processing. Adoption of effective anonymization techniques and data quality assurance practices will further help us deal with the concern.

- Effective feature engineering:** Features represent the problem-specific knowledge in pieces extracted from the data; the effectiveness of any ML model depends on the features fed into it. Many studies identified the importance of effective feature engineering and the challenges in gathering the same (Tsintzira et al., 2020; Shen and Chen, 2020; Patel et al., 2008; Wan et al., 2019; Ivers et al., 2019). Specifically, software engineering researchers have notified that identifying and extracting relevant features beyond code quality metrics is non-trivial. For example, Ivers et al. (2019) discusses that identifying features that establishes a relationship among different code elements is a significant challenge for ML implementations applied on source code analysis. Sharma et al. (2021) have shown in their study that smell detection using ML techniques perform poorly especially for design smells where multiple code elements and their properties has to be observed.

⁹ <https://2021.esec-fse.org/track/fse-2021-artifacts>.

¹⁰ <https://mlflow.org/>.

Mitigation: Recent advancements in the field of large language models (LLMs) trained on huge corpus of code and text have significantly eased the task for researchers. For example, tasks such as generating code embeddings and fine-tuning are supported natively by the LLMs. However, encoding code features specific to downstream tasks is required often and making the task easier requires a significant push from the research community.

- **Skill gap:** Wan et al. (2019) identified that ML software development requires an extended set of skills beyond software development including ML techniques, statistics, and mathematics apart from the application domain. Similarly, Hall and Bowes (2012) also reports a serious lack of ML expertise in academic software engineering efforts. Other authors (Patel et al., 2008) have emphasized the importance of domain knowledge to design effective ML models.

Mitigation: Raising awareness and training sessions customized for the audience is considered the mitigation strategy for skill gap. Software engineering conferences organize tutorials that typically helps new researchers in the field. Availability of various hands-on courses and lecture series from known universities also help bringing the gap.

- **Hardware resources:** Given the need of large training datasets and many hidden layers, often ML training requires high-end processing units (such as GPUs and memory) (Wan et al., 2019; Giray, 2021). A user-survey study (Wan et al., 2019) highlights the need to special hardware for ML training. Such requirements poses a challenge to researchers constrained with limited hardware resources.

Mitigation: ML development is resource hungry. Certain DL models (such as models based on RNN) consume excessive hardware resources. The need for a large-scale hardware infrastructure is increasing with the increase in size of the captured features and the training samples. To address the challenge, infrastructure at institution and country level are maintained in some countries; however, a generic and widely-applicable solution is needed for more globally-inclusive research. Additionally, efforts in the direction of proposed pretrained models, various data pruning techniques, and effective preprocessing techniques are expected to reduce the need of large infrastructure requirements.

6. Threats to validity

The first internal threats to validity relates to the concern of covering all the relevant articles in the selected domain. It is prohibitively time consuming to search each machine learning technique during the literature search. To mitigate the concern, we defined our scope i.e., studies that use ML techniques to solve a software engineering problem by analyzing source code. We also carefully defined inclusion and exclusion criteria for selecting relevant studies. We carry out an extensive manual search process on commonly used digital libraries with the help of a comprehensive set of search terms. Furthermore, we identified a set of frequently occurring keywords in the articles obtained initially for each category individually and carried out another round of literature search with the help of newly identified keywords to enrich the search results.

Another threat to validity is the validity of data extraction and their interpretation applicable to the generated summary and metadata for each selected study. We mitigated this threat by dividing the task of summarization to all the authors and cross verifying the generated information. During the manual summarization phase, metadata of each paper was reviewed by, at least, two authors.

External validity concerns the generalizability and reproducibility of the produced results and observations. We provide a spreadsheet (Sharma et al., 2022) containing all the metadata for all the articles selected in each of the phases of article selection. In addition, inspired by previous surveys (Allamanis et al., 2018a; Hort et al.,

2021), we have developed a website¹¹ as a *living documentation and literature survey* to facilitate easy navigation, exploration, and extension. The website can be easily extended as the new studies emerge in the domain; we have made the repository¹² open-source to allow the community to extend the living literature survey.

