

Code Smells in Machine Learning Systems

Jiri Gesi
fjiriges@uci.edu
University of California, Irvine
Irvine, California, USA

Siqi Liu
sliu17@uci.edu
University of California, Irvine
Irvine, California, USA

Jiawei Li
jiaw128@uci.edu
University of California, Irvine
Irvine, California, USA

Iftekhhar Ahmed
iftekha@uci.edu
University of California, Irvine
Irvine, California, USA

Nachiappan Nagappan
nnagappan@acm.org
Facebook
Seattle, WA, USA

David Lo
davidlo@smu.edu.sg
Singapore Management University
Singapore

Eduardo Santana de Almeida
esa@rise.com.br
Federal University of Bahia
Salvador, Brazil

Pavneet Singh Kochhar
pavneet.kochhar@microsoft.com
Microsoft Research
Vancouver, Canada

Lingfeng Bao
fjiriges@uci.edu
Zhejiang University
Zhejiang, China

ABSTRACT

As Deep learning (DL) systems continuously evolve and grow, assuring their quality becomes an important yet challenging task. Compared to non-DL systems, DL systems have more complex team compositions and heavier data dependency. These inherent characteristics would potentially cause DL systems to be more vulnerable to bugs and, in the long run, to maintenance issues. Code smells are empirically tested as efficient indicators of non-DL systems. Therefore, we took a step forward into identifying code smells, and understanding their impact on maintenance in this comprehensive study. This is the first study on investigating code smells in the context of DL software systems, which helps researchers and practitioners to get a first look at what kind of maintenance modification made and what code smells developers have been dealing with. Our paper has three major contributions. First, we comprehensively investigated the maintenance modifications that have been made by DL developers via studying the evolution of DL systems, and we identified nine frequently occurred maintenance-related modification categories in DL systems. Second, we summarized five code smells in DL systems. Third, we validated the prevalence, and the impact of our newly identified code smells through a mixture of qualitative and quantitative analysis. We found that our newly identified code smells are prevalent and impactful on the maintenance of DL systems from the developer's perspective.

KEYWORDS

Deep learning, Code Smell, Code Quality, Empirical analysis

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Arxiv,

© 2022 Association for Computing Machinery.
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

ACM Reference Format:

Jiri Gesi, Siqi Liu, Jiawei Li, Iftekhhar Ahmed, Nachiappan Nagappan, David Lo, Eduardo Santana de Almeida, Pavneet Singh Kochhar, and Lingfeng Bao. 2022. Code Smells in Machine Learning Systems. In *Proceedings of (Arxiv)*. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

In the past few years, Deep Learning (DL) systems, a branch of machine learning (ML), has become an inseparable part of billions of peoples' lives worldwide, from personal banking to communication, from entertainment to transportation, and more [2, 4]. Due to such ever-increasing dependence, ensuring DL system quality is of utmost importance. Failure to do so has already resulted in catastrophic consequences [1].

As DL systems evolve and grow in size and complexity, continuous maintenance in the form of performance improvement, mandatory upgrades, and fixing bugs is necessary to ensure its correctness and continuous availability during its lifetime [44]. However, maintenance of DL systems, similar to non-ML systems, can be hindered due to poor design and implementation choices. Compared to non-ML systems, DL systems are even more affected by maintenance issues since they are prone to the maintenance issues pertaining to both non-ML components and DL components as DL systems are combinations of both non-ML and DL components [42].

Over the years, researchers have investigated the indicators of maintenance issues and methods to identify and quantify their impact [11, 23, 25, 42]. Code smells is one such indicator, which is related to long term maintenance issues [23, 25, 42]. Prior research have investigated when and why code smells are introduced [46] and how they evolve over time [13, 17, 35, 39, 46]. Code smells' impact on software comprehensibility [10], fault-proneness, change-proneness [20, 28, 29], and code maintainability [22, 34, 43, 48, 49] has also been demonstrated.

However, the majority of these studies focus on non-ML code smells with only a few focusing on ML code smells [25] and none focuses on DL-specific code smells. Since DL and traditional software development is significantly different in terms of workflow and engineering practices [21], as well as DL's data dependent

behavior [12, 47], it is safe to assume that along with previously known code smells, there are code smells that are unique to DL systems which have not yet identified.

A study conducted by Hadhemi et al. [25] is closest to our work, where they studied code smells in DL systems. However, they investigated the prevalence of Python code smells; and analyzed code smells that were designed for non-DL general-purpose source code [18]. We posit that generic Python code smells provide only a partial picture, and there are DL-specific code smells that require further investigation. For example, Fig.1 shows an example of Jumbled Model Architecture (JMA) code smell where a Variational Autoencoder (VAE) [30] is extracted into encoding, sampling, and decoding¹. Intuitively, jumbled VAE impedes the understandability of model architecture and makes future maintenance difficult, this refactoring helps to ameliorate that. Due to the already proven impact of code smells on various aspects of non-DL software, it is safe to assume that code smells will have a similar, if not more detrimental effect on the long-term maintainability and overall quality of DL systems. Making it of utmost importance to get a complete picture of the unique code smells in DL systems and understand their impact. The first step towards achieving this goal is by identifying DL-specific code smells derived from real-world modifications applied to DL projects by developers.

In this study, we identify and analyze maintenance related modifications done by developers on 59 open source DL projects that were previously investigated by other researchers [25].

By employing a combination of PythonChangeMiner [7], GitProc [16] and manual analysis, we collected 426 maintenance related code changes from these 59 projects, where each change has at least three other similar occurrences among the projects. Next, using qualitative analysis, multiple coders independently coded collected changes into nine groups and extracted five frequently occurring code smells. Next, we validated the prevalence and severity of code smells by conducting a survey of 235 OSS DL developers. The survey analysis results show that our identified new code smells are often seen and have a significant impact on system maintenance activities.

In this paper, we answer the following research questions:

RQ1: What kinds of modifications do developers make frequently in DL systems?

RQ2: How prevalent are code smells in DL systems?

RQ3: How do practitioners perceive the identified code smells in DL systems?

The remainder of the paper is structured as follows. Section 2 provides an overview of related work. Section 3 details our methodology, with Section 4 presenting our findings. Section 5 places our results in the broader context of work to date and outlines the implications for DL practitioners and researchers. Section 6 lists the threats to the validity of our results. Section 7 concludes with a summary of the key findings and an outlook on our future work.

2 RELATED WORK

Code smells were introduced by Martin Fowler [23] to describe the design and implementation flaws in source code. These flaws

do not make the software system behave incorrectly or crash but make it harder to understand, and maintain [20]. Research communities have investigated the impact of code smells in non-ML software systems such as how code smells impact fault-proneness and change-proneness [20, 28, 29], it's impact on maintainability [22, 34, 43, 48, 49], when and why code smells are introduced [46], how they evolve over time [13, 17, 35, 39, 46], and how to detect code smells using different techniques [33, 36, 37, 40].

However, whether these code smells can capture all code smells relevant to DL systems is still an open question since existing research shows that there are significant differences between DL and traditional software systems. Wan et al. showed that the incorporation of DL into a software system significantly impacts the requirement analysis, system design, testing, and process management [47]. Scully et al. presented a set of unique anti-patterns in DL system development and highlighted a number of areas where technical debts unique to DL systems exist [42]. Researchers also identified differences in the development process for DL systems due to the team formation and dependence on data which necessitates steps such as data understanding, data cleaning, model training, model deployment, and monitoring [3, 9, 12, 21]. All these differences can potentially introduce unique poor designs or implementations in source code, also known as code smells.

