



Delft University of Technology

Fine-grained just-in-time defect prediction

Pascarella, Luca; Palomba, Fabio; Bacchelli, Alberto

DOI

[10.1016/j.jss.2018.12.001](https://doi.org/10.1016/j.jss.2018.12.001)

Publication date

2019

Document Version

Accepted author manuscript

Published in

Journal of Systems and Software

Citation (APA)

Pascarella, L., Palomba, F., & Bacchelli, A. (2019). Fine-grained just-in-time defect prediction. *Journal of Systems and Software*, 150, 22-36. <https://doi.org/10.1016/j.jss.2018.12.001>

Important note

To cite this publication, please use the final published version (if applicable). Please check the document version above.

Copyright

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

Takedown policy

Please contact us and provide details if you believe this document breaches copyrights. We will remove access to the work immediately and investigate your claim.

Fine-Grained Just-In-Time Defect Prediction

Luca Pascarella,¹ Fabio Palomba,² Alberto Bacchelli²

¹*Delft University of Technology, The Netherlands* — ²*University of Zurich, Switzerland*
l.pascarella@tudelft.nl, palomba@ifi.uzh.ch, bacchelli@ifi.uzh.ch

Abstract

Defect prediction models focus on identifying defect-prone code elements, for example to allow practitioners to allocate testing resources on specific subsystems and to provide assistance during code reviews. While the research community has been highly active in proposing metrics and methods to predict defects on *long-term* periods (*i.e.*, at release time), a recent trend is represented by the so-called *short-term* defect prediction (*i.e.*, at commit-level). Indeed, this strategy represents an effective alternative in terms of effort required to inspect files likely affected by defects. Nevertheless, the granularity considered by such models might be still too coarse. Indeed, existing commit-level models highlight an *entire* commit as defective even in cases where only *specific* files actually contain defects.

In this paper, we first investigate to what extent commits are partially defective; then, we propose a novel *fine-grained just-in-time* defect prediction model to predict the specific files, contained in a commit, that are defective. Finally, we evaluate our model in terms of (i) performance and (ii) the extent to which it decreases the effort required to diagnose a defect. Our study highlights that: (1) defective commits are frequently composed of a mixture of defective and non-defective files, (2) our fine-grained model can accurately predict defective files with an AUC-ROC up to 82% and (3) our model would allow practitioners to save inspection efforts with respect to standard just-in-time techniques.

Preprint <https://doi.org/10.5281/zenodo.1880517>. Appendix <https://doi.org/10.5281/zenodo.1886063>.

Keywords: Just-in-Time Defect Prediction, Empirical Software Engineering, Mining Software Repositories

1. Introduction

During software maintenance and evolution, developers constantly modify the source code to introduce new features or fix defects [40]. These modifications, however, may lead to the introduction of new defects [37], thus developers must carefully verify that the performed modifications do not introduce new defects in the code. This task is usually performed directly during development (*e.g.*, by running test cases) [77] or when changes are reviewed [3]. An efficient way to allocate inspection and testing resources to the portion of source code more likely to be defective is represented by *defect prediction* [24], which involves the construction of statistical models to predict the defect-proneness of software artifacts, by mostly exploiting information regarding the source code or the development process.

The problem of defect prediction has attracted the attention of many researchers in the past decade, who tried to address it by (i) conducting empirical studies on the factors making artifacts more defect-prone (*e.g.*, [6, 36, 60, 61, 63, 68, 79, 81]) and (ii) proposing novel prediction models aimed at accurately predicting the defect-proneness of the source code (*e.g.*, [8, 25, 51, 57, 62, 70]).

Most of the existing techniques evaluate the defectiveness of software artifacts perform *long-term* predictions. Analyzing the information accumulated in previous soft-

ware releases, these models predict which artifacts are going to be more prone to defect *in future releases*. For instance, Basili *et al.* investigated the effectiveness of Object-Oriented metrics [12] in predicting post-release defects [6], while other approaches consider process metrics (*e.g.*, the entropy of changes [25]) or developer-related factors [8, 57] for the same purpose.

Kamei *et al.* reported that these long-term defect prediction models—despite their good accuracy—may have a limited usefulness in practice because they do not provide developers with immediate feedback [34], thus not avoid the introduction of defects during the commit of artifacts on the repository. To overcome this limitation, a recent trend is the investigation of *just-in-time* prediction models, *i.e.*, techniques exploiting the characteristics of a commit to perform *short-term* predictions of the likelihood of a commit introducing a defect. With this solution, a developer can limit the effort required to diagnose problems since s/he focuses on the committed artifacts only [34]. Among the studies investigating *just-in-time* prediction models, Kamei *et al.* [34, 32] defined 14 metrics characterizing a commit under five perspectives, demonstrating how such metrics can be successfully exploited for predicting defective commits either in the case the model is trained using previous data of the project [34] and in the case the training information come from different projects [32]. Other approaches proposed the use of deep-learning [88],

textual analysis [5], and unsupervised methodologies [89].

It is reasonable to think that, in a real-world scenario, a commit may be *partially* defective, *i.e.*, it may be composed of both defective and non-defective files. In this case, despite the advantages provided by *just-in-time* defect prediction, a developer might still need to spend a considerable effort to locate the files of a commit that are actually defective. For instance, during a Modern Code Review (MCR) the reviewers iterate several times over the proposed set of changes and the amount of time spent finding a subset of defective files might substantially increase [3]. In this paper, we aim at making a further step ahead in the context of *just-in-time* defect prediction by investigating the original problem at a finer granularity. Particularly, our goal is to investigate the prominence of partially defective commits and, should they be a significant amount, devise a defect prediction model to identify the *defective files within a commit*.

To this aim, we firstly performed an exploratory study to characterize defective commits and evaluate whether fine-grained solutions are actually needed. In the second place, we built a *fine-grained just-in-time* defect prediction model adapting 24 basic features previously defined in the papers by Kamei *et al.* [34] and Rahman and Devanbu [70]. Finally, we assessed the performance of the model in terms of (i) accuracy of the predictions and (ii) effort developers can save using our model with respect to state-of-the-art *just-in-time* prediction models. The study was conducted considering 10 major open source systems and 160,515 commits. Key findings of our investigations revealed that (i) almost 43% of defective commits are composed of a mixture of both defective and non-defective resources, (ii) the devised *fine-grained* model obtained an AUC-ROC up to 82% when locating defective files in a commit, and (iii) our model is more cost-effective than the state of the art *just-in-time* model.

Structure of the paper. Section 2 reports background, related work, and a concept of the envisioned solution. Section 3 reports the methodology used to address our research goal as well as possible threats that might influence our findings, while Section 4 presents the results of the study. Finally, Section 5 concludes the paper.

2. Background and Related Work

In this section we introduce the terminology used through the paper, discuss the related work, and motivate our study.

2.1. Terminology

Throughout the paper, we frequently refer to the following five concepts:

Defect. To define a condition in which a software system does not meet its requirements, we use the term *defect*, among all the possible terms (*e.g.*, *bug* and *fault* [24]).

Defect-Inducing/-Fixing Change. We identify two events in the life of a defect: (i) the *defect-inducing change*, *i.e.*, the code change that inserts the defect into a project and (ii) the *defect-fixing change*, *i.e.*, the code change that fixes the defect.

Commit. In most modern collaborative software projects, authors develop code relying on version system control tools such as Git.¹ Such tools track changes as *commits*, which are documented changes that involve one or more files.

Non-/Partially/Fully Defective Commit. We define three classes of commits: *non-defective commits* (when all the committed files are changed without introducing any defect), *fully defective commits* (when all the committed files are changed introducing defects), and *partially defective commits* (when a subset of the committed files are changed introducing a defect). The top part of Figure 1 depicts a part of the history of an example software system, with the activities made on the versioning system after a system’s release. A set of commits $C = \{c_x, \dots, c_{x+4}\}$ are performed by developers to evolve the system; all the commits change the same three files (A.c, B.c, and C.c). In Figure 1, we see examples of non-defective commits (c_x and c_{x+3}), partially defective commits (c_{x+1} , c_{x+4}), and a fully defective commit (c_{x+2}).

2.2. Related Work

In this section, we discuss the related work that inspired and guided this study, considering long- and short-term defect prediction.

2.2.1. Long-term Defect Prediction

Long-term defect prediction pertains to models able to classify defect-prone files in future releases of a software project. Several studies addressed this problem in the recent years (a relatively recent survey has been compiled by Hall *et al.* [24]). Basically, the proposed models differ for the source of information used to predict the defectiveness of a class: the main distinction is between static, product information and historical, process data.

Product Information. Structural data are computed with metrics such as the McCabe’s cyclomatic complexity [47] or the Chidamber and Kemerer (CK) [12] metrics. These product metrics have been investigated in several studies [2, 16, 17, 22, 26, 41, 56, 23] and researchers have shown how such metrics can provide useful contribution in the prediction of defective classes. For instance, Nagappan *et al.* [56] found that a model based on code metrics may achieve up to 83% of accuracy in the identification of defect-prone classes.