7. Conclusions

With the increasing presence of ML techniques in software engineering research, it has become challenging to have a comprehensive overview of its advancements. This survey aims to provide a detailed overview of the studies at the intersection of source code analysis and ML. We have selected 494 studies spanning since 2011 covering 12 software engineering categories. We present a comprehensive summary of the selected studies arranged in categories, subcategories, and their corresponding involved steps. Also, the survey consolidates useful resources (datasets and tools) that could ease the task for future studies. Finally, we present perceived challenges and opportunities in the field. The presented opportunities invite practitioners as well as researchers to propose new methods, tools, and techniques to make the integration of ML techniques for software engineering applications easy, flexible, and maintainable.

Looking ahead: In the recent past, we have witnessed game-changing advancements and all-around adoption of Large language models (LLMs) (Zhao et al., 2023). LLMs such as GPTx (Brown et al., 2020; Radford and Narasimhan, 2018) and BERT (Devlin et al., 2018) learn generic language representation. They help ML models perform better with limited training (i.e., fine-tuning) for a targeted downstream task. Universal contextual representation learned from huge corpora (such as all available textbooks and publicly available articles on the internet) makes them suitable for various natural language tasks.

Similarly, language models for code, such as CodeBERT (Feng et al., 2020), CodeT5 (Wang et al., 2021a), CodeGraphBERT (Guo et al., 2020), and Llama 2 (Touvron et al., 2023) are gaining popularity rapidly among software engineering researchers. Such pre-trained models are trained with generic objectives with large corpora of code and natural language. The models learn the syntax, semantics, and fundamental relationships among the concepts and entities that make fine-tuning the model for a specific software engineering task easier (in terms of training time). These models are not only extensively used in software engineering research (Lu et al., 2021; Chen et al., 2021a; Liu, 2019; Jain et al., 2021; Phan et al., 2021) already but also will be shaping the software engineering research for the years to come.

CRediT authorship contribution statement

Tushar Sharma: Conceptualization, Methodology, Writing – original draft, Investigation. **Maria Kechagia:** Data curation, Writing – original draft. **Stefanos Georgiou:** Data curation. **Rohit Tiwari:** Software, Data curation. **Indira Vats:** Data curation. **Hadi Moazen:** Data curation. **Federica Sarro:** Writing – review & editing.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Data availability

Replication package can be found on GitHub- <https://github.com/tushartushar/ML4SCA>.

¹¹ <http://www.tusharma.in/ML4SCA>.

¹² <https://github.com/tushartushar/ML4SCA>.

Acknowledgments

We would like to thank Mootez Saad and Abhinav Reddy Mandli for helping us categorize the ML techniques. We also thank anonymous reviewers who helped us significantly improve our manuscript. Maria Kechagia and Federica Sarro are supported by the ERC grant no. 741278 (EPIC).