Despite the clear differences between DL and traditional software systems, only a few studies have investigated code smells in the context of DL systems. Hadhemi et al. [25] investigated the prevalence of Python code smells in DL systems along with investigating the differences in the distribution of code smells between DL and traditional systems. The code smells they investigated are:

Long Parameter List (LPL) [23]: A method or a function that has a large number of parameters.

Long Method (LM) [23]: A method or a function that is extremely long.

Long Scope Chaining (LSC) [19]: A method or a function that has a deeply nested closure.

Large Class (LC) [15]: A class that has a large number of source code lines.

Long Message Chain (LMC) [15]: An expression for accessing an object using the dot operators through a long sequence of attributes or method calls.

Long Base Class List (LBCL) [15]: When a class extends too many base classes due to the multiple inheritances that Python language supports, it makes code hard to understand.

Long Lambda Function (LLF) [15]: An anonymous function that is extremely long and complex in terms of conditions and parameters.

Long Ternary Conditional Expression (LTCE) [15]: A ternary conditional expression that is extremely long.

Complex Container Comprehension (CCC) [15]: One-line comprehension list, set, or dictionary that contains a large number of clauses and filter expressions.

Multiply-Nested Container (MNC) [15]: a container (including set, list, tuple, dict) that is deeply nested.

As it can be seen for the definitions, these code smells were designed for traditional general-purpose Python code [18]. However, in a DL system, there is general-purpose code, along with model architecture, data preparation, and pipeline related code. Hence, we

¹This refactoring commit is collected from the "NiftyNet" [5] open-source software project.

```

class VAE(nn.Module):
    def __init__(self):
        super(VAE, self).__init__()

        self.fc1 = nn.Linear(784, 400)
        self.fc21 = nn.Linear(400, 20)
        self.fc22 = nn.Linear(400, 20)
        self.fc3 = nn.Linear(20, 400)
        self.fc4 = nn.Linear(400, 784)

    def forward(self, x):
        h1 = F.relu(self.fc1(x.view(-1, 784)))
        mu, logvar = self.fc21(h1), self.fc22(h1)

        std = torch.exp(0.5*logvar)
        eps = torch.randn_like(std)
        z = mu + eps*std

        h3 = F.relu(self.fc3(z))
        decodedZ = torch.sigmoid(self.fc4(h3))

        return decodedZ, mu, logvar

class VAE(nn.Module):
    def __init__(self):
        super(VAE, self).__init__()

        self.fc1 = nn.Linear(784, 400)
        self.fc21 = nn.Linear(400, 20)
        self.fc22 = nn.Linear(400, 20)
        self.fc3 = nn.Linear(20, 400)
        self.fc4 = nn.Linear(400, 784)

    def encode(self, x):
        h1 = F.relu(self.fc1(x))
        return self.fc21(h1), self.fc22(h1)

    def sample(self, mu, logvar):
        std = torch.exp(0.5*logvar)
        eps = torch.randn_like(std)
        return mu + eps*std

    def decode(self, z):
        h3 = F.relu(self.fc3(z))
        return torch.sigmoid(self.fc4(h3))

    def forward(self, x):
        mu, logvar = self.encode(x.view(-1, 784))
        z = self.sample(mu, logvar)
        return self.decode(z), mu, logvar

```

Figure 1: Jumbled Model Architecture code smell refactoring for Variational Autoencoder model

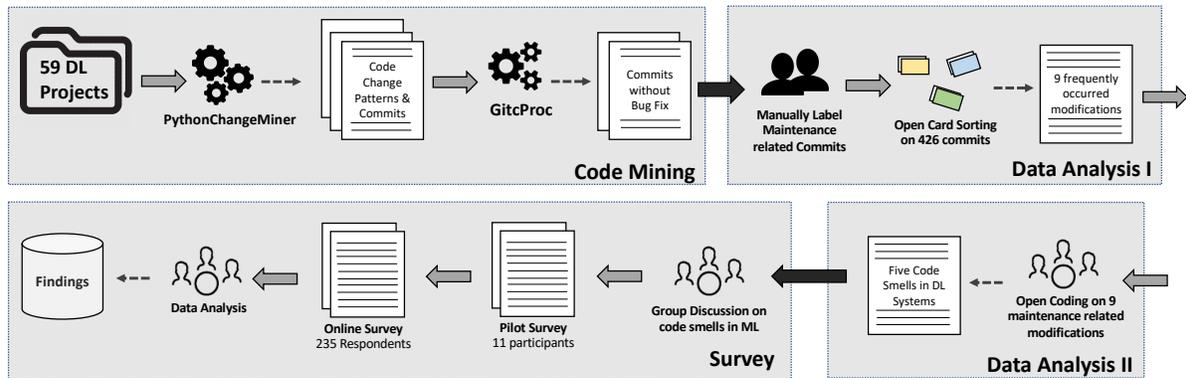


Figure 2: Schematic diagram of research methodology

posit that there are other code smells that are unique to DL specific code (i.e., model architecture, data preparation and pipeline, etc.)

Prior research in the context of non-ML systems indicated that developers have differentiated opinions about code smells, their prevalence, and effect [49]. However, existing research in DL did not investigate how developers perceive code smells in the context of DL systems. As result, questions such as how prevalent these smells are, and how developers perceive their impact remains unanswered. We aim to fill the gap and answer the questions in this work.

3 METHODOLOGY

We used a mixed method approach consisting of mining software repositories and qualitative analysis. Figure 2 shows the process that we follow in this study. We start by code mining to gather recurring code change patterns, then apply open card coding to identify new code smells, and finally, conduct a large-scale survey

to validate the prevalence and impact of the newly identified code smells.

3.1 Code Mining

Our first step was collecting recurring code changes in 59 open source DL systems. These projects were investigated by Hadhemi et al. [25] in their study and we wanted to investigate whether there are other codes smells unique to DL in these systems besides generic Python code smells, thus we used the same dataset.

3.1.1 Data Collection. We started by obtaining 90,301 commits from the 59 DL open source projects downloaded on May 20, 2020. Next, we used PythonChangeMiner [7] to detect and group commits with similar change patterns. PythonChangeMiner mines the history of a given repository using the PyDriller framework [6] and builds change graphs for matching functions in each changed file for a commit. To achieve this, both versions of the file (before and after the change) are parsed into Abstract Syntax Trees (ASTs) [38],