¹<https://git-scm.com/>

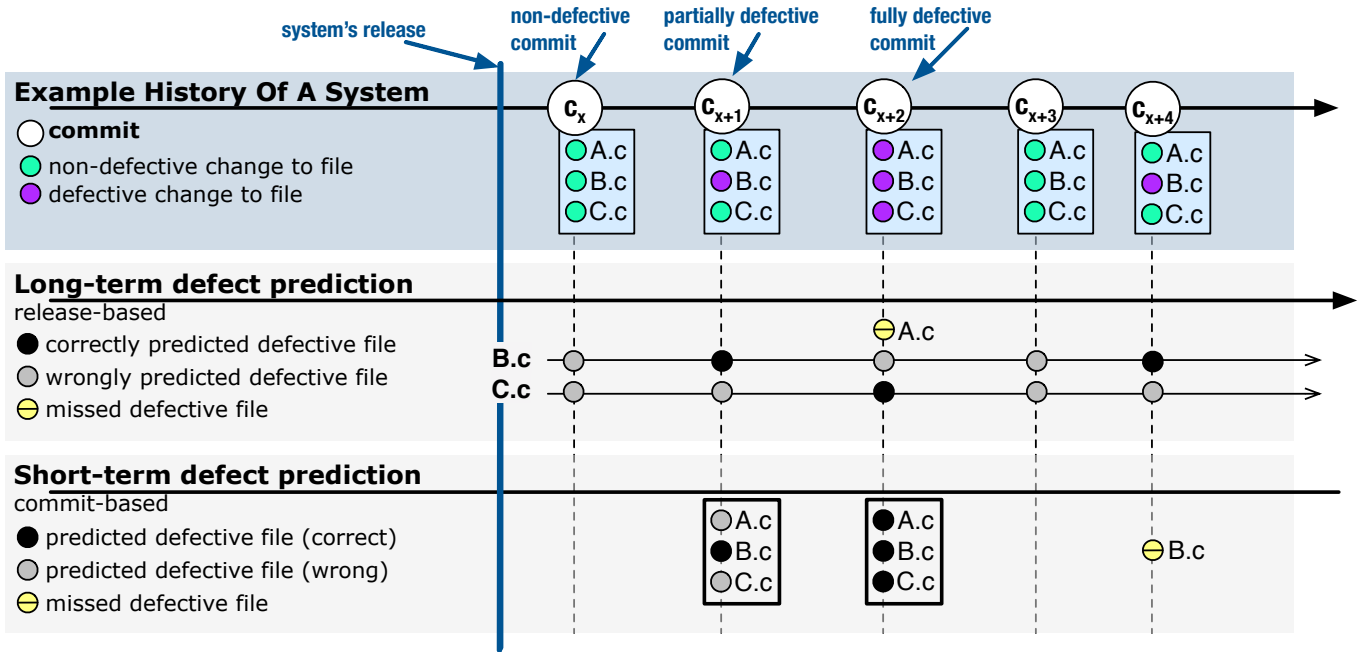


Figure 1: An example set of (defective) commits after release to files in a software system.

Process Information. Historical data are computed with metrics such as relative code churn, entropy of changes, or developer-related factors [25, 22, 24, 54, 55, 57]. Also in this case, researchers provided empirical evidence on the value of such metrics for defect prediction. For instance, Moser *et al.* [54] performed a comparative study analyzing static- and historical-based predictors, concluding that metrics computed over the history of projects are better predictors and can significantly improve the performance of defect prediction models.

Combining Information. Later on, D’Ambros *et al.* [14] found that combined techniques work better than models based on single unique set of metrics. On the basis of this result, di Nucci *et al.* [57] defined a combined model based on a mixture of static and historical metrics able to outperform the prediction capabilities of single models. Finally, Menzies *et al.* [51] introduced the concept of *local* defect prediction, an approach in which classes that will be used for training the classifier are firstly clustered into homogeneous groups to reduce the differences among such classes and obtain higher prediction accuracy.

We build on top of this line of work, by considering the features that are better able to predict defects at file-level—a key attribute that we include in our model.

2.2.2. Short-term Defect Prediction

Short-term defect prediction refers to models able to classify defect-prone at commit time. Previous work motivated the introduction of this new strategy with the need of having tools able to locate defects in the shortest possible time [52]. While Madeyski *et al.* [44] proposed the

idea of continuous defect prediction, Mockus *et al.* [52] addressed the problem by proposing a model based on the change-proneness of files to predict defects at commit-level.

Characteristics of Defect-Introducing Changes.

Other studies (*e.g.*, Sliwerski *et al.* [75] and Eyolfson *et al.* [18]) tried to localize defect-introducing changes in open source projects by means of correlation between the defectiveness of a commit and the experience of developers. Sliwerski *et al.* [75] also discovered that defect-introducing changes are generally a part of large transactions. The “unnaturalness” of defective code was subsequently confirmed by Ray and Hellendoorn *et al.* [72], who discovered that source code presenting defects is characterized by higher entropy than non-defective code. They also found that source code entropy might be a valid and simpler way to complement the effectiveness of static analysis tools (*e.g.*, CHECKSTYLE²) in recommending to developers the areas of source code where to focus inspection activities.

Ad Hoc Models. Jiang *et al.* [29] proposed the concept of “personalized” defect prediction by proposing a technique able to create a different model for each developer. Their results report that such a technique outperforms existing *just-in-time* defect prediction models. Along with this line, Xia *et al.* [86] improved the aforementioned technique by Jiang *et al.* using a multi-objective genetic algorithm that firstly builds a defect prediction model for each developer, and then combine these models assigning different weights with the aim of maximizing

²<http://checkstyle.sourceforge.net>

F-Measure and cost-effectiveness. With respect to these papers, our approach has not the goal to build a prediction model for each developer, but instead that of providing feedback on defect-prone classes within the scope of a commit: further analysis will evaluate the possible benefits provided by the creation of personalized models in the context of fine-grained just-in-time defect prediction.

Just-In-Time. The studies by Kamei *et al.* [34, 32] are great source of inspiration for our work. They proposed a *just-in-time* quality assurance technique that predict defects at commit-level trying to reduce the effort of a reviewer [34]. Later on, they also evaluated how *just-in-time* models perform in the context of cross-project defect prediction [32]. Their main findings report good accuracy for the models in terms of both precision and recall, but also in terms of saved inspection effort. Our work is complementary to these papers. In particular, we start from their basis of detecting defective commits and complement this model with the attributes necessary to filter only those files that are defect-prone and should be more thoroughly reviewed. Rahman *et al.* [71], Yang *et al.* [88], and Barnett *et al.* [5] proposed the usage of alternative techniques for *just-in-time* quality assurance, such as cached history, deep learning, and textual analysis, reporting promising results. We did not investigate these further in the current paper, but studies can be designed and carried out to determine if and how these techniques can be used within the model we present in this paper to further increase its accuracy.

2.3. Motivating Example

We discuss an example in which a developer uses long- vs. short-term defect prediction models while inspecting a commit, in order to show some of the limitations of these approaches.

The top part of Figure 1 depicts an example history of a software system, with the activities made on the versioning system after a system’s release. A set of commits $C = \{c_x, \dots, c_{x+4}\}$ are performed by developers to evolve the system. For sake of clarity, suppose that the files A.c, B.c, and C.c are always changed in the considered commits after the system’s release. The small circles in the top bar represent changes made to files in each commit and the colors represent whether these changes introduce a defect (purple) or not (dark green) in these files. In addition, a black box surrounds all the files in the same commit. In the following we describe the behavior of the two aforementioned prediction models:

Long-term Defect Prediction. Based on the information gathered before the system’s release, a long-term defect prediction model would mark certain files as defect-prone for the entire period leading to the issue of the next release. In our example, the model marks the files B.c and C.c as defect-prone and A.c as non-defective. The model classifies both B.c and C.c as defective starting

from the system’s release onward, in Figure 1 we depict this behavior with a horizontal small arrows, thus showing that B.c and C.c are considered as defective in *every commit*. Indeed, the model does not provide any information about the exact commit that will likely lead to the introduction of a defect. This model would issue warnings about these files on each commit involving them. In our example, this represents an unjustified extra-effort for the developer inspecting the commit. As found in previous research, this unjustified extra-effort derived from using a tool can reduce the developers’ confidence in the prediction [65], thus leading to miss important defects in future commits involving actual defect-prone artifacts. Finally, we see that the model does not classify as defective the code in file A.c, also when a defect is introduced in commit c_{x+2} (the missed defective file is depicted with a yellow circle).

Short-term Defect Prediction. As an alternative, a reviewer may adopt a short-term defect prediction model such as the *just-in-time* one proposed by [34]. In this scenario, a developer is pointed to analyze more in depth only the files referring to a commit marked as potentially defective by the model. However, the number of resources to inspect might be still high depending on the number of files committed and the wasted effort on how many are defect-free. For instance, in the commit c_{x+1} shown in Figure 1, only the file B.c introduces a defect, while the others are defect-free, yet a warning from the tool would be issued; the developer may need to analyze some non-defective files before finding the actual defect. Thus, while short-term solutions can significantly reduce the reviewers’ effort, they might still produce extra-effort in cases a commit is partially defective. Furthermore, in our example, file B.c is again defective in commit c_{x+4} , but it is not marked as such, since the model does not recognize the commit as defective (in the figure, the missed defective file is depicted with a yellow circle).

The goal of our work is to make the first steps in supporting software developers during the inspection of a commit (*e.g.*, in a code review), by striving to overcoming the aforementioned limitations of existing defect prediction models in this context. The next section details our research questions and the research method.