References

- Abbas, Raja, Albaloshi, Fawzi Abdulaziz, Hammad, Mustafa, 2020. Software change proneness prediction using machine learning. In: 2020 International Conference on Innovation and Intelligence for Informatics, Computing and Technologies (3ICT). IEEE, pp. 1–7.
- Abdalkareem, Rabe, Mujahid, Suhaib, Shihab, Emad, 2020. A machine learning approach to improve the detection of ci skip commits. *IEEE Trans. Softw. Eng.*
- Abdeljaber, Osama, Avci, Onur, Kiranyaz, Serkan, Gabbouj, Moncef, Inman, Daniel J., 2017. Real-time vibration-based structural damage detection using one-dimensional convolutional neural networks. *J. Sound Vib.* 388, 154–170.
- Abuhamad, Mohammed, AbuHmed, Tamer, Mohaisen, Aziz, Nyang, DaeHun, 2018. Large-scale and language-oblivious code authorship identification. In: Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security. CCS '18, ISBN: 9781450356930, pp. 101–114. <http://dx.doi.org/10.1145/3243734.3243738>.
- Abunadi, Ibrahim, Alenezi, Mamdouh, 2015. Towards cross project vulnerability prediction in open source web applications. In: Proceedings of the International Conference on Engineering & MIS 2015. ICEMIS '15, Association for Computing Machinery, New York, NY, USA, ISBN: 9781450334181, <http://dx.doi.org/10.1145/2832987.2833051>.
- Aggarwal, Simran, 2019. Software code analysis using ensemble learning techniques. In: Proceedings of the International Conference on Advanced Information Science and System. AISS '19, ISBN: 9781450372916, <http://dx.doi.org/10.1145/3373477.3373486>.
- Agnihotri, Mansi, Chug, Anuradha, 2020. Application of machine learning algorithms for code smell prediction using object-oriented software metrics. *J. Stat. Manag. Syst.* 23 (7), 1159–1171. <http://dx.doi.org/10.1080/09720510.2020.1799576>.
- Ahmad, Wasi, Chakraborty, Saikat, Ray, Baishakhi, Chang, Kai-Wei, 2020. A transformer-based approach for source code summarization. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics. pp. 4998–5007. <http://dx.doi.org/10.18653/v1/2020.acl-main.449>.
- Ahmed, Umar Z., Kumar, Pawan, Karkare, Arme, Kar, Purushottam, Gulwani, Sumit, 2018. Compilation error repair: For the student programs, from the student programs. In: Proceedings of the 40th International Conference on Software Engineering: Software Engineering Education and Training. In: ICSE-SEET '18, ISBN: 9781450356602, pp. 78–87. <http://dx.doi.org/10.1145/3183377.3183383>.
- Al-Jamimi, H.A., Ahmed, M., 2013. Machine learning-based software quality prediction models: State of the art. In: 2013 International Conference on Information Science and Applications (ICISA). pp. 1–4. <http://dx.doi.org/10.1109/ICISA.2013.6579473>.
- Al Qasem, Osama, Akour, Mohammed, Alenezi, Mamdouh, 2020. The influence of deep learning algorithms factors in software fault prediction. *IEEE Access* 8, 63945–63960.
- AL-Shaaby, A., Aljamaan, Hamoud I., Alshayeb, M., 2020. Bad smell detection using machine learning techniques: A systematic literature review. *Arab. J. Sci. Eng.* 45, 2341–2369.
- Alazba, Amal, Aljamaan, Hamoud, 2021. Code smell detection using feature selection and stacking ensemble: An empirical investigation. *Inf. Softw. Technol.* 138, 106648.
- Alleem, Saiqa, Capretz, Luiz Fernando, Ahmed, Faheem, et al., 2015. Comparative performance analysis of machine learning techniques for software bug detection. In: Proceedings of the 4th International Conference on Software Engineering and Applications, number 1. AIRCC Press, Chennai, Tamil Nadu, India, pp. 71–79.