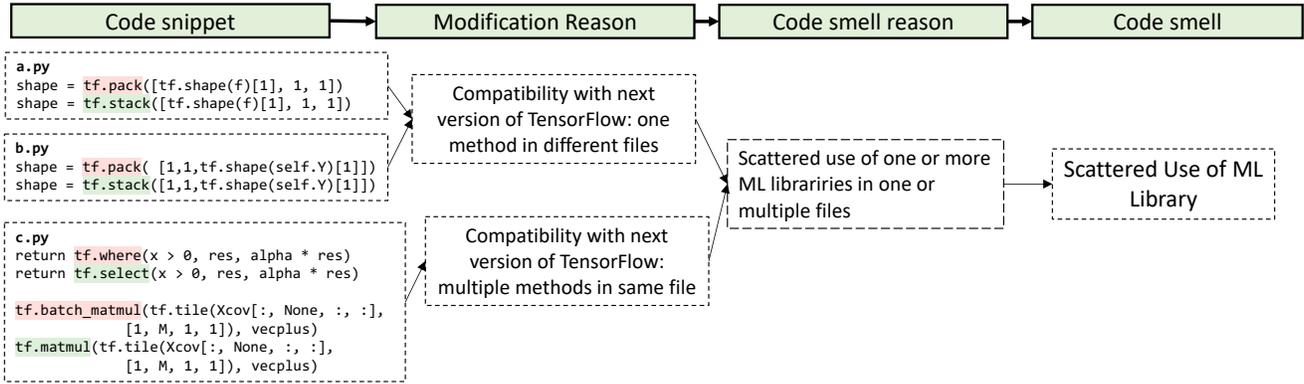


Figure 3: An example of using qualitative analysis for determining code smells in deep learning systems

which are then traversed to create the structure of *fine-grained Program Dependence Graph* (*fgPDG*). Then, the obtained *fgPDG* are analyzed to find all node pairs before and after the change using GumTree [31], resulting in grouped change pattern categories. Figure 4 shows an example of a changing pattern identified in several projects that developers switched from using built-in copying to creating a deep copy of an object using a *copy* module. To make sure that our analyzed patterns are common across multiple projects and not specific to a project, we extracted code changes that happened at least three times within all commits across multiple projects. We identified 1,942 commits matching this criterion.

```

model = optimizer.target.copy()
model = copy.deepcopy(optimizer.target)

(a) Code snippet 1

merged_config = self._default_config.copy()
merged_config = copy.deepcopy(self._default_config)

(b) Code snippet 2

votes = rule.get_class_votes(X, self).copy()
votes = copy.deepcopy(rule.get_class_votes(X, self))

(c) Code snippet 3
  
```

Figure 4: An example of a changing pattern identified in several projects on GitHub. The developers switched from using built-in copying to creating a deep copy of an object using a *copy* module of a Standard library.

Research shows that refactoring (non-bug-fixing and non-program behavior-altering commits) is performed to remove code smells [24]. Since our identified 1,942 patterns contained both bug-fixing, non-bug-fixing commits, we removed bug-fixing commits from our analysis as they alter program behavior. We used GitCProc [16] for this purpose, which identifies the bug-fixing commits based on the presence of specific words in the commit message. Words such as *error*, *bug*, *defect* and *fix* are considered while identifying bug-fix commits by GitCProc. After removing bug-fix related commits, 1,335 non-bug-fixing commits were left, which come from all 59 projects. Next, the first and second authors independently went through the commits to identify the commits related to maintenance. They relied on the commit message and compared the code before and after the update for deciding whether the commit

was maintenance related or not. They initially used 10% (134) of the commits and independently labeled them. After initial labeling, the inter-rater agreement was 0.61, which according to Landis et al.[32] is considered as a substantial level of agreement. After an initial disagreement on some of the commits, the authors discussed their approach and had a complete agreement regarding the label of commits initially disagreed. Then the two authors labeled the remaining 1,201 commits together. This resulted in selecting 426 maintenance related commits where each commit had at least three occurrences across multiple projects.

3.1.2 Modification Category Creation. Our next step was to group these commits based on the modification reasons. To do so, we followed descriptive coding [41] which is used for identifying topics from data. The result of descriptive coding is categorized groups based on identified topics. Two authors jointly conducted the descriptive coding on the selected 426 commits. They relied on the commit message and compared the code before and after the update for identifying the reasons for making the changes. This resulted in grouping the commits into nine modification categories. We selected descriptive coding technique for the following reasons: (1) we can get an overview of recurring changes that are indicative of poor maintainability; (2) we can obtain the context of these modifications.

3.1.3 Code Smell Categorization. Our primary goal was to extract code smells from the frequently occurred modifications. For this purpose, in the next step first and the second authors checked if the modification reasons mentioned in Table 1 met the following criterion: (1) whether the modification reason is general (common to many DL systems), (2) if there is a general solution to the root cause that required the modification. If both criteria are met, they considered the modification reason as a code smell. Figure 3 shows an example of qualitative analysis to determine whether a modification is a code smell.

Two rounds of descriptive coding were conducted. In the first round, the first and second authors independently investigated all modification reasons and created a list of code smell candidates based on the previously mentioned two criteria. After discussing they curated a list of 12 code smells and reached an inter-rater agreement of 83.2%. In the second round, these 12 code smells were

presented to all authors, and after discussion, everyone agreed on five new code smells and the remaining seven were discarded as they did not meet the previously mentioned criteria completely.

Since the collected commits consisted of both new code smells and pre-existing code smells that were identified by Hadhemi et al [25] (listed in Section 2), we check the prevalence of generic Python code smells among the 427 commits. Since Pysmell [18] can identify these smells, instead of applying quantitative techniques, we relied on Pysmell for this purpose. We ran Pysmell before and after applying each of the 426 commits and calculated the number of fixed generic Python code smells. If the count of code smells decreased, we labeled that commit as Python code smell fixing commit. Through this analysis, we identified eight Python code smells in our data.

3.2 Survey

We delivered a survey to gain an understanding of the prevalence and severity of the newly identified code smells and gather the developer’s perspective about them.

3.2.1 Protocol. We based our questions on the identified code smells from code change pattern mining. Our questionnaire included questions about the following topics (the complete questionnaire is available as supplemental material²):

- **Demographics:** We asked questions about organizations, geographical locations, and ML-related working experiences for this part of the survey.
- **Self-perception:** We let respondents self-identify their professional categories (“I think of myself as a/an...” like researchers, engineers, scientists, etc). We used the answers to classify all respondents into four groups based on the result’s keywords: data scientist/engineer, Machine Learning (ML) engineer, software engineer, and project manager based on their self-perception. ML engineers sit at the intersection of software engineering and data science, whose job is applying ML techniques and developing DL models. Data scientists/engineers are the group of people who create and maintain optimal data pipeline architecture, study and understand the data, and clean data. All respondents who are working on data-related jobs are grouped. Software engineers are those who build the software system and deploy the DL models.
- **Perception on code smells:** We asked respondents whether they have encountered the code smells. To clarify any possible confusion, we provided a definition and a simple example for each code smell. If they responded “yes”, we also asked them to what extent the code smell impacts their DL system maintenance (Very Serious, Serious, Moderate, Scarcely, and Not At All).

We followed a pilot protocol [14] while designing the survey. We designed a pilot version and sent them to a small subset of developers (11 developers). Based on the feedback, we rephrased some questions to make them easier to understand. We simplified and merged some questions to ensure that participants could finish the survey in 7 minutes. The responses from the pilot survey were used solely for improving the survey questions and were not included in the final results. We also translated our original survey to a Chinese version to support respondents who read Chinese before

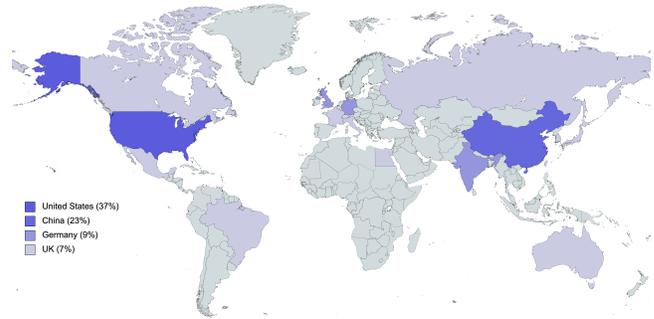


Figure 5: Countries in which survey respondents reside. The legend presents the top 4 countries with most respondents

distributing the survey. two of the authors (one of them is a native English speaker and the other a native Chinese speaker) discussed the survey and performed the translation together.