3. Methodology

This section defines the overall goal of our study, motivates our research questions, and outlines our method.

3.1. Research Questions

The *goal* of the study is to investigate how frequently commits are only partially defective and to devise a defect prediction model able to identify the files with the changes that are more likely to introduce a defect. We set up our work around three research questions. The first

Table 1: Characteristics of the subject software systems.

Project	KLOC	Developers	Commits	Defective Commits
Accumulo	102	66	8,747	1,399
Angular-js	87	1,589	8,467	1,525
Bugzilla	239	99	9,788	3,621
Gerrit	79	38	22,232	4,001
Gimp	102	216	38,240	8,412
Hadoop	291	92	15,556	2,606
JDeodorant	70	9	1,105	199
Jetty	88	70	12,638	3,286
JRuby	129	322	38,894	8,945
OpenJPA	822	25	4,848	1502
Overall	2,009	2,526	160,515	35,496

one is a preliminary analysis aimed at assessing the extent to which a defect prediction model is actually able to estimate the defect-proneness of files within a commit. To this aim, we investigate the ratio of commits that contain both defective and not defective files. Should the frequency of partially commits be low, standard just-in-time models, such as the one devised by Kamei et al. [34] would be sufficient, while in case there should be a notable percentage of commits presenting both defective and non-defective files, then defect prediction models working at a finer granularity than standard just-in-time ones would be desirable.

RQ₁. *What is the ratio of partially defective commits?*

Once assessed the actual need for finer grained solutions, we devise a defect prediction model to predict defective files at commit scope.

RQ₂. *To what extent can we predict defect-inducing changes at file-level in a commit?*

In addition to assessing the model as a whole, we also evaluate which features provide the highest contribution to the achieved performance.

RQ₃. *What are the features of the devised model that the most to its performance?*

Finally, we are also interested in understanding how much effort could be saved when using the proposed model, comparing it to the just-in-time defect prediction model proposed by Kamei et al. [34] as our baseline.

RQ₄. *How much effort can be saved using a fine-grained just-in-time defect prediction model with respect to a standard just-in-time model?*

In the following sections, we describe the steps we perform to answer our three research questions.

3.2. Subject Systems

To conduct our analysis, we focused on open-source software systems and defined multiple selection criteria: We selected software systems (i) written in the most common programming languages (C, C++, Java, JavaScript, Ruby, and Perl, i.e., the most popular programming languages [10]), (ii) having different size and scope, and (iii) having a change history composed of at least 1,000 commits. We preferred open-source systems where a versioning system is used to track all changes. The access to the source code history enables the computation of metrics with static analysis tools. Moreover, to increase the generalizability of our research, we selected software projects having different domains and programming languages. Note that our selection is not intended to be statistically significant, but rather we just aim at selecting a various set of systems to assess the performance of our prediction model in different contexts (e.g., when considering projects having different change history sizes). In practice, we started from the entire list of open source projects available on GITHUB; then we filtered out systems not implemented in the considered programming languages and with less than 1,000 commits in their history. Successively, among 2,362,287 project candidates, we considered only the most popular projects for a given domain or scope; finally, we randomly selected the ten open-source software systems reported in Table 1. For each system, the table reports size (in terms of KLOCs), number of developers, number of commits, and the information on the number of defective commits.

3.3. RQ₁ - Investigating Defective Commits

To answer our first research question, we analyze the ratio of the defective files (i.e., source code, configuration, and auxiliary files) contained in defective commits. To this aim, for each commit c_i of the change history of a system S , we identify the set $defectiveFiles(c_i)$ composed of the defective resources contained in c_i . To the best of our knowledge, there is not a publicly available dataset reporting this information: Previous work defined datasets of defective commits [34], without providing details on which of the resources in a certain commit were actually defective. For this reason, we build our own dataset as detailed in the following.

Data Extraction. To automatically identify the set of defective files in each of the commits of the considered systems, we rely on the SZZ algorithm [76, 84]. SZZ exploits the annotation/blame feature of a versioning system to *estimate* the lines of code of a file that induced a certain defect, thus retrieving files that are defect-inducing in each commit. More formally, the algorithm implements the following steps:

1. For each file f_i (where $i = 1..n$) involved in a defect fixing commit dfc , the algorithm

$prevVersion(commit, file)$ extracts the last version of the file before the defect fixing commit: $prevVersion(df_c, f_i)$;

- Starting from the commit $prevVersion(df_c, f_i)$, for each line of code in f_i changed to fix the defect in df_c , the algorithm uses `git blame` to detect the file revision where the last change to that line occurred. We identify comments and empty lines using island parsing [53] and we exclude f_i if no other code is touched. This step outputs the commits in which a defect in file f_i is introduced.

The SZZ algorithm takes as input the list of defects that are *already fixed* by developers, excluding the open ones,³ but the analysis and the effect of considering open issues will be considered in future work (*e.g.*, exploiting tools such as RELINK [85]).

Data Analysis. Once extracted the defective files involved in defective commits, we answer **RQ₁** in two ways. First, we measure how many defective commits are *partially* defective, *i.e.*, they contain a mixture of both defective and non-defective resources. This analysis allow us to understand the magnitude of the problem investigated: If the vast majority of defective commits is composed of only defective artifacts, then standard *just-in-time* defect prediction models would suffice; conversely, if a significant part of defective commits is *partially* defective, then the introduction of fine-grained solutions might be worthwhile. Second, we further analyzed the set of *partially* defective commits, by measuring the ratio between defective and non-defective files they contain. More formally, we computed the $defectiveFiles_{dc}$ ratio as follow:

$$defectiveFiles_{dc} = \frac{\#defectiveFiles(dc)}{\#files(dc)} \quad (1)$$

where $\#defectiveFiles(dc)$ represents the number of defective files in the defective commit dc , and $\#files(dc)$ the total number of files in dc . This analysis helped us to understand the intrinsic characteristics of *partially* defective commits. Also in this case, if the resulting ratio is high (most files are defective in *partially* defective commits), then the adoption of fine-grained solutions would be not worthwhile.

3.4. **RQ₂** - The Fine-Grained JIT Model

To answer our second research question, we build a *fine-grained just-in-time* defect prediction model and evaluate its performance. In the following, we describe (i) the independent variables, *i.e.*, the metrics on which the model relies, (ii) the dependent variable, *i.e.*, the characteristic that the model have to predict, (iii) the machine learner performing the predictions, and (iv) the validation methodologies to estimate the accuracy.

³Open issues might be not verified by developers (*i.e.*, they might be not real defects).

i. Independent Variables. This step consists in extracting and quantifying the characteristics of each file involved in a commit. To this purpose, we considered the 24 basic features shown in Table 2. These features represent a modified version of those previously proposed by Kamei et al. [34] and Rahman and Devanbu [70]. We adapted the previous metrics to work at file-level in a commit. The column ‘Description’ in Table 2 details the implementation of the metrics in our context. The choice of the independent variable is driven by two goals: (i) to understand the value of standard *just-in-time* measures in a fine-grained context; (ii) to investigate whether metrics originally proposed in the context of long-term defect prediction to predict defective files may also provide useful contributions when employed in the prediction of defective files contained in a change set.

Furthermore, the chosen metrics help us to characterize commits under different perspectives, thus allowing us to evaluate which metric types are more relevant in our context. Specifically, we selected metrics to measure (i) the developers’ experience (*e.g.*, the experience of the committer [34]), (ii) structural and process factors of the files in the commit (*e.g.*, the lines of code added or the number of previous changes of a committed file [70]), and (iii) factors related to the neighbors’ of a committed file, which have been shown to be relevant for predicting the defectiveness of files [70]. Although other metrics have been proposed in the contexts of both code review (*e.g.*, by McIntosh et al. [48] and Kononenko et al. [38]) and defect prediction (*e.g.*, [14, 57]), the selected metrics better allow us to verify the role of a larger set of metrics that have been previously adopted for traditional short- and long-term defect prediction. Further studies can be conducted to investigate the addition of other metrics in our context.