- Aleti, Aldeida, Martinez, Matias, 2021. E-APR: mapping the effectiveness of automated program repair techniques. *Empir. Softw. Eng.* 26 (5), 1–30.
- Alhusain, Sultan, Coupland, Simon, John, Robert, Kavanagh, Maria, 2013. Towards machine learning based design pattern recognition. In: 2013 13th UK Workshop on Computational Intelligence (UKCI). IEEE, pp. 244–251.
- Ali, Nasir, Sharafi, Zohreh, Guhéneuc, Yann-Gaël, Antoniol, Giuliano, 2015. An empirical study on the importance of source code entities for requirements traceability. *Empir. Softw. Eng.* 20 (2), 442–478.
- Ali Alatwi, Huda, Oh, Tae, Fokoue, Ernest, Stackpole, Bill, 2016. Android malware detection using category-based machine learning classifiers. In: Proceedings of the 17th Annual Conference on Information Technology Education. SIGITE '16, ISBN: 9781450344524, pp. 54–59. <http://dx.doi.org/10.1145/2978192.2978218>.
- Alikhashashneh, E.A., Raje, R.R., Hill, J.H., 2018. Using machine learning techniques to classify and predict static code analysis tool warnings. In: 2018 IEEE/ACS 15th International Conference on Computer Systems and Applications (AICCSA). pp. 1–8. <http://dx.doi.org/10.1109/AICCSA.2018.8612819>.
- Aljamaan, Hamoud, Alazba, Amal, 2020. Software defect prediction using tree-based ensembles. In: Proceedings of the 16th ACM International Conference on Predictive Models and Data Analytics in Software Engineering. pp. 1–10.
- Allamanis, Miltiadis, Barr, Earl T., Bird, Christian, Sutton, Charles, 2015a. Suggesting accurate method and class names. In: Proceedings of the 2015 10th Joint Meeting on Foundations of Software Engineering. In: ESEC/FSE 2015, ISBN: 9781450336758, pp. 38–49. <http://dx.doi.org/10.1145/2786805.2786849>.
- Allamanis, Miltiadis, Barr, Earl T., Devanbu, Premkumar, Sutton, Charles, 2018a. A survey of machine learning for big code and naturalness. *ACM Comput. Surv.* (ISSN: 0360-0300) 51 (4), <http://dx.doi.org/10.1145/3212695>.
- Allamanis, Miltiadis, Brookschen, Marc, Khademi, Mahmoud, 2018b. Learning to represent programs with graphs. In: International Conference on Learning Representations.
- Allamanis, Miltiadis, Peng, Hao, Sutton, Charles, 2016. A convolutional attention network for extreme summarization of source code.
- Allamanis, M., Sutton, C., 2013a. Mining source code repositories at massive scale using language modeling. In: 2013 10th Working Conference on Mining Software Repositories (MSR). pp. 207–216. <http://dx.doi.org/10.1109/MSR.2013.6624029>.
- Allamanis, Miltiadis, Sutton, Charles, 2013b. Mining source code repositories at massive scale using language modeling. In: 10th Working Conference on Mining Software Repositories (MSR). pp. 207–216. <http://dx.doi.org/10.1109/MSR.2013.6624029>.
- Allamanis, Miltiadis, Tarlow, Daniel, Gordon, Andrew D., Wei, Yi, 2015b. Bimodal modelling of source code and natural language. In: Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37. ICML '15, pp. 2123–2132.
- Allix, Kevin, Bissyandé, Tegawendé F., Klein, Jacques, Le Traon, Yves, 2016. AndroZoo: Collecting millions of android apps for the research community. In: Proceedings of the 13th International Conference on Mining Software Repositories. MSR '16, ISBN: 978-1-4503-4186-8, pp. 468–471. <http://dx.doi.org/10.1145/2901739.2903508>.
- Alon, Uri, Brody, Shaked, Levy, Omer, Yahav, Eran, 2019a. code2seq: Generating sequences from structured representations of code.