3.2.2 Respondent Selection. We aimed to get a sufficient number of practitioners from diverse backgrounds working on open source DL development and maintenance. Thus, we collected active contributors’ emails in the 59 DL projects by using GitHub REST APIs. In total, we collected 1,157 email addresses and successfully distributed them to 1,061 contributors. We kept the survey anonymous, but the respondents could choose to receive a summary of the study.

In total, we received 265 responses. After excluding incomplete surveys, 235 responses were considered valid. The countries and the corresponding number of respondents are shown in Fig. 5. The survey respondents who met our criteria are distributed across 15 countries and six continents. The majority of our respondents currently work in North America, Asia, and Europe, with the United States and China being the top two countries. Respondents’ software development experience varies from 1 to 23 years with an average of 5.25 years, and their DL development experience varies from 1 to 10 years with an average of 3.13 years.

3.2.3 Survey Data Analysis. To analyze the responses, we used descriptive statistics.

For the 235 valid responses to the question related to whether they have encountered our identified code smells, we normalized the frequency of each code smell by computing the percentage of respondents who have encountered code smells. If a high proportion of respondents reported that they have encountered a certain code smell, we consider this smell as more common. We did the same for the impact level of the code smell question. We also analyzed the responses based on roles. We mainly analyzed the responses from the top three categories of respondents which belonged to software engineers, ML engineers and data scientist/engineers since the number of project manager respondents were too few.

To check if there is a significant difference between new identified code smells in terms of impact, we adopted the Scott-Knott test [26]. Scott-Knott test divides the measurement averages into statically distinct groups by hierarchical clustering analysis. However, the limitations of the Scott-Knott test are that it assumes the data are in a normal distribution and it may create groups with trivially different from each other. Thus, we adopted its normality and effect

²<https://github.com/codesmell-material/codeSmell>

size-aware variant Scott-Knott effect size difference (ESD) test [45]. The Scott-Knott ESD test (1) corrects the normal distribution of the input data and (2) merges any two statistically different groups of negligible effects. A detailed description of the Scott-Knott ESD test can be found in [45].

4 RESULTS

In this section, we report the answers to our targeted research questions and findings that emerged from the data.

4.1 Maintenance Related Modifications in Deep Learning

RQ 1: What kinds of modifications do developers make frequently in DL systems?

To answer this question, we mined 59 open source DL project repositories and identified 426 maintenance-related modification commits. By using descriptive coding, we categorized the selected commits into nine modification categories (explained in Section 3). The modification categories, identified modification reasons, and their corresponding distributions are shown in Table 1.

Our manual analysis revealed that, as expected, some of the frequent modifications are specific to DL systems and others are not. For example, the most frequent (21%) modification category named **Change function declaration** which involves renaming functions, changing lambda functions to normal functions, and modifying function signatures is not specific to DL. **Extract class/function** category which includes changes pertaining to separating new class, isolating independent parts of code, and splitting long functions is also not specific to DL.

We also found that three of the modification categories are specific to DL.

Update/replace ML library: This recurring modification is the second most frequent category of modification (19%). Similar to API update/replace in traditional systems, developers usually use third-party DL libraries and frameworks to implement DL functionalities. However, DL libraries are usually updated more frequently than traditional libraries [42] and DL developers need to fix either deprecated or outdated functions to keep up with the updates. For example, the code snippet in Figure 3, shows that developers had to replace API names to resolve the compatibility issue with a newer version of TensorFlow.

Data preparation modification: This recurring modification is performed on the data preprocessing steps. We found that 8% of overall modifications in our dataset belonged to this group. Since a substantial part of code in DL systems is written for data preparation, and feeding to DL model any changes to the data source, preparation steps, or the model architecture requires this category of modification.

Model architecture modification: This recurring modification is done on DL model architecture related code. In order to resolve model degradation problems, developers iteratively train models or deploy new model architectures. We also found that developers make modifications to improve the model architecture by untangling the components. This group of modifications is 6% of our analyzed commits.

Observation 1: One third (33%) of the maintenance related modifications in DL systems are specific to DL systems and are related to the data, model, and library.

Interestingly, our results highlight another category of modification that is not specific to DL, but contains some DL specific changing reasons:

Replace hard-coded value: This is the recurring modification where developers replace hard-coded values with variables. Similar to traditional software, hard-coded values make it difficult to maintain software systems. We found developers frequently replace hard-coded model path, hyper-parameters, and learning rate with variables. 13% of our identified commits fall into this category.

Remove redundant debugging code: Developers frequently remove unnecessary debugging code in DL systems. The software engineering community has developed a number of tools, IDEs, and techniques to help catch bugs. Unfortunately, practitioners for DL systems do not enjoy the same robust set of debugging tools available for traditional software while debugging DL models due to the opaqueness of DL models and strong coupling between model and software components [42]. Thus, many DL developers resort to using print statements for debugging. 16% of maintenance-related modifications were grouped into this category.

Move code: In this category of recurring modification, developers move code between files and positions. Developers often put model training, testing, and validation related code in the same file. Later on, they end up moving each of the training, testing, and validation to separate files. We found that 6% of the modification commits belong to this category.

Remove dispensable dependency: This is the recurring modification where developers remove unused or unnecessary dependencies. Resolving dependency compatibility problems or versioning conflicts can be time consuming. As a result, developers are usually reluctant to remove dispensable dependencies until they have to. This kind of modification consists of 2% of modifications commits in our dataset.

4.2 Code smells in Deep Learning System

RQ 2: How prevalent are code smells in DL systems?

Through manual analysis of the maintenance related modifications done on real-world projects, we identified five code smells in DL systems (details in Section 3). Table 2 shows the five code smells along with their signs and symptoms ordered based on the frequency of occurrence in projects from high to low.

Scattered Use of ML Library: This smell is about implementing third-party ML libraries/frameworks in a non-cohesive manner throughout the project. As a result, whenever these libraries/frameworks update, developers have to modify multiple positions in single or multiple files. Such scattered use of ML library requires additional effort from the developer while maintaining the source code. 32 out of 59 (54%) projects have at least one commit showing this problem.

Unwanted Debugging Code: This smell was derived from the recurring pattern of leaving unwanted or unnecessary code in the DL system and we found 24 out of 59 DL projects have this code smell. DL systems tend to be more complicated than a traditional system and developers use debugging code for getting data shape

Table 1: Summary of modification commits and their distribution

Modification Categories	Modification reasons	Percentages of selected commits (%)
<i>Change function declaration</i>	Rename functions; Change lambda function to normal functions; Change function signatures; Convert public function to private; Convert private function to public.	21%
<i>Update/replace ML library</i>	Update deprecated functions with new ML library; Resolve python version compatibility issue; Resolve ML library compatibility issue; Switch to new ML library; Employ another ML library to improve model performance.	19%
<i>Remove redundant debugging code</i>	Clean up no longer used debugging of redundant code.	16%
<i>Replace hard-coded value</i>	Replace hard-coded model names, learning rate, parameters, numbers, and etc to variables.	13%
<i>Extract class/function</i>	Create a separate new class/function to remove old duplicate code; Isolate independent parts of code; Split long function/class.	9%
<i>Data preparation modification</i>	Resolve Data API compatibility issue; Separate data preparation code; Clean up data loader.	8%
<i>Move code</i>	Move code to proper files or positions; Simplify deep nested closure functions or containers.	6%
<i>Model architecture modification</i>	Rewrite model architecture source code; Deep learning layers and parameter modification; Replace with a new model; Separate model parts.	6%
<i>Remove dispensable dependency</i>	Remove unused or dispensable DL library/frameworks; Resolve dependent conflicts.	2%