From a methodological standpoint, the process metrics adapted from Rahman and Devanbu [70] (*i.e.*, COMM, ADEV, DDEV, ADD, DEL, OWN, MINOR, SCTR, NADEV, NDDEV, NCOMM, NSCTR, OXEP, and EXP) were always evaluated considering the commits up to the commit of interest. Similar adjustments were applied for the metrics proposed by Kamei et al. [34]. For instance, the NUC metric represents the number of unique changes to the files modified in a commit. In our case, we adjust NUC to represent the number of times a single file involved in a commit is modified alone up to the considered commit. Descriptions of how we adapted the Kamei et al. [34] and Rahman and Devanbu [70] metrics are reported in Table 2.

ii. Dependent Variable. The characteristic to measure is the defectiveness of files contained in a commit. To this aim, we exploited the dataset built in the context of **RQ₁** (*i.e.*, we used the output of the SZZ algorithm as a dependent variable to predict).

iii. Machine Learner. In this stage, we needed to se-

Table 2: List of the independent/predicting variables adapted from Rahman *et al.* [70]* and Kamei *et al.* [34]**

Acronym	Name	Description	Ref.
COMM	Commit Count	Number of changes to the file up to the considered commit	*
ADEV	Active Dev Count	Number developers who modified to the file up to the considered commit	*
DDEV	Distinct Dev Count	Cumulative number of distinct developers contributed to the file up to the considered commit	*
ADD	Normalized Lines Added	Normalized number of lines added to the file in the considered commit	*
DEL	Normalized Lines Deleted	Normalized number of lines removed to the file in the considered commit	*
OWN	Owner's Contributed Lines	Boolean value indicating whether the commit is done by the owner of the file	*
MINOR	Owner's Contributed Lines	Number of contributors who contributed less than 5% of the file up to the considered commits	*
SCTR	Changed Code Scattering	Number of packages modified by the committer in the commit	*
NADEV	Neighbor's Active Dev Count	Number of developers who changed the files involved in commits where the file has been modified	*
NDDEV	Neighbor's Distinct Dev Count	Cumulative number of distinct developers who changed the files involved in commits where the file has been modified	*
NCOMM	Neighbor's Commit Count	Number of commits made to files involved in commits where the file has been modified	*
NSCTR	Neighbor's Commit Count	Number of different packages touched by the developer in commits where the file has been modified	*
OEXP	Neighbor's Commit Count	Percentage of lines authored in the project	*
EXP	All Committer's Experience	Mean of the experiences of all the developers	*
ND	Number of modified directories	Number of modified directories	**
Entropy	Distribution of modified code across each file	Entropy of changes of the file up to the considered commit	**
LA	Lines of code added	Number of lines added to the file in the considered commit (absolute number of the ADD metric)	**
LD	Lines of code deleted	Number of lines removed to the file in the considered commit (absolute number of the DEL metric)	**
LT	Lines of code in a file before the change	Lines of code in the file before the change	**
AGE	Average interval between the last and the current change	The average time interval between the last and the current change	**
NUC	Number of unique changes to the modified files	Number of times the file has been modified alone up to considered commit	**
CEXP	Experience of the committer	Number of commits made on the file by the committer up to the considered commit	**
REXP	Recent developer experience (last x months)	Number of commits made on the file by the committer in the last month	**
SEXP	Developer experience on a subsystem	Number of commits made by the developer in the package containing the file	**

lect a machine learning classifier able to use the independent variables to infer the defectiveness of files in a change set [20]. To this aim, we tested different classifiers (using the validation methodologies described later in this section), *i.e.*, *Binary Logistic Regression* [39], *J-48* [45], *ADTree* [19], *Multilayer Perceptron* [83], *Naive Bayes* [31], and *Random Forest* [43]. As a result, we found that the *Random Forest* technique [43] is the one having the highest performance, in line with previous findings [30, 73]. A complete report of such analysis is available in our online appendix [66].

Such classifiers builds several decision trees, each of them containing nodes representing a condition on a certain feature that splits the dataset into two. A condition is chosen based on the so-called *Mean Decrease in Impurity* (MDI) [21], a metric able to measure the extent to which the value of a feature can correctly discriminate the dependent variable. It is important to point out that the selected classifier automatically performs a feature selection, thus avoiding the well-known problem of multi-collinearity [58] that occurs when two or more independent variables correlate with each other, possibly affecting the performance of the classifier.

iv. Validation Methodologies. The final step to answer **RQ₂** is related to the validation of the model. Commonly used techniques such as *ten-fold cross* [15, 80], or *leave-one-out cross-validation* [74] are not suitable for the validation of just-in-time defect prediction models because the data points (*i.e.*, the commits) follow a certain time order: *Time-insensitive* validation strategies might cause a model to be trained using future data that should not be known at the time of the prediction [80]. For this reason, we adopt a *time-sensitive* analysis where the defectiveness of a commit c_i is evaluated by a model trained using the data coming from the previous three months of

history of the system considered. In other words, while the training set is composed of three-month data, the test set is represented by each commit singularly. Doing so, we exclude the first three months of change history, because of the lack of data needed to perform a proper validation [80]. Our choice of considering three-month periods is based on: (i) choices made in previous work [25, 57, 80] and (ii) the results of an empirical assessment we performed on such a parameter. The empirical assessment showed that the best performance for the devised model is achieved by using three-month periods. In particular, we experimented with time windows of one, two, three, and six months. The complete results are available in our replication package [66].

Afterward, we measure the performance of the model using *precision* and *recall* [4]:

$$precision = \frac{|TP|}{|TP + FP|} \quad (2)$$

$$recall = \frac{|TP|}{|TP + FN|} \quad (3)$$

where TP , FP , and FN are:

- TRUE POSITIVES (TP): elements that are correctly retrieved by the *fine-grained just-in-time* prediction model (*i.e.*, defective files correctly classified as such);
- FALSE POSITIVES (FP): elements that are wrongly classified by the *fine-grained just-in-time* prediction model (*i.e.*, non-defective files misclassified as defective by the model);
- FALSE NEGATIVES (FN): elements that are not retrieved by the *fine-grained just-in-time* prediction model (*i.e.*, defective files misclassified as non-defective by the model).

In addition, to have a unique value that synthesizes precision and recall we also measure the *F-measure*, *i.e.*, the harmonic mean of precision and recall:

$$F - Measure = 2 \cdot \frac{Precision \cdot Recall}{Precision + Recall} \quad (4)$$

While the metrics described so far have been widely used in the past to evaluate defect prediction models [24], most of the classifiers output a probability ranging between 0 and 1 representing the likelihood of a code component to be part of a certain class (*i.e.*, in our case, to be *defective* or *non-defective*). The threshold used to discriminate the two classes (in most cases—as well as in this work—such threshold is set to 0.5) influences the computation of both precision and recall, and as a consequence of F-Measure. ROC plots the true positive rates against the false positive rates for all possible thresholds between 0 and 1; the diagonal represents the expected performance of a random classifier. AUC computes the area below the ROC and allows us to have a comprehensive measure for comparing different ROCs: An area of 1 represents a perfect classifier (all the defective methods are recognized without any error), whereas for a random classifier an area close 0.5 is expected (since the ROC for a random classifier tends to the diagonal). To have a detailed view of the performance of the model in the different cases found in **RQ₁**, in Section 4.2 we report the evaluation metrics achieved when ran the model over the set of (i) all the defective commits in the dataset, (ii) partially defective commits only, and (iii) fully defective commits only.

3.5. **RQ₃** - Investigating the Importance of the Features

While in **RQ₂** we provide an overview of the accuracy of the devised model in predicting defective files within a commit, **RQ₃** has the goal of investigating which features contribute the most to the prediction capabilities. To address this point, we use an *information gain* algorithm [69] to quantify the gain provided by each independent variable to the prediction of defective files within commits. Formally, let M be the devised *fine-grained just-in-time* prediction model, let $F = \{f_1, \dots, f_n\}$ be the set of features used by M , the *information gain* algorithm [69] applies the following formula to compute the difference in entropy:

$$InfoGain(M|f_i) = H(M) - H(M|f_i) \quad (5)$$

where the function $H(M)$ indicates the entropy of the model that includes the feature f_i , while the function $H(M|f_i)$ measures the entropy of the model that does not include f_i . Entropy is computed as follow:

$$H(M) = - \sum_{i=1}^n prob(f_i) \log_2 prob(f_i) \quad (6)$$

The algorithm quantifies the degree of uncertainty in M that is reduced by considering the feature f_i . In our

work, we employ the *Gain Ratio Feature Evaluation* algorithm [69], which ranks f_1, \dots, f_n in descending order based on the contribution provided by each feature to the decisions made by M . More specifically, the output of the algorithm is represented by a ranked list in which the features having the highest expected reduction in entropy are placed at the top. During this step, we also verified—through the `evaluateAttribute` function of the WEKA⁴ implementation of the algorithm—whether a certain feature mainly contributes to the identification of defective or non-defective files, *i.e.*, if there exists a positive or negative relationship between the feature and the defect-proneness of files.

3.6. **RQ₄** - Measuring the Saved Effort

For **RQ₄** we investigate the potential benefits in terms of saved effort that the *fine-grained just-in-time* defect prediction model provides to a developer analyzing the committed files to discover possible defects (*e.g.*, in a code review). Specifically, we perform an effort-aware validation as recommended by Ostrand et al. [59]. In this formulation, a technique is assessed on the fraction of defects it can detect while varying the effort required to locate them. As done in previous work [34], we first rank the files to inspect according to their probability of being defective, as it is assigned by the automated classifier (in our case, *Random Forest*); then we measure the percentage of defects that a developer would identify as the effort spent in analyzing the suggested defective files increases. To approximate such an effort, we use the *number of lines of code to inspect*; this metric has been shown to be a surrogate measure of the effort needed for testing or reviewing a module, as code and cognitive complexity are strongly related to size [56]. Thus, size can be considered as a lightweight and efficient solution to estimate the developer’s effort in inspecting a code change [2, 33, 49, 50].

We compare our model to the traditional *just-in-time* defect prediction model proposed by Kamei et al. [34]. The selection of this baseline is driven by experimental tests, where we found that this approach works better than the twelve unsupervised techniques proposed by Yang et al. [89]. In particular, the model by Kamei et al. [34] achieves an AUC-ROC 6% higher than the best unsupervised technique, which was the one that predicts a commit as defective in case of a number of committed files higher than eight. We report the results of this additional analysis in our online appendix [66].