- Alon, Uri, Zilberstein, Meital, Levy, Omer, Yahav, Eran, 2018. A general path-based representation for predicting program properties. *SIGPLAN Not.* (ISSN: 0362-1340) 53 (4), 404–419. <http://dx.doi.org/10.1145/3296979.3192412>.
- Alon, Uri, Zilberstein, Meital, Levy, Omer, Yahav, Eran, 2019b. Code2vec: Learning distributed representations of code. *Proc. ACM Program. Lang.* 3 (POPL), <http://dx.doi.org/10.1145/3290353>.
- Alrajeh, Dalal, Kramer, Jeff, Russo, Alessandra, Uchitel, Sebastian, 2015. Automated support for diagnosis and repair. *Commun. ACM* (ISSN: 0001-0782) 58 (2), 65–72. <http://dx.doi.org/10.1145/2658986>.
- Alsolai, Hadeel, Roper, Marc, 2020. A systematic literature review of machine learning techniques for software maintainability prediction. *Inf. Softw. Technol.* (ISSN: 0950-5849) 119, 106214. <http://dx.doi.org/10.1016/j.infsof.2019.106214>.
- Altarawy, Doaa, Shahin, Hossameldin, Mohammed, Ayat, Meng, Na, 2018. Lascad: Language-agnostic software categorization and similar application detection. *J. Syst. Softw.* 142, 21–34.
- Alves, H., Fonseca, B., Antunes, N., 2016. Experimenting machine learning techniques to predict vulnerabilities. In: 2016 Seventh Latin-American Symposium on Dependable Computing (LADC). pp. 151–156. <http://dx.doi.org/10.1109/LADC.2016.32>.
- Amal, Boukhdhir, Kessentini, Marouane, Bechikh, Slim, Dea, Josselin, Said, Lamjed Ben, 2014. On the use of machine learning and search-based software engineering for ill-defined fitness function: A case study on software refactoring. In: Le Goues, Claire, Yoo, Shin (Eds.), *Search-Based Software Engineering*. ISBN: 978-3-319-09940-8, pp. 31–45.
- Amorim, L., Costa, E., Antunes, N., Fonseca, B., Ribeiro, M., 2015. Experience report: Evaluating the effectiveness of decision trees for detecting code smells. In: 2015 IEEE 26th International Symposium on Software Reliability Engineering (ISSRE). pp. 261–269. <http://dx.doi.org/10.1109/ISSRE.2015.7381819>.
- Amorim, L.A., Freitas, M.F., Dantas, A., de Souza, E.F., Camilo-Junior, C.G., Martins, W.S., 2018. A new word embedding approach to evaluate potential fixes for automated program repair. In: 2018 International Joint Conference on Neural Networks (IJCNN). pp. 1–8. <http://dx.doi.org/10.1109/IJCNN.2018.8489079>.
- Aniche, M., Maziero, E., Durelli, R., Durelli, V., 2020. The effectiveness of supervised machine learning algorithms in predicting software refactoring. *IEEE Trans. Softw. Eng.* 1. <http://dx.doi.org/10.1109/TSE.2020.3021736>.
- Arar, Ömer Faruk, Ayan, K'urşat, 2015. Software defect prediction using cost-sensitive neural network. *Appl. Soft Comput.* 33, 263–277.
- Arcelli Fontana, Francesca, Zanoni, Marco, 2017. Code smell severity classification using machine learning techniques. *Knowl.-Based Syst.* (ISSN: 0950-7051) 128, 43–58. <http://dx.doi.org/10.1016/j.knosys.2017.04.014>.
- Aribandi, Vamsi Krishna, Kumar, Lov, Bhanu Murthy Neti, Lalita, Krishna, Aneesh, 2019. Prediction of refactoring-prone classes using ensemble learning. In: Gedeon, Tom, Wong, Kok Wai, Lee, Minho (Eds.), *Neural Information Processing*. ISBN: 978-3-03-36802-9, pp. 242–250.
- Azcona, David, Arora, Piyush, Hsiao, I-Han, Smeaton, Alan, 2019. User2code2vec: Embeddings for profiling students based on distributional representations of source code. In: Proceedings of the 9th International Conference on Learning Analytics & Knowledge. In: LAK19, ISBN: 9781450362566, pp. 86–95. <http://dx.doi.org/10.1145/3303772.3303813>.