Table 2: Summary of Newly Identified Code Smells in Deep Learning Systems

Code Smell	Signs and symptoms	Percentages of commits (%)	Percentages of Projects (%)
<i>Scattered Use of ML Library</i>	Scattered use of ML API in multiple files, once an ML API needs to be modified/updated, and the practitioners must modify places across several files.	13%	54%
<i>Unwanted Debugging Code</i>	A debugging code fragment, method, or class is no longer used, but still is left in the source code	17%	41%
<i>Deep God File</i>	A file contains multiple components of a DL system, such as model training, testing, etc	9%	37%
<i>Jumbled Model Architecture</i>	DL model architecture parts are cobbled together and difficult to understand and maintain.	6%	19%
<i>Dispensable Dependency</i>	An installed DL library or framework is no longer used or can be replaced by other existing ones	2%	9%

or printing current status to understand the code. However, left uncleaned these debugging codes can impede maintainability. If many people are working on a project, individuals are more reluctant to remove code that they do not thoroughly understand since no one wants to be responsible for errors. With these redundant codes left in the system, the code will be more difficult to understand, especially for DL systems.

Deep God File: This smell was derived from the recurring pattern where developers kept separating DL parts into multiple files after they had initially put some or all of them into one big file. We found 22 projects (37%) with this code smell. *Deep God File* usually starts small, but over time, they get bloated as practitioners may find it mentally easier to place programs into existing files.

Jumbled Model Architecture: This smell was derived from the recurring pattern when DL practitioners programmed the DL

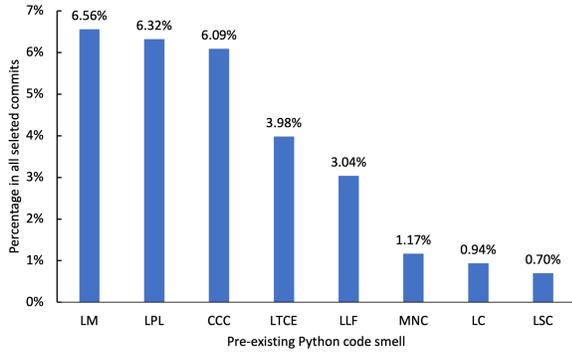


Figure 6: Prevalence of identified existing code smells

models, they do not clearly divide the different functional parts of the model. Instead, all parts of the model are jumbled together, which makes model code difficult to understand. We found 11 (19%) projects with this code smell.

Dispensable Dependency: This smell was derived from the recurring pattern where some redundant dependencies are left in DL systems and we noticed five out of 59 projects have modification commits to remove dispensable dependencies. Many DL libraries have repetitive functions, so some practitioners might try similar functions in each library and use the one with the best performance. However, this process adds some unnecessary dependencies to the entire system.

4.3 Prevalence of Python Smells

We used PySmell to analyze the Python code smells in our selected commits and identified eight Python code smells that were investigated by Hadhemi et al. [25] (shown in Figure 6). The most frequently fixed Python code smells in our dataset are LM, LPL, and CCC and their respective fixing commit percentages are 6.56%, 6.32%, and 6.09%. LTCE and LLF code smell fixing occupy 3.98% and 3.04% percentage of all selected commits. And the fixing commit percentages for MNC, LC, and LSC are 1.17%, 0.94%, and 0.70%.

According to Table 2, the percentage of selected commits to fix code smell of *Scattered Use of ML Library* and *Unwanted Debugging code* are respectively 13% and 17%. And for *Deep God File* and *Jumbled Model Architecture* code smells, the percentage of commits are 9% and 6%. The lowest percentage of commits is *Dispensable Dependency* code smell, which is 2%. By comparing the percentage of commits containing the code smells between our identified code smells and existing Python code smells we see that newly identified code smells are more frequent compared to generic Python code smells.

Observation 2: Newly identified code smells occur more frequently in our sample than generic Python code smells.

4.4 Code Smells Validation

RQ 3: How do practitioners perceive the identified code smells in DL systems?

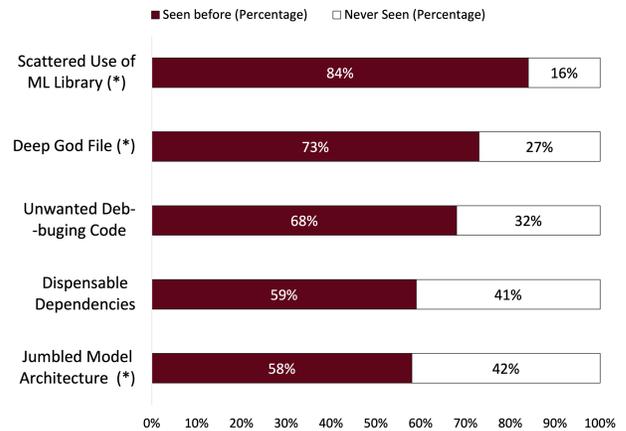
To answer this question, we analyzed the survey results. We asked respondents to what extent they have encountered these code smells and their perception about the impact of these code

smells on making DL systems difficult to maintain. The aggregated results are shown in Fig. 7.

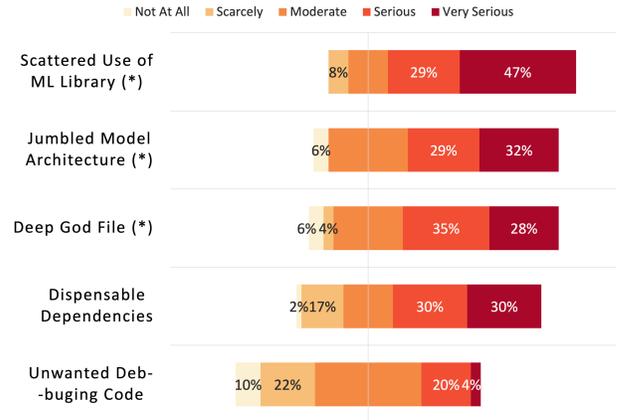
According to the aggregated results shown in Figure. 7-(a), respondents are familiar with the code smells we identified. 84% of the respondents expressed that they have seen *Scattered Use of ML Library* code smell before, which matches our repositories mining result that the most frequently occurred code smell is *Scattered Use of ML Library*. The respondents were also familiar with the other code smells. The ranking through mining was closely matched with the survey’s ranking as *Unwanted Debugging Code*, and *Deep God File* were among the top three code smells in both rankings.

According to the combined result from all the participants in Figure. 7-(b), the most impactful code smell is *Scattered Use of ML Library*. More than 60% of survey respondents reported that these code smells seriously impact their DL systems’ maintenance. According to the developers, the other two most impactful code smells are *Jumbled model architecture* and *Deep God File*. Among them, *Deep God File* is also the second most frequent code smell (Figure. 7).

Observation 3: According to both mining and developers’ responses, *Scattered Use of ML Library* and *Deep God File* are two of the most frequent and impactful code smells.