To perform a fair comparison, the baseline relies on the same predictors used by Kamei *et al.* in their experiments and is trained using the best performing classifier (*i.e.*, *Random Forest*, the same used by our approach). We also empirically evaluate the performance of the several classifiers, namely *Binary Logistic Regression* [39], *J-48* [45], *ADTree* [19], *Multilayer Perceptron* [83], *Naive Bayes* [31],

⁴<https://www.cs.waikato.ac.nz/ml/weka/>

and *Random Forest* [43], when applied on the model by Kamei et al. [34]. Also in this case, *Random Forest* classifier outperforms the others.

As the baseline can only assign defect probabilities to commits (because of the commit-level granularity), we assume that the same probability holds for all files within that commit. In other words, if a commit is considered defective by Kamei *et al.*'s technique, then all the files within that commit are considered as potentially defective and have the same probability to require further inspection by a developer. To determine in which sequence the developer would inspect the files in the same commit, we use the *alphabetical* order, because it is the normal order offered by both IDEs and code review tools [7]. Once we assign the probabilities/order to all the files, we rank them in descending order and compare it with the ranking provided by our technique. It is worth noting that we expect our technique to outperform this baseline, as by definition it aims at lowering the granularity of the information presented to developers. Nevertheless, we still consider this comparison useful because we can verify whether and how much our approach actually meet the expected goal.

Finally, we perform a comparison with the *optimal* approach that ranks all the actual defective files first, starting from the smallest to the largest. In this way, we can investigate how far our technique is with respect to optimal scenario as well as how much it improves upon existing *just-in-time* approaches.

Data Analysis. To quantify the differences between our model and the baselines, we use the P_{opt} and P_k evaluation metrics [49]. P_{opt} is defined as the Δ_{opt} between the effort-based cumulative lift charts of the optimal model and the devised prediction model. Similarly, Δ_k is defined as the Δ_k between our technique and the one by Kamei et al. [34]. Larger values of P_{opt} and P_k indicate smaller differences between the compared techniques. Such values are normalized in the range [0,1] to ease their interpretation [34].

3.7. Threats to validity

The results of our study may be affected by a number of threats.

Threats to construct validity. As for factors threatening the relation between theory and observation, in our context, these are mainly concerned with the measurements we performed. Above all, we rely on the results of the SZZ algorithm [76] to answer our research questions. Although the intrinsic imprecisions of SZZ [13] still represent a threat for the validity of our results, it is the most effective algorithm available in literature.

To compute the CEXP, REXP, SEXP metrics, we mined commits to count the number of modifications applied by a developer in different time windows. However, it might be possible that the actual author of a commit is not the same person as the committer. That may be especially

true in large projects where sometimes developers (*e.g.*, newcomers) can modify the source code but do not have rights to perform a push onto the repository. This potential problem might have influenced the way the metrics are computed and used within the devised prediction model. To verify the extent to which this represents an actual issue for our analyses, we quantified in how many cases there was a mismatch between author and committer in the analyzed commits. Specifically, for each commit of the considered projects, we ran the command `git show --format=full`⁵ to obtain the full set of information available for the commit. That includes data on both author and committer email addresses. Thus, we could compute the number of times in which the two email addresses differ, *i.e.*, in how many cases the author of a change was not the actual committer. Out of the 160,515 total commits considered in our study, we found 4,173 mismatches, meaning that we are not accurate in only 2.6% of the cases. Based on this result, we can argue that such mismatches represent corner cases rather than systematic problems that threatens our analysis. To further verify the impact of this potential threat, we completely re-ran our study excluding those 4,173 commits. However, we did not observe any difference for the results achieved when including the commits. That indicates that mismatches between authors and committers do not influence our findings. A complete overview of this additional analysis is available in our online appendix [66].

Threats to conclusion validity. Although the metrics used to evaluate the performance of the *fine-grained just-in-time* defect prediction model, (*i.e.*, precision, recall, F-measure, and AUC-ROC), are widely used in the field [14], future studies can be conducted to validate our model from a different angle, *e.g.*, by evaluating its industrial impact.

A possible threat concerning the results achieved in **RQ₁** is related to the co-presence of production and test files within a commit, which may lead to a over-estimation of the number of partially defective commits. We conducted an additional analysis to assess the effect of excluding test files on our findings. We could not find differences with respect to the results reported in the original submission (a complete report of this additional analysis is available in our online appendix [66]). These results are in line with recent work: Even test code may be defective [82] and test files have the same proneness of production files to be affected by functional issues [78]. It seems reasonable to keep test files in our analysis/approach to maintain developers' awareness also on bugs in tests.

Another threat regards how we assess the cost-effectiveness of the models experimented. As done in previous research [9, 64, 2, 33, 49, 50], we measure the inspection cost in terms of lines of code to be inspected

⁵<https://git-scm.com/docs/git-show>

by a reviewer. LOC has been evaluated as a valid proxy measure [2] since it is correlated with code and cognitive complexity [56]. However, this is an approximation. Function points (FPs) [46] represent an alternative that we do not consider in this study, since it requires setting parameters that only original developers/managers or expert effort estimation consultants might properly set and a third-party analysis done by the authors of this paper would introduce noise/bias. Future work can be designed and conducted to investigate how much size approximates defect inspection effort.

In the context of **RQ₂**, we adopt a time-sensitive validation strategy where a single commit c_i represents the test set and the data of the previous three months form the training set. We select this strategy because this is the most similar to a real-case scenario where developers use the devised approach as soon as a new commit is performed, for example in a code review. While other researchers adopted slight variations of this strategy (*e.g.*, Tan et al. [80] used a gap between training and test sets to add in the training set defective commits that were discovered and fixed), we preferred it for its stronger ecological validity.

We statistically compare the differences between our model and the standard *just-in-time* model proposed by Kamei et al. [34]. We do not perform statistical tests with the Bonferroni correction [1]: This is a conscious decision taken on the basis of the findings by Perneger [67], who explained why such a correction is unnecessary and deleterious for sound statistical inference. Finally, we assess the model for the presence of multi-collinearity [58], relying on *Random Forest*, which can automatically remove non-relevant features.

As a final note, we compare our model with the one proposed by Kamei et al. [34] in the context of **RQ₄** (the cost-effectiveness analysis) but not in **RQ₂** (the accuracy analysis). On the one hand, the model by Kamei et al. [34] targets a different problem (*i.e.*, detecting defective commits rather than defective files within commits), thus it cannot be fairly compared with the proposed model in terms of accuracy. This statement is supported by experimental data, which showed that the model by Kamei et al. [34] achieved an overall F-Measure of 31% and AUC-ROC of 53% when employed in our context (by considering all the files within an identified defective commit as defective). On the other hand, the comparison performed regarding cost-effectiveness allows us to understand and quantify the gain provided by our approach against state of the art.

Threats to external validity. The main issue concerns the generalizability of the results. To alleviate this issue, we take into account a variety of projects having different characteristics, scope, and size. Nevertheless, future studies must be devised to replicate and extend our investigation on a larger set of systems, possibly taking into

consideration industrial projects as well.

Table 3: Results for RQ1 on partially defective commits

Systems	Ratio		
	Partially defective commits	Defective files	Avg. files per commit
Accumulo	46%	44%	4.1
Angular-js	51%	38%	2.2
Bugzilla	47%	37%	5.4
Gerrit	38%	43%	3.3
Gimp	44%	45%	4.3
Hadoop	49%	38%	3.1
JDeodorant	39%	47%	3.4
Jetty	53%	46%	3.8
JRuby	45%	42%	3.5
OpenJPA	40%	37%	4.0
Overall	43%	42%	3.7

4. Results and Analysis

In this section, we present the results of the study by research question.

4.1. **RQ₁**. What is the ratio of partially defective commits?

The analysis of the results associated to the first research question aims to understand the prominence of *partially* defective commits, hence the importance of devising a *fine-grained* solution for just-in-time defect prediction. Table 3 reports the results for each considered system: The second column reports the percentage of *partially* defective commits contained in the considered systems, the third column shows the percentage of defective files for each projects (computed using Formula 1), and the fourth column reports the average number of files per commit in the considered systems. The last row (“Overall”) represents the average ratio computed taking into account all the projects as a single dataset.

Among all the defective commits investigated we found that 43% of them are *partially* defective, *i.e.*, they contain a mixture of both defective and non-defective files, while 57% of defective commits only contain one resource. Thus, while standard *just-in-time* models can be adopted in most cases, there still exists a consistent part of defective commits for which they cannot provide developers with detailed information.

Investigating the *partially* defective commits more in depth, we found that on overall only 42% of committed files are defective; this is quite surprising, since it implies that less than the half of the elements in a *partially* defective commit is actually defective. Considering the perspective of a developer who has to inspect the files in a change set, she might spend more than half of the time inspecting non-defective resources before finding an actual defect.

For instance, let us consider the commit `a0641ea475` belonging to the ANGULAR.JS project.⁶ In this case, the developer committed 9 different files with the aim of making configurable the errors to show in case of wrong usage of the tool. However, there was only one defective file in the whole change set, *i.e.*, the `minErr.js` one. As a consequence, the usage of *coarse-grained just-in-time prediction* model such as the one proposed by Kamei *et al.* might not provide the adequate support in these cases. The observations made until now still hold when considering the “best” scenario reported in the table, *i.e.*, the one of the JDEODORANT project, where we found that 47% of the resources in a defective commit is affected by a problem, enforcing a developer to inspect many non-defective resources before diagnosing the defect.