- Chen, Qiuyuan, Hu, Han, Liu, Zhaoyi, 2019b. Code summarization with abstract syntax tree. In: Gedeon, Tom, Wong, Kok Wai, Lee, Minho (Eds.), *Neural Information Processing*. ISBN: 978-3-030-36802-9, pp. 652–660.
- Chen, Jinyin, Hu, Keke, Yu, Yue, Chen, Zhuangzhi, Xuan, Qi, Liu, Yi, Filkov, Vladimir, 2020a. Software visualization and deep transfer learning for effective software defect prediction. In: Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering. ICSE '20, ISBN: 9781450371216, pp. 578–589. <http://dx.doi.org/10.1145/3377811.3380389>.
- Chen, Fuxiang, Kim, Mijung, Choo, Jaegul, 2021a. Novel natural language summarization of program code via leveraging multiple input representations. In: Findings of the Association for Computational Linguistics: EMNLP 2021. Association for Computational Linguistics, Punta Cana, Dominican Republic, pp. 2510–2520. <http://dx.doi.org/10.18653/v1/2021.findings-emnlp.214>, URL <https://aclanthology.org/2021.findings-emnlp.214>.
- Chen, Z., Kommarusch, S.J., Tufano, M., Pouchet, L., Poshyvanyk, D., Monperrus, M., 2019c. SEQUENCER: Sequence-to-sequence learning for end-to-end program repair. *IEEE Trans. Softw. Eng.* 1. <http://dx.doi.org/10.1109/TSE.2019.2940179>.
- Chen, Xinyun, Liu, Chang, Shin, Richard, Song, Dawn, Chen, Mingcheng, 2016. Latent attention for if-then program synthesis. In: Proceedings of the 30th International Conference on Neural Information Processing Systems. NIPS '16, ISBN: 9781510838819, pp. 4581–4589.
- Chen, Xinyun, Liu, Chang, Song, Dawn, 2018. Towards synthesizing complex programs from input-output examples.
- Chen, Xinyun, Liu, Chang, Song, Dawn, 2019d. Execution-guided neural program synthesis. In: International Conference on Learning Representations.
- Chen, Yang, Santosa, Andrew E., Yi, Ang Ming, Sharma, Abhishek, Sharma, Asankhya, Lo, David, 2020b. A machine learning approach for vulnerability curation. In: Proceedings of the 17th International Conference on Mining Software Repositories. Association for Computing Machinery, New York, NY, USA, ISBN: 9781450375177, pp. 32–42, URL <https://doi.org/10.1145/3379597.3387461>.
- Chen, Mark, Tworek, Jerry, Jun, Heewoo, Yuan, Qiming, Pinto, Henrique Ponde de Oliveira, Kaplan, Jared, Edwards, Harri, Burda, Yuri, Joseph, Nicholas, Brockman, Greg, et al., 2021b. Evaluating large language models trained on code. arXiv preprint [arXiv:2107.03374](https://arxiv.org/abs/2107.03374).
- Chen, M., Wan, X., 2019. Neural comment generation for source code with auxiliary code classification task. In: 2019 26th Asia-Pacific Software Engineering Conference (APSEC). pp. 522–529. <http://dx.doi.org/10.1109/APSEC48747.2019.00076>.
- Chen, Qiuyuan, Xia, Xin, Hu, Han, Lo, David, Li, Shaping, 2021c. Why my code summarization model does not work: Code comment improvement with category prediction. *ACM Trans. Softw. Eng. Methodol. (TOSEM)* 30 (2), 1–29.
- Chen, Long, Ye, Wei, Zhang, Shikun, 2019e. Capturing source code semantics via tree-based convolution over API-enhanced AST. In: Proceedings of the 16th ACM International Conference on Computing Frontiers. CF '19, ISBN: 9781450366854, pp. 174–182. <http://dx.doi.org/10.1145/3310273.3321560>.
- Chen, Q., Zhou, M., 2018. A neural framework for retrieval and summarization of source code. In: 2018 33rd IEEE/ACM International Conference on Automated Software Engineering (ASE). pp. 826–831. <http://dx.doi.org/10.1145/3238147.3240471>.
- Chernis, Boris, Verma, Rakesh, 2018. Machine learning methods for software vulnerability detection. In: Proceedings of the Fourth ACM International Workshop on Security and Privacy Analytics. IWSPA '18, ISBN: 9781450356343, pp. 31–39. <http://dx.doi.org/10.1145/3180445.3180453>.