(a) Results of code smells frequency of occurrence



(b) Results of code smells impact on deep learning systems

Figure 7: Aggregated survey results form all respondents

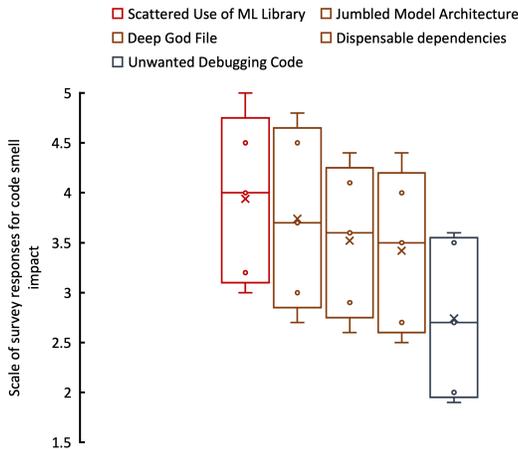


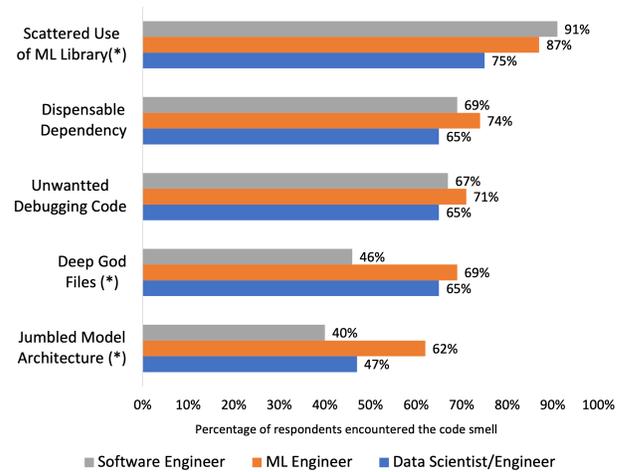
Figure 8: Scott-Knott effect size difference (ESD) test results

When we grouped the perceived frequency of code smells based on respondents' roles shown in Fig. 9-(a), the most common code smells for ML engineers, software engineer and data scientist/engineer respondents were *Scattered Use of ML Library*, *Dispensable Dependency*, *Unwanted Debugging Code*. However, software engineer met *Scattered Use of ML Library* more often, but ML engineer encountered *Dispensable Dependency* and *Unwanted Debugging Code* more often. It is also reasonable that software engineers and Data engineers encountered less *Jumbled Model Architecture* since they are not primarily maintaining models, but 62% ML engineer respondents encountered *Jumbled Model Architecture* code smell as they are primarily working with models.

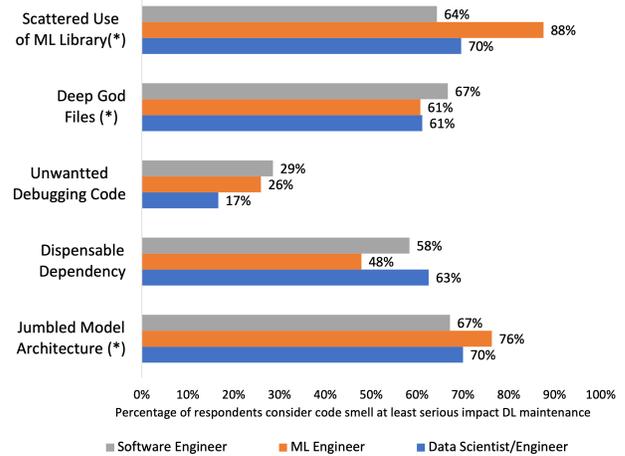
We looked into the impact of these code smells for each role, shown in Fig. 9-(b), which only shows the percentage of respondents who identified the code smell having serious or very serious impact on system maintenance. The *Scattered Use of ML Library* code smell is considered as the most severe by all three roles, especially by ML Engineers since 88% of ML Engineer respondents think this code smell has a "serious impact" on their system maintenance. Similarly, *Jumbled Model Architecture* is considered as a severe code smell by all three roles, even though it's not common for software and data scientist/engineer. In our analysis, we found that *Unwanted Debugging Code* is a common code smell, but most of respondents do not think it is a severe issue.

Observation 4: Different roles encounter code smells differently and they also have varied opinions about the impact of the code smell.

We conducted Scott-Knott ESD test on the responses collected for pertaining to the impact of code smells on DL maintenance to check if there is significant difference among all newly identified code smells. Figure 8 shows that Scott-Knott ESD categorized five code smells into three different groups. *Scattered Use of ML Library* is categorized in the first group as the most impactful code smell; *Jumbled Model Architecture*, *Deep God File*, and *Dispensable Dependency* code smells are categorized into the second group; and *Unwanted debugging Code* is categorized into the third group.



(a) Survey responses for different roles about code smell Occurrence (shows percentage of respondents have seen code smell before)



(b) Survey responses for different roles about code smell impact (only show the percentage of respondents think the code smell has serious or very serious impact)

Figure 9: Separated survey results for different roles

5 DISCUSSION AND IMPLICATIONS

In this section, we discuss the results presented in the previous section and present mitigation strategies, probable root causes for code smells, and practical implications of our study for researchers, educators, tool builders, and developers.

5.1 Mitigation Strategies

Scattered Use of ML Library: DL practitioners should import module as an alias to shorten the ML API message chain. That way, when an ML API is updated, maintainers no longer need to modify the usage of this API call throughout the whole project. Instead, they only need to change the code in module importing parts.

Jumbled model architecture: We suggest DL practitioners clearly separate the parts of the DL model with different functions, so that the model code is easier to understand and maintain.

Deep God File: We recommend developers place each part of code into a proper file, and have clear boundaries, such as placing model architecture, training, testing and validating program in separate files. If there is a Deep God File, DL practitioners can employ extract class or move function refactoring operations to separate components in such files.

Unwanted Debugging Code: We advocate DL practitioners remove unused debugging code in a timely manner.

Dispensable dependencies: We encourage practitioners remove unnecessary dependencies in DL systems since it takes a lot of time and effort to resolve dependency and library version conflicts in DL systems.

5.2 Probable root causes

ML teams are composed of different roles with overlapping tasks [12]. We posit that code smells might be a product of such overlapping tasks since the overlap in responsibility leads to unclear maintenance responsibility. A general thought is that this problem is not unique to DL systems but applies to regular systems as well. That is absolutely correct. Nevertheless, this problem is exacerbated in DL systems because of the significantly distinct roles of various team members. Along with creating confusion and dissatisfaction, uncertain responsibilities can result in dropped or mishandled source code and catastrophic consequences down the line. Practitioners need to ensure the maintenance task’s boundaries for different roles in DL systems.

Differences in job responsibilities among team members can be another reason for accumulating code smells over time. For example, ML engineers mainly concentrate on model development rather than software deployment and maintenance. To obtain a better model performance, they may try different ML libraries and add all tried library dependencies into the system at the same time. Even though ML engineers finally end up requiring only a few of the imported libraries, the unused but imported DL libraries and their dependencies remain in the system. Such unnecessary dependencies introduce additional problems to the software engineers who try to build and maintain the DL system. Since the process used at the ML developer’s end is opaque to the software engineers, it becomes difficult, even impossible in certain cases, for software engineers to remove any unused DL library dependencies. As a result, all the unused dependencies are left in the DL system and the quality of the system as a whole suffers [42]. Projects in the industry have started investigating ways to overcome these challenges. One approach is hybrid teams that include ML engineers, data scientists, and DevOps engineers [8]. Further work is needed to help DL systems identify and remove the unused dependencies. Improving cross-team communication, reducing the opaqueness in the development process used within the sub-groups along with ensuring documentation are some of the possible steps to mitigate this to some extent.