With the aim of further understanding the characteristics of defective commits, we also computed the Kendall’s τ correlation [35] between the number of files per commit and the number of defective files. This is a non-parametric statistical test used to measure the ordinal association between two measured quantities, with a value ranging between -1 and +1.⁷ In our case, the correlation between number of files per commit and number of defective files turned to be equals to 0.42, thus indicating a positive concordance between the two variables. This confirms previous findings reporting that the more resources a developer changes the higher the chances to introduce defects [28].

In conclusion, the results show the need of fine-grained techniques to reduce the number of resources to inspect in a defective commit.

Result 1: 42% of defective commits in our subjects are *partially* defective, *i.e.*, composed of both files that are changed without introducing defects and files that are changed introducing defects. Further, in almost 43% of the changed files a defect is introduced, while the remaining files are defect-free.

4.2. RQ₂. To what extent can the model predict defect-inducing changes at file-level?

To answer our second research question we evaluate the effectiveness of the prediction model described in Section 3.4 based on a machine learning algorithm built using the *Random Forest* classifier. For sake of clarity, we report the results of both RQ₂ and RQ₃ in three separated tables that have a similar structure. The columns “RQ₂” report the evaluation metrics, *i.e.*, *precision*, *recall*, *F-measure*, and *AUC-ROC*, for each system. Table 4 is obtained evaluating our model considering indiscriminately all commits in the history of the projects, instead Table 5 considers

⁶<https://github.com/angular/angular.js/pull/15881>

⁷(i) -1 represents a perfect negative linear relationship, (ii) +1 a perfect positive linear relationship, and (iii) the values in between indicate the degree of linear dependence between the two measured quantities

Table 4: Results of the RQ₂ considering all commits in the history of the subject software systems.

Systems	RQ2				RQ3
	Precision	Recall	F-measure	AUC-ROC	Cost-effectiveness
Accumulo	71%	66%	69%	82%	16% (L)
Angular-js	73%	62%	68%	79%	18% (L)
Bugzilla	65%	65%	65%	73%	4% (M)
Gerrit	69%	62%	65%	72%	8% (L)
Gimp	61%	59%	60%	69%	16% (L)
Hadoop	67%	58%	63%	73%	7% (L)
JDeodorant	75%	61%	68%	74%	11% (L)
Jetty	60%	65%	62%	77%	17% (L)
JRuby	64%	61%	62%	70%	21% (L)
OpenJPA	63%	60%	61%	72%	7% (M)
Overall	67%	62%	65%	76%	13% (L)

Table 5: Results of the RQ₂ considering only partially defective commits in the history of the subject software systems.

Systems	RQ2				RQ3
	Precision	Recall	F-measure	AUC-ROC	Cost-effectiveness
Accumulo	76%	69%	73%	85%	19% (L)
Angular-js	75%	63%	69%	77%	23% (L)
Bugzilla	69%	68%	68%	74%	4% (M)
Gerrit	77%	66%	72%	77%	9% (L)
Gimp	65%	63%	64%	69%	17% (L)
Hadoop	68%	61%	64%	76%	11% (L)
JDeodorant	83%	68%	76%	79%	16% (L)
Jetty	64%	69%	67%	83%	20% (L)
JRuby	66%	63%	65%	73%	28% (L)
OpenJPA	66%	64%	65%	71%	7% (M)
Overall	72%	65%	69%	77%	16% (L)

Table 6: Results of the RQ₂ considering only fully defective commits in the history of the subject software systems.

Systems	RQ2				RQ3
	Precision	Recall	F-measure	AUC-ROC	Cost-effectiveness
Accumulo	65%	63%	64%	72%	7% (M)
Angular-js	71%	61%	66%	80%	13% (L)
Bugzilla	61%	62%	61%	71%	3% (S)
Gerrit	62%	58%	60%	66%	6% (M)
Gimp	57%	58%	57%	68%	13% (L)
Hadoop	66%	56%	60%	64%	6% (M)
JDeodorant	74%	59%	67%	73%	9% (L)
Jetty	61%	64%	62%	71%	11% (L)
JRuby	62%	59%	61%	69%	18% (L)
OpenJPA	60%	60%	60%	67%	6% (M)
Overall	63%	61%	62%	70%	10% (L)

only *partially* defective commits and Table 6 represents only *fully-defective* commits.

Looking at the full-inclusive results of Table 4, we observe that the precision ranges between 60% and 75% (overall=67%), the recall between 58% and 66% (overall=62%), while the overall F-measure is equal to 65%. Interesting are the results in terms of precision: in a context where the recommendations are given when developers are committing their changes on the repository, having a tool able to pinpoint the files that are likely defective can avoid the introduction of a consistent number of defects in a system. Assuming that developers can recognize a defect should they get a true positive warning from our model, the adoption of our model has the potential to be useful in practice, since its precision is higher than 60% in most of the cases. The recall values tell us that our model locates more than half of the defects actually present in the subject systems.

Considering also the AUC-ROC we observe that the model obtains levels between 69% and 82% (overall=76%). The worst case observed in our dataset regards the GIMP project, where our model achieves the lowest F-measure (60%). Investigating the likely causes behind this result, we found that our model was not able to work on the other projects because of the peculiar characteristics of the *C* programming language used. In particular, 39% of the change sets were composed of interacting files (*e.g.*, a file *C* including functions from other files): consistent modifications to files including several external files were often performed, while minor defective changes were performed on the other files. As a consequence, the metrics computed (*e.g.*, the normalized number of lines of code added) were not effective.

We do not observe large decays between the overall metric values and the highest/lowest ones (*i.e.*, the difference is always within 21%). This means that the *fine-grained just-in-time* model is consistent across the projects.

When analyzing the results obtained by the models built only considering *partially* and *fully* defective commits (Table 5 and Table 6), we observe that the former outperforms the latter by 7% in terms of both F-measure and AUC-ROC. Since the goal of this paper is to assess the extent to which a prediction model can identify defective files within commits, we consider the performance of the technique built on *partially* defective commits as encouraging because the proposed approach is actually able to meet the intended goal. At the same time, we consider the performance degradation noticed on *fully* defective commits as reasonable. Unfortunately, we are not able to speculate on the specific reasons causing such degradation. Likely, the addition of further independent variables able to characterize the defectiveness of commits as a whole (*e.g.*, the metrics devised by Kamei *et al.* [34]) can be beneficial to improve the performance of the model further. A future research effort can be devoted to the potential combination between just-in-time and fine-grained just-in-time

Table 7: Gain Provided by Each Feature To The Prediction Model.

Variable name	Expected Entropy Reduction	Shape
CEXP	0.76	defective
LA	0.71	defective
NCOMM	0.68	non-defective
REXP	0.64	defective
ND	0.58	non-defective
SCTR	0.55	defective
Entropy	0.49	non-defective
OWN	0.48	non-defective
SEXP	0.43	defective
LD	0.33	non-defective
MINOR	0.31	defective
DEL	0.28	non-defective
ADD	0.21	defective
COMM	0.19	defective
NSCTR	0.17	defective
DDEV	0.12	defective
NDDEV	0.10	non-defective
OEXP	0.09	non-defective
EXP	0.09	defective
ADEV	0.06	non-defective
NADEV	0.04	non-defective
NUC	0.04	non-defective
LT	0.03	defective
AGE	0.02	defective

models. In any case, our results show that in the majority of the cases the model can provide further recommendations also when considering fully-defective commits. Finally, the model including all commits inherits pros and cons observed in the cases of the models built on *partially* and *fully* defective commits only. In other words, it can predict *partially* defective commits better than *fully* defective ones, having higher performance on the former and lower on the latter; that shifts performance in the middle to the individual models.

Result 2: The proposed model achieves an overall *AUC-ROC* of 76% and obtains stable performance across the considered projects.

4.3. *RQ*₃. What are the features of the devised model that the most to its performance?

Table 7 reports the results achieved when applying the *Gain Ratio Feature Evaluation* algorithm [69] to understand which are the most relevant features that allow the model to identify defect-inducing changes within the files of a commit. In particular, for each variable we report (i) the expected entropy reduction it gives to the model and (ii) the shape of the relationship with the dependent variable, *i.e.*, whether the feature contributes more to the prediction of defective or non-defective files.