- Chidamber, S.R., Kemerer, C.F., 1994. A metrics suite for object oriented design. *IEEE Trans. Softw. Eng. (ISSN: 0098-5589)* 20 (6), 476–493. <http://dx.doi.org/10.1109/32.295895>.
- Choi, Y., Kim, S., Lee, J., 2020. Source code summarization using attention-based keyword memory networks. In: 2020 IEEE International Conference on Big Data and Smart Computing (BigComp). pp. 564–570. <http://dx.doi.org/10.1109/BigComp48618.2020.00011>.
- Choudhary, Garvit Rajesh, Kumar, Sandeep, Kumar, Kuldeep, Mishra, Alok, Catal, Cagatay, 2018. Empirical analysis of change metrics for software fault prediction. *Comput. Electr. Eng.* 67, 15–24.
- Chug, A., Dhall, S., 2013. Software defect prediction using supervised learning algorithm and unsupervised learning algorithm. In: Confluence 2013: The Next Generation Information Technology Summit (4th International Conference). pp. 173–179. <http://dx.doi.org/10.1049/cp.2013.2313>.
- Clemente, C.J., Jaafar, F., Malik, Y., 2018. Is predicting software security bugs using deep learning better than the traditional machine learning algorithms? In: 2018 IEEE International Conference on Software Quality, Reliability and Security (QRS). pp. 95–102. <http://dx.doi.org/10.1109/QRS.2018.00023>.
- Compton, Rhys, Frank, Eibe, Patros, Panos, Koay, Abigail, 2020. Embedding java classes with code2vec: Improvements from variable obfuscation. In: Proceedings of the 17th International Conference on Mining Software Repositories. MSR '20, ISBN: 9781450375177, pp. 243–253. <http://dx.doi.org/10.1145/3379597.3387445>.
- Cortes-Coy, Luis Fernando, Vásquez, M., Aponte, Jairo, Poshyvanyk, D., 2014. On automatically generating commit messages via summarization of source code changes. In: 2014 IEEE 14th International Working Conference on Source Code Analysis and Manipulation. pp. 275–284.
- Cruz, Daniel, Santana, Amanda, Figueiredo, Eduardo, 2020a. Detecting bad smells with machine learning algorithms: an empirical study. In: Proceedings of the 3rd International Conference on Technical Debt. pp. 31–40.
- Cruz, Daniel, Santana, Amanda, Figueiredo, Eduardo, 2020b. Detecting bad smells with machine learning algorithms: An empirical study. In: Proceedings of the 3rd International Conference on Technical Debt. TechDebt '20, ISBN: 9781450379601, pp. 31–40. <http://dx.doi.org/10.1145/3387906.3388618>.
- Cui, Jianfeng, Wang, Lixin, Zhao, Xin, Zhang, Hongyi, 2020. Towards predictive analysis of android vulnerability using statistical codes and machine learning for IoT applications. *Comput. Commun. (ISSN: 0140-3664)* 155, 125–131. <http://dx.doi.org/10.1016/j.comcom.2020.02.078>.
- Cummins, C., Petoumenos, P., Wang, Z., Leather, H., 2017. Synthesizing benchmarks for predictive modeling. In: 2017 IEEE/ACM International Symposium on Code Generation and Optimization (CGO). pp. 86–99. <http://dx.doi.org/10.1109/CGO.2017.7863731>.
- Cunha, Warteruzannan Soyer, Armijo, Guisella Angulo, de Camargo, Valter Vieira, 2020. Investigating non-usually employed features in the identification of architectural smells: A machine learning-based approach. In: Proceedings of the 14th Brazilian Symposium on Software Components, Architectures, and Reuse. ISBN: 9781450387545, pp. 21–30.
- Cvitkovic, Milan, Singh, Badal, Anandkumar, Animashree, 2019. Open vocabulary learning on source code with a graph-structured cache. In: Chaudhuri, Kamalika, Salakhutdinov, Ruslan (Eds.), In: Proceedings of Machine Learning Research, vol. 97, pp. 1475–1485.
- Dam, Hoa Khanh, Pham, Trang, Ng, Shien Wee, Tran, Truyen, Grundy, John, Ghose, Aditya, Kim, Taeksu, Kim, Chul-Joo, 2019. Lessons learned from using a deep tree-based model for software defect prediction in practice. In: Proceedings of the 16th International Conference on Mining Software Repositories. MSR '19, pp. 46–57. <http://dx.doi.org/10.1109/MSR.2019.00017>.