5.3 Implication

Implications for researchers, tool builders and educators: Our results show that DL systems have a wide variety of code smells. However, when we looked for code smell related work for DL,

we found limited studies. We encourage researchers to investigate many more kinds of code smells in DL systems.

Tool builders can focus on making the code smell detection tools seamlessly integrated into the existing DL development pipeline without causing major disruptions. This is important because research shows that if a workflow is disrupted, practitioners tend to stop using the tool [27].

The large variety of code smells in DL systems is also good news for educators. Educators can illustrate many design principles by showing both well-designed programs and those that exhibit code smells. Using DL systems as subject case studies is guaranteed to provide a variety of code smells. Moreover, students might also prefer examples from the DL domain given the rise and allure of DL programming.

Implications for DL developers: As Table 2 shows, identified code smells are distributed in a big percentage of DL systems. Thus, it is important that developers educate themselves about the kinds of code smells that occur in DL systems, and how to mitigate them. Or even better, being conscious about code smells when programming in the first place and avoid them altogether.

6 THREATS TO VALIDITY

Our refactoring pattern mining was performed on 59 projects carefully selected by Hadhemi et al. [25]. However, these are open source projects, which means the results may not be generalizable to all DL projects, particularly closed-source projects. Nonetheless, the majority of the DL projects use Python, so we believe our code mining on these Python projects still provides significant insights on code smells in DL systems. This is our first step towards building an empirical body of knowledge. With further replication across different contexts by different research teams, we can build a body of knowledge to generalize the results.

The manual analysis applied throughout the study could have introduced unintentional bias. First, we manually identified the commits that were related to maintenance activities based on commit messages and comparing the code before and after an update. Another manual analysis was conducted while grouping the frequently occurring change categories into code smells. This could have introduced bias or mistakes due to the lack of domain expertise. To address this concern, two researchers individually labeled a significant portion of the data. We established a high inter-rater agreement of 0.61 and 0.83 respectively for the two manual analyses, which according to Landis et al.[32], is considered as a substantial level of agreement and we believe we have minimized this threat.

We ran PythonChangeMiner to obtain frequently changed patterns, and then we used GitcProc to exclude bug fixes. Relying on these tools can be a threat to validity. However, these tools or variant of them has been validated in other studies. We also performed a manual investigation of any refactored code that has not been labeled by GitcProc to identify if there is any systemic error. Through our manual analysis, we did not see any evidence of the systemic error.

There is a possibility that our participants misunderstood the survey questions. To mitigate this threat, we conducted a pilot study with 11 developers with different background experiences and updated the survey based on the feedback. In order to clarify

any confusion, we provided definitions for each of the smells. Additionally, we translated the original survey to simplified Chinese to help native Chinese readers to reduce any confusion. Our survey’s language selection and translation process may be subject to bias. It might cause the group of respondents who can read Chinese and English to be over represented. However, it is important to mention that we chose to present our survey in English and Chinese because these are the top two most used languages in software development. Our survey could also have translation errors that cause the questions to deviate from the original meaning. To mitigate these risks, two of the authors (one of them is a native English speaker and the other a native Chinese speaker) discussed the survey and performed the translation together.

7 CONCLUSIONS AND FUTURE WORK

We investigated frequently occurring modifications in DL open source software repositories and identified nine modifications along with five code smells in this work. We also validated the code smells with DL practitioners through a survey. Participants identified the most impactful smells; however, surprisingly, the most frequent code smells are not necessarily the most impactful ones.

Our findings also open up new directions for future research. In addition to the future directions already presented in the discussion and implication sections, future research entails exploring the evolution of the identified code smells and their effect on DL systems’ overall quality.