Four key factors give the highest contribution to the performance of the model, *i.e.*, experience of the committer, lines of code added, neighbor’s commit count, and

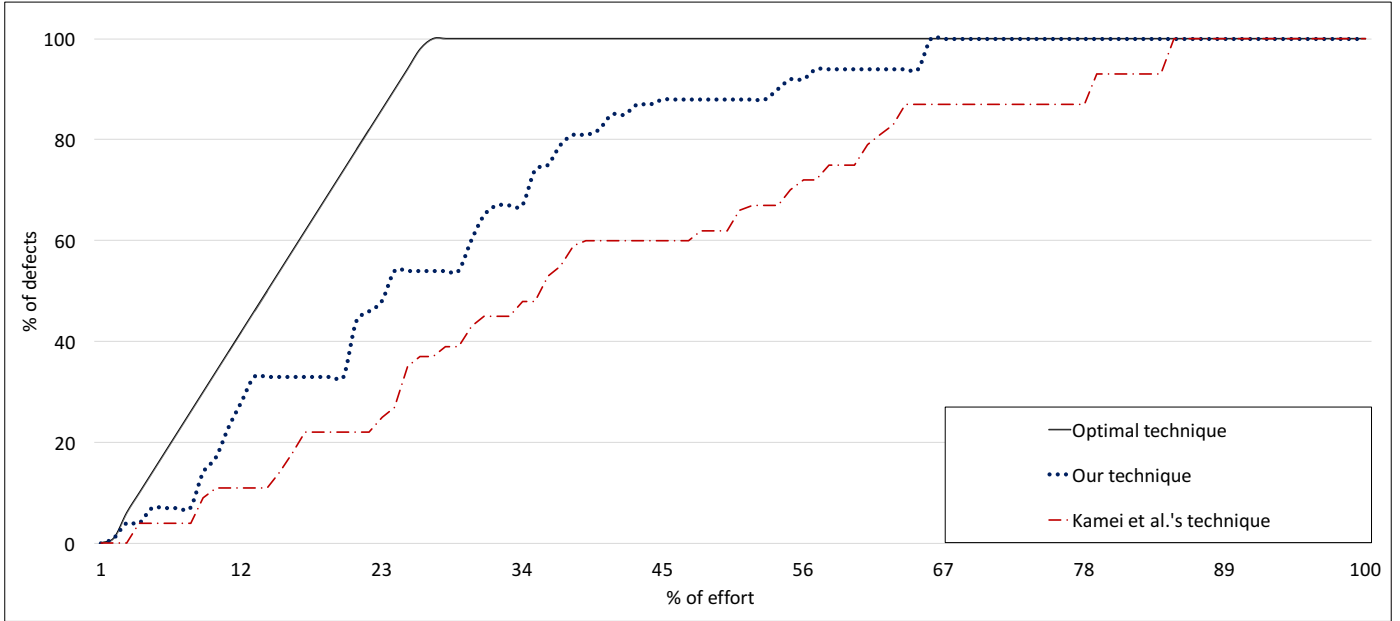


Figure 2: Results achieved for **RQ4**.

recent experience of the committed. All of them have an expected entropy reduction higher than 0.6, which means that they are highly relevant to discriminate defective files. The number of modified directories is also pretty relevant, with an entropy reduction of 0.58, while the result of the changed code scattering (entropy reduction=0.55) confirms that non-focused modifications tend to have an adverse effect on source code quality [57]. Other factors, even though lower in ranking, can characterize defective files within a commit, such as entropy or owner’s contributed lines. The remaining independent variables tend to be less/poorly related to defective files (*e.g.*, the average interval between the last and the current change).

Looking more in depth into the results, we observe that the model relies on different types of information to discriminate defective and non-defective files. For example, experience-related metrics have a positive relationship with the dependent variable, meaning that they mostly help in detecting defective files. It seems reasonable to think that this is due to a substantial difference in the behavior of developers expert and non-expert of a particular piece of code, which the model can correctly interpret for the identification of defect-inducing changes. Similarly, the amount of lines added mainly indicates the defectiveness of an artifact. On the other hand, several factors contribute more to the prediction of non-defective files. Among them, the NCOMM feature, which represents the neighbor’s commit count, has a strong impact on the predictions made by the model on non-defective files: this indicates that the number of changes applied to files connected to a specific file can particularly characterize the lack of defects.

To sum up, this analysis let emerge that defective commits are strongly explained by developer-related factors

(thus corroborating the need for methodologies and tools for an efficient allocation of resources) and that developers should perform small changes. Our results are in partial agreement with the findings by Kamei et al. [34]: indeed, only the experience of the committer is a powerful predictor in both traditional and fine-grained just-in-time defect prediction. Instead, when comparing our results with those of Rahman and Devanbu [70], we confirm that process metrics are generally better predictors than product ones. More in general, we observe that no single family of metrics (*i.e.*, product, process, or developer-related) provides the best predictors; this is in line with recent findings reporting the importance of exploiting a combination of metrics to improve the performance of prediction models [11, 57].

Result 3: Developer-related factors are those that generally provide the highest contribution to the prediction of defective files within commits. Similarly, also the amount of lines added influences the prediction.

4.4. **RQ4.** *How much effort can be saved using a fine-grained just-in-time defect prediction model with respect to a standard just-in-time model?*

This analysis is intended to provide evidence on the effort developers can save using our model to guide the inspection of commits for defects. We consider the state of the art *just-in-time* model and the optimal results for comparison. Figure 2 plots the effort-based cumulative lift charts of the experimented techniques. The devised *fine-grained just-in-time* solution presents a larger curve with respect to the technique of Kamei et al. [34]: this confirms the ability of our model to work better than the baseline,

and thus it can optimize the effort required by a developer to locate actual defects. For instance, our results show that 54% of all defects can be identified by investigating an effort of 24% in terms of lines of code to inspect. With the same budget, only 27% of them can be found by relying on the technique by Kamei et al. [34]. This observation is confirmed when considering the P_k metric, that is equal to 0.43. Thus, the devised technique seems to represent a more viable solution for predicting defects at commit-level. A clear example can be found in the JRUBY project and is represented by the pull request #4371⁸ where two reviewers needed to inspect four committed files. In one of these, *i.e.*, the `mx_jruby.py` file composed of 171 lines of code, a defect was identified. On this commit, our model correctly marked the file as defective. At the same time, the model by Kamei *et al.* also correctly pointed out the defectiveness of the commit. However, while the effectiveness of the models is the same their code review cost is largely different: using the proposed model a reviewer should have focused on the `mx_jruby.py` file only, while using the baseline one she should have potentially investigated all the four files, leading to the analysis of a total of 1,947 LOCs (*i.e.*, +1,776 LOCs).

If we consider the differences between our technique and the optimal model, the P_{opt} is 0.31. This means that, as expected, the optimal model outperforms ours. Nonetheless, we can also see that the difference is closer when considering a reduced effort budget, *i.e.*, in cases where developers have limited time to dedicate to defect fixing activities. For example, let consider a hypothetical limited budget of 10%: in this case, using our technique, it is possible to identify 23% of the defective files, while with the optimal technique 35%. This indicates that, at least in the first phases, our model can be considered as a valid solution to speed up the identification of defects. At the same time, we argue that more research on the topic would be needed, as there is still room for further improvements.

When considering the individual projects, we observe that the BUGZILLA and OPENJPA projects follow a different trend. Further analyzing these cases, we found that the limited improvement with respect to the baseline was due to the characteristics of the commits in that repositories. Even though in \mathbf{RQ}_1 we found that average commit defectiveness ratio of the systems was 37% in both the cases, often the commits on these repositories contain few resources, *i.e.*, the average number of resources per commit was 2.3 and 3.4 for the two systems, respectively. This aspect limited the difference in the inspection costs achieved by the experimented models since in cases of defective commits composed of few resources the lines of code to inspect with the two models is similar. Nevertheless, also in these systems, the *fine-grained* solution outperforms the baseline: this result indicates that the proposed model may be useful also on systems which follow restrictive commit

policies or whose developers tend to commit fewer changes (*e.g.*, to avoid the introduction of tangled changes [27]).

Finally, it is worth remarking that the results achieved on the entire set of defective commits were also confirmed when considering partially and fully defective commits independently.

Result 4: Considering an effort-budget of 24%, 54% of the defects can be identified by our technique. In comparison with the state of the art, we observe that the devised technique represents a more viable solution to locate defects at commit-time. Our model improves upon the baseline also with a few files per commit.

5. Conclusion

Many defect prediction models have been proposed to locate defect-prone files or commits exploiting *long-term* or *short-term* techniques, respectively. Nevertheless, such models suffer from limitations due to the coarse-grained granularity of the predictions performed, which hinder their practical applicability (*e.g.*, in code review). For this reason, we investigated the possibility to devise a fine-grained *just-in-time* defect prediction model to locate defective files contained in a commit. The study considered 10 open-source systems written in different programming languages and having different size and scope. In total we analyzed 160,515 commits of which 35,496 defective.

The main contributions made by this paper are:

1. An empirical validation aimed at understanding the prominence of partially defective commits, *i.e.*, commits containing both defective and non-defective files on a set of 10 different open source software projects. The results highlight that almost half of defective commits contain both defect-inducing and defect-free changes.
2. A fine-grained *just-in-time* defect prediction model and its empirical evaluation, which showed performance up to 82% in terms of AUC-ROC.
3. An assessment of the *cost-effectiveness* of our model and its comparison with the standard *just-in-time* model proposed by Kamei et al. [34], with evidence that our model is more cost-effective.
4. An online appendix [66] that reports all the additional analyses mentioned in the paper.

Based on the results, our future agenda includes the replication of our study on a larger set of systems, possibly performing an in-depth study in an industrial context. At the same time, future studies can be designed and conducted to investigate (i) the role of other independent variables, *e.g.*, those reported by McIntosh et al. [48], on the performance of fine-grained defect prediction, (ii)

⁸<https://github.com/jruby/jruby/pull/4371> (associated with the commit 9c921e)

the model in the context of cross-project defect prediction, and (iii) the benefits provided by the usage of personalized defect prediction [29, 86] as well as more sophisticated ensemble techniques [87]. Moreover, we plan to evaluate the extent to which standard just-in-time approaches working at commit-level can be combined with the fine-grained solution we proposed, *e.g.*, through a multi-stage classification process where the defective commits are identified first and then the specific defective files are detected. Furthermore, the effectiveness of our model should be evaluated in-field, through a controlled study with practitioners to incorporate in our model some of the guidelines suggested by Lewis et al. [42] to make defect prediction more actionable in practice and support human activities (*e.g.*, by introducing a graphical user interface supporting code reviewers when diagnosing defect-prone code components).