- D'Ambros, Marco, Lanza, Michele, Robbes, Romain, 2012. Evaluating defect prediction approaches: A benchmark and an extensive comparison. *Empir. Softw. Eng. (ISSN: 1382-3256)* 17 (4–5), 531–577. <http://dx.doi.org/10.1007/s10664-011-9173-9>.
- Dantas, Altino, de Souza, Eduardo F., Souza, Jerffeson, Camilo-Junior, Celso G., 2019. Code naturalness to assist search space exploration in search-based program repair methods. In: Nejati, Shiva, Gay, Gregory (Eds.), *Search-Based Software Engineering*. ISBN: 978-3-030-27455-9, pp. 164–170.
- De Lucia, Andrea, Di Penta, Massimiliano, Oliveto, Rocco, Panichella, Annibale, Panichella, Sebastiano, 2014. Labeling source code with information retrieval methods: an empirical study. *Empir. Softw. Eng.* 19 (5), 1383–1420.
- Dejaeger, Karel, Verbraeken, Thomas, Baesens, Bart, 2012. Toward comprehensible software fault prediction models using bayesian network classifiers. *IEEE Trans. Softw. Eng.* 39 (2), 237–257.
- Devlin, Jacob, Bapna, Rudy, Singh, Rishabh, Hauseknecht, Matthew, Kohli, Pushmeet, 2017a. Neural program meta-induction. In: Proceedings of the 31st International Conference on Neural Information Processing Systems. NIPS '17, ISBN: 9781510860964, pp. 2077–2085.
- Devlin, Jacob, Chang, Ming-Wei, Lee, Kenton, Toutanova, Kristina, 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint [arXiv:1810.04805](https://arxiv.org/abs/1810.04805).
- Devlin, Jacob, Uesato, Jonathan, Bhupatiraju, Surya, Singh, Rishabh, Mohamed, Abdellrahman, Kohli, Pushmeet, 2017b. RobustFill: Neural program learning under noisy I/O. In: Proceedings of the 34th International Conference on Machine Learning - Volume 70. ICML '17, pp. 990–998.
- Dewanjan, Seema, Rao, Rajwant Singh, Mishra, Alok, Gupta, Manjari, 2021. A novel approach for code smell detection: An empirical study. *IEEE Access* 9, 162869–162883.
- Dhamayanthi, N., Lavanya, B., 2019. Improvement in software defect prediction outcome using principal component analysis and ensemble machine learning algorithms. In: Hemanth, Jude, Fernando, Xavier, Lafata, Pavel, Baig, Zubair (Eds.), *International Conference on Intelligent Data Communication Technologies and Internet of Things (ICICI) 2018*. ISBN: 978-3-03-03146-6, pp. 397–406.
- Di Martino, Sergio, Ferrucci, Filomena, Gravino, Carmine, Sarro, Federica, 2011. A genetic algorithm to configure support vector machines for predicting fault-prone components. In: Caivano, Danilo, Oivo, Markku, Baldassarre, Maria Teresa, Visaggio, Giuseppe (Eds.), *Product-Focused Software Process Improvement*. Springer Berlin Heidelberg, Berlin, Heidelberg, ISBN: 978-3-642-21843-9, pp. 247–261.
- Di Nucci, D., Palomba, F., Tamburri, D.A., Serebrenik, A., De Lucia, A., 2018. Detecting code smells using machine learning techniques: Are we there yet? In: 2018 IEEE 25th International Conference on Software Analysis, Evolution and Reengineering (SANER). pp. 612–621. <http://dx.doi.org/10.1109/SANER.2018.8330266>.
- Dong, Li, Lapata, Mirella, 2016. Language to logical form with neural attention. In: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). Association for Computational Linguistics, Berlin, Germany, pp. 33–43. <http://dx.doi.org/10.18653/v1/P16-1004>, URL <https://aclanthology.org/P16-1004>.
- Dos Santos, Geanderson Esteves, Figueiredo, E., Veloso, Adriano, Viggiani, Markos, Ziviani, N., 2020. Understanding machine learning software defect predictions. *Autom. Softw. Eng.* 27, 369–392.
- Du, Xiaoning, Chen, Bihuan, Li, Yuekang, Guo, Jianmin, Zhou, Yaqin, Liu, Yang, Jiang, Yu, 2019. LEOPARD: Identifying vulnerable code for vulnerability assessment through program metrics. In: 2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE). pp. 60–71. <http://dx.doi.org/10.1109/ICSE.2019.00024>.