REFERENCES

- [1] “2 killed in driverless tesla car crash, officials say,” <https://www.nytimes.com/2021/04/18/business/tesla-fatal-crash-texas.html>, accessed: 2020-07-1.
- [2] “20 natural language processing examples for businesses,” <https://www.wonderflow.ai/blog/20-natural-language-processing-examples-for-businesses>, accessed: 2020-07-1.
- [3] “The amazon machine learning process,” <https://docs.aws.amazon.com/machine-learning/latest/dg/the-machine-learning-process.html>, accessed: 2020-07-1.
- [4] “Mobileye is bringing its autonomous vehicle test fleets to at least four more cities in 2021,” <https://techcrunch.com/2021/01/11/mobileye-is-bringing-its-autonomous-vehicle-test-fleets-to-at-least-four-more-cities-in-2021/>, accessed: 2020-07-1.
- [5] “NiftyNet,” <https://github.com/NifTK/NiftyNet>, accessed: 2020-07-1.
- [6] “Pydriller,” <https://github.com/ishepard/pydriller>, accessed: 2020-07-1.
- [7] “Pythonchangeminer: a tool for mining change patterns in python projects,” <https://github.com/JetBrains-Research/code-change-miner>, accessed: 2020-07-1.
- [8] “Solving enterprise machine learning’s five main challenges,” <https://info.algorithmia.com/ml-challenges-ebook>, accessed: 2020-07-1.
- [9] “What is the team data science process?” <https://docs.microsoft.com/en-us/azure/machine-learning/team-data-science-process/overview>, accessed: 2020-07-1.
- [10] M. Abbes, F. Khomh, Y.-G. Gueheneuc, and G. Antoniol, “An empirical study of the impact of two antipatterns, blob and spaghetti code, on program comprehension,” in *2011 15th european conference on software maintenance and reengineering*. IEEE, 2011, pp. 181–190.
- [11] I. Ahmed, C. Brindescu, U. A. Mannan, C. Jensen, and A. Sarma, “An empirical examination of the relationship between code smells and merge conflicts,” in *Empirical Software Engineering and Measurement (ESEM), 2017 ACM/IEEE International Symposium on*. IEEE, 2017, pp. 58–67.
- [12] S. Amershi, A. Begel, C. Bird, R. DeLine, H. Gall, E. Kamar, N. Nagappan, B. Nushi, and T. Zimmermann, “Software engineering for machine learning: A case study,” in *2019 IEEE/ACM 41st International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP)*. IEEE, 2019, pp. 291–300.
- [13] R. Arcoverde, A. Garcia, and E. Figueiredo, “Understanding the longevity of code smells: preliminary results of an explanatory survey,” in *Proceedings of the 4th Workshop on Refactoring Tools*, 2011, pp. 33–36.
- [14] I. Brace, *Questionnaire design: How to plan, structure and write survey material for effective market research*. Kogan Page Publishers, 2018.
- [15] W. J. Brown, R. C. Malveau, H. W. McCormick III, and T. J. Mowbray, “Antipatterns: Refactoring software, architectures, and projects in crisis.(1998),” *Google Scholar Google Scholar Digital Library Digital Library*.
- [16] C. Casaluovo, Y. Suchak, B. Ray, and C. Rubio-González, “Giteproc: a tool for processing and classifying github commits,” in *Proceedings of the 26th ACM SIGSOFT International Symposium on Software Testing and Analysis*, 2017, pp. 396–399.
- [17] A. Chatzigeorgiou and A. Manakos, “Investigating the evolution of bad smells in object-oriented code,” in *2010 Seventh International Conference on the Quality of Information and Communications Technology*. IEEE, 2010, pp. 106–115.
- [18] Z. Chen, L. Chen, W. Ma, and B. Xu, “Detecting code smells in python programs,” in *2016 International Conference on Software Analysis, Testing and Evolution (SATE)*. IEEE, 2016, pp. 18–23.
- [19] Z. Chen, L. Chen, W. Ma, X. Zhou, Y. Zhou, and B. Xu, “Understanding metric-based detectable smells in python software: A comparative study,” *Information and Software Technology*, vol. 94, pp. 14–29, 2018.
- [20] M. D’Ambros, A. Bacchelli, and M. Lanza, “On the impact of design flaws on software defects,” in *2010 10th International Conference on Quality Software*. IEEE, 2010, pp. 23–31.
- [21] E. de Souza Nascimento, I. Ahmed, E. Oliveira, M. P. Palheta, I. Steinmacher, and T. Conte, “Understanding development process of machine learning systems: Challenges and solutions,” in *2019 ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM)*. IEEE, 2019, pp. 1–6.
- [22] I. Deligiannis, I. Stamelos, L. Angelis, M. Roumeliotis, and M. Shepperd, “A controlled experiment investigation of an object-oriented design heuristic for maintainability,” *Journal of Systems and Software*, vol. 72, no. 2, pp. 129–143, 2004.
- [23] M. Fowler, “Refactoring: Improving the design of existing code,” in *11th European Conference. Jyväskylä, Finland, 1997*.
- [24] K. Herzig, S. Just, and A. Zeller, “It’s not a bug, it’s a feature: how misclassification impacts bug prediction,” in *2013 35th International Conference on Software Engineering (ICSE)*. IEEE, 2013, pp. 392–401.
- [25] H. Jebnoun, H. Ben Braiek, M. M. Rahman, and F. Khomh, “The scent of deep learning code: An empirical study,” in *Proceedings of the 17th International Conference on Mining Software Repositories*, 2020, pp. 420–430.
- [26] E. Jelihovschi, J. C. Faria, and I. B. Allaman, “The scottknott clustering algorithm,” *Universidade Estadual de Santa Cruz-UDESC, Ilheus, Bahia, Brasil*, 2014.
- [27] B. Johnson, Y. Song, E. Murphy-Hill, and R. Bowdidge, “Why don’t software developers use static analysis tools to find bugs?” in *2013 35th International Conference on Software Engineering (ICSE)*. IEEE, 2013, pp. 672–681.
- [28] F. Khomh, M. Di Penta, and Y.-G. Gueheneuc, “An exploratory study of the impact of code smells on software change-proneness,” in *2009 16th Working Conference on Reverse Engineering*. IEEE, 2009, pp. 75–84.
- [29] F. Khomh, M. Di Penta, Y.-G. Guéhéneuc, and G. Antoniol, “An exploratory study of the impact of antipatterns on class change-and fault-proneness,” *Empirical Software Engineering*, vol. 17, no. 3, pp. 243–275, 2012.
- [30] D. P. Kingma and M. Welling, “Auto-encoding variational bayes,” *arXiv preprint arXiv:1312.6114*, 2013.
- [31] T. Lam, N. Hauser, A. Götz, P. Hathaway, F. Franceschini, H. Rayner, and L. Zhang, “Gumtree—an integrated scientific experiment environment,” *Physica B: Condensed Matter*, vol. 385, pp. 1330–1332, 2006.
- [32] J. R. Landis, *A general methodology for the measurement of observer agreement when the data are categorical*, 1975.
- [33] M. Lanza and R. Marinescu, *Object-oriented metrics in practice: using software metrics to characterize, evaluate, and improve the design of object-oriented systems*. Springer Science & Business Media, 2007.
- [34] W. Li and R. Shatnawi, “An empirical study of the bad smells and class error probability in the post-release object-oriented system evolution,” *Journal of systems and software*, vol. 80, no. 7, pp. 1120–1128, 2007.
- [35] A. Lozano, M. Wermelinger, and B. Nuseibeh, “Assessing the impact of bad smells using historical information,” in *Ninth international workshop on Principles of software evolution: in conjunction with the 6th ESEC/FSE joint meeting*, 2007, pp. 31–34.
- [36] N. Moha, Y.-G. Guéhéneuc, L. Duchien, and A.-F. Le Meur, “Decor: A method for the specification and detection of code and design smells,” *IEEE Transactions on Software Engineering*, vol. 36, no. 1, pp. 20–36, 2009.
- [37] M. J. Munro, “Product metrics for automatic identification of “bad smell” design problems in java source-code,” in *11th IEEE International Software Metrics Symposium (METRICS’05)*. IEEE, 2005, pp. 15–15.
- [38] H. A. Nguyen, T. N. Nguyen, D. Dig, S. Nguyen, H. Tran, and M. Hilton, “Graph-based mining of in-the-wild, fine-grained, semantic code change patterns,” in *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*. IEEE, 2019, pp. 819–830.
- [39] D. Rapu, S. Ducasse, T. Girba, and R. Marinescu, “Using history information to improve design flaws detection,” in *Eighth European Conference on Software Maintenance and Reengineering, 2004. CSMR 2004. Proceedings*. IEEE, 2004, pp. 223–232.
- [40] D. Sahin, M. Kessentini, S. Bechikh, and K. Deb, “Code-smell detection as a bilevel problem,” *ACM Transactions on Software Engineering and Methodology (TOSEM)*, vol. 24, no. 1, pp. 1–44, 2014.
- [41] J. Saldaña, *The coding manual for qualitative researchers*. Sage, 2015.

- [42] D. Sculley, G. Holt, D. Golovin, E. Davydov, T. Phillips, D. Ebner, V. Chaudhary, M. Young, J.-F. Crespo, and D. Dennison, "Hidden technical debt in machine learning systems," in *Advances in neural information processing systems*, 2015, pp. 2503–2511.
- [43] D. I. Sjøberg, A. Yamashita, B. C. Anda, A. Mockus, and T. Dybå, "Quantifying the effect of code smells on maintenance effort," *IEEE Transactions on Software Engineering*, vol. 39, no. 8, pp. 1144–1156, 2012.
- [44] I. Sommerville and R. Thomson, "An approach to the support of software evolution," *The Computer Journal*, vol. 32, no. 5, pp. 386–398, 1989.
- [45] C. Tantithamthavorn, S. McIntosh, A. E. Hassan, and K. Matsumoto, "An empirical comparison of model validation techniques for defect prediction models," *IEEE Transactions on Software Engineering*, vol. 43, no. 1, pp. 1–18, 2016.
- [46] M. Tufano, F. Palomba, G. Bavota, R. Oliveto, M. Di Penta, A. De Lucia, and D. Poshyvanyk, "When and why your code starts to smell bad," in *2015 IEEE/ACM 37th IEEE International Conference on Software Engineering*, vol. 1. IEEE, 2015, pp. 403–414.
- [47] Z. Wan, X. Xia, D. Lo, and G. C. Murphy, "How does machine learning change software development practices?" *IEEE Transactions on Software Engineering*, 2019.
- [48] A. Yamashita and M. Leon, "Exploring the impact of inter-smell relations on software maintainability: An empirical study," in *2013 35th International Conference on Software Engineering (ICSE)*. IEEE, 2013, pp. 682–691.
- [49] A. Yamashita and L. Moonen, "Do code smells reflect important maintainability aspects?" in *2012 28th IEEE international conference on software maintenance (ICSM)*. IEEE, 2012, pp. 306–315.