Acknowledgments

This project has received funding from the European Union’s Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No 642954. Bacchelli and Palomba gratefully acknowledge the support of the Swiss National Science Foundation through the SNF Project No. PP00P2_170529.

References

- [1] ABDI, H. 2007. Bonferroni and šidák corrections for multiple comparisons. *Encyclopedia of measurement and statistics* 3, 103–107.
- [2] ARISHOLM, E., BRIAND, L. C., AND JOHANNESSEN, E. B. 2010. A systematic and comprehensive investigation of methods to build and evaluate fault prediction models. *Journal of Systems and Software* 83, 1, 2–17.
- [3] BACCHELLI, A. AND BIRD, C. 2013. Expectations, outcomes, and challenges of modern code review. *Proceedings - International Conference on Software Engineering*, 712–721.
- [4] BAEZA-YATES, R. AND RIBEIRO-NETO, B. 1999. *Modern Information Retrieval*. Addison-Wesley.
- [5] BARNETT, J. G., GATHURU, C. K., SOLDANO, L. S., AND MCINTOSH, S. 2016. The relationship between commit message detail and defect proneness in java projects on github. In *Proceedings of the 13th International Conference on Mining Software Repositories*. MSR ’16. ACM, New York, NY, USA, 496–499.
- [6] BASILI, V. R., BRIAND, L. C., AND MELO, W. L. 1996. A validation of object-oriented design metrics as quality indicators. *IEEE Transactions on software engineering* 22, 10, 751–761.
- [7] BAUM, T., SCHNEIDER, K., AND BACCHELLI, A. 2017. On the optimal order of reading source code changes for review. In *33rd International Conference on Software Maintenance and Evolution*. ICSME 2017. 329–340.
- [8] BELL, R., OSTRAND, T., AND WEYUKER, E. 2013. The limited impact of individual developer data on software defect prediction. *Empirical Software Engineering* 18, 3, 478–505.
- [9] CANFORA, G., LUCIA, A. D., PENTA, M. D., OLIVETO, R., PANICHELLA, A., AND PANICHELLA, S. 2015. Defect prediction as a multiobjective optimization problem. *Softw. Test. Verif. Reliab.* 25, 4, 426–459.
- [10] CASS, S. AND BULUSU, P. 2018. Interactive: The top programming languages 2018. *IEEE Spectrum* <https://spectrum.ieee.org/static/interactive-the-top-programming-languages-2018>.
- [11] CATOLINO, G., PALOMBA, F., DE LUCIA, A., FERRUCCI, F., AND ZAIDMAN, A. 2018. Enhancing change prediction models using developer-related factors. *Journal of Systems and Software* 143, 14–28.
- [12] CHIDAMBER, S. R. AND KEMERER, C. F. 1994. A metrics suite for object oriented design. *IEEE Transactions on software engineering* 20, 6, 476–493.
- [13] DA COSTA, D. A., MCINTOSH, S., SHANG, W., KULESZA, U., COELHO, R., AND HASSAN, A. E. 2017. A framework for evaluating the results of the szz approach for identifying bug-introducing changes. *IEEE Transactions on Software Engineering* 43, 7, 641–657.
- [14] D’AMBROS, M., LANZA, M., AND ROBBES, R. 2012. Evaluating defect prediction approaches: a benchmark and an extensive comparison. *Empirical Software Engineering* 17, 4-5, 531–577.
- [15] DEVIJVER, P. A. AND KITTLER, J. 1982. *Pattern Recognition: A Statistical Approach*.
- [16] DI NUCCI, D., PALOMBA, F., OLIVETO, R., AND DE LUCIA, A. 2017. Dynamic selection of classifiers in bug prediction: An adaptive method. *IEEE Transactions on Emerging Topics in Computational Intelligence* 1, 3, 202–212.
- [17] EL EMAM, K., MELO, W., AND MACHADO, J. C. 2001. The prediction of faulty classes using object-oriented design metrics. *Journal of Systems and Software* 56, 1, 63–75.
- [18] EYOLFSON, J., TAN, L., AND LAM, P. 2011. Do time of day and developer experience affect commit bugginess? In *Proceedings of the 8th Working Conference on Mining Software Repositories*. ACM, 153–162.
- [19] FREUND, Y. AND MASON, L. 1999. The alternating decision tree learning algorithm. In *Proceedings of the Sixteenth International Conference on Machine Learning*. ICML ’99. Morgan Kaufmann Publishers Inc., 124–133.
- [20] FRIEDMAN, J., HASTIE, T., AND TIBSHIRANI, R. 2001. *The elements of statistical learning*. Vol. 1. Springer series in statistics Springer, Berlin.
- [21] GRABMEIER, J. L. AND LAMBE, L. A. 2007. Decision trees for binary classification variables grow equally with the gini impurity measure and pearson’s chi-square test. *Int. J. Bus. Intell. Data Min.* 2, 2, 213–226.
- [22] GRAVES, T. L., KARR, A. F., MARRON, J. S., AND SIY, H. 2000. Predicting fault incidence using software change history. *IEEE Transactions on software engineering* 26, 7, 653–661.
- [23] GYIMÓTHY, T., FERENC, R., AND SIKET, I. 2005. Empirical validation of object-oriented metrics on open source software for fault prediction. *IEEE Transactions on Software Engineering (TSE)* 31, 10, 897–910.
- [24] HALL, T., BEECHAM, S., BOWES, D., GRAY, D., AND COUNSELL, S. 2012. A systematic literature review on fault prediction performance in software engineering. *IEEE Transactions on Software Engineering* 38, 6, 1276–1304.
- [25] HASSAN, A. E. 2009. Predicting faults using the complexity of code changes. In *Proceedings of the 31st International Conference on Software Engineering*. ICSE ’09. IEEE Computer Society, Washington, DC, USA, 78–88.
- [26] HERRAIZ, I., GONZALEZ-BARAHONA, J. M., AND ROBLES, G. 2007. Towards a theoretical model for software growth. In *Mining Software Repositories, 2007. ICSE Workshops MSR’07. Fourth International Workshop on*. IEEE, 21–21.
- [27] HERZIG, K., JUST, S., AND ZELLER, A. 2016. The impact of tangled code changes on defect prediction models. *Empirical Software Engineering* 21, 2, 303–336.
- [28] HINDLE, A., GERMAN, D. M., AND HOLT, R. 2008. What do large commits tell us?: a taxonomical study of large commits. In *Proceedings of the 2008 international working conference on Mining software repositories*. ACM, 99–108.
- [29] JIANG, T., TAN, L., AND KIM, S. 2013. Personalized defect prediction. In *Proceedings of the 28th IEEE/ACM International Conference on Automated Software Engineering*. IEEE Press, 279–289.
- [30] JIANG, Y., CUKIC, B., AND MENZIES, T. 2008. Can data transformation help in the detection of fault-prone modules? In *Proceedings of the 2008 Workshop on Defects in Large Software Systems*. DEFECTS ’08. ACM, 16–20.

- [31] JOHN, G. H. AND LANGLEY, P. 1995. Estimating continuous distributions in bayesian classifiers. In *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence*. UAI'95. Morgan Kaufmann Publishers Inc., 338–345.
- [32] KAMEI, Y., FUKUSHIMA, T., MCINTOSH, S., YAMASHITA, K., UBAYASHI, N., AND HASSAN, A. E. 2016. Studying just-in-time defect prediction using cross-project models. *Empirical Software Engineering* 21, 5, 2072–2106.
- [33] KAMEI, Y., MATSUMOTO, S., MONDEN, A., MATSUMOTO, K.-I., ADAMS, B., AND HASSAN, A. E. 2010. Revisiting common bug prediction findings using effort-aware models. In *Software Maintenance (ICSM), 2010 IEEE International Conference on*. IEEE, 1–10.
- [34] KAMEI, Y., SHIHAB, E., ADAMS, B., HASSAN, A. E., MOCKUS, A., SINHA, A., AND UBAYASHI, N. 2013. A large-scale empirical study of just-in-time quality assurance. *IEEE Transactions on Software Engineering* 39, 6, 757–773.
- [35] KENDALL, M. G. 1938. A new measure of rank correlation. *Biometrika* 30, 1/2, 81–93.
- [36] KHOMH, F., PENTA, M. D., GUÉHÉNEUC, Y.-G., AND ANTONIOL, G. 2012. An exploratory study of the impact of antipatterns on class change- and fault-proneness. *Empirical Softw. Engg.* 17, 3, 243–275.
- [37] KIM, S., WHITEHEAD JR, E. J., AND ZHANG, Y. 2008. Classifying software changes: Clean or buggy? *IEEE Transactions on Software Engineering* 34, 2, 181–196.
- [38] KONONENKO, O., BAYSAL, O., GUERROUJ, L., CAO, Y., AND GODFREY, M. W. 2015. Investigating code review quality: Do people and participation matter? In *Software Maintenance and Evolution (ICSME), 2015 IEEE International Conference on*. IEEE, 111–120.
- [39] LE CESSIE, S. AND VAN HOUWELINGEN, J. 1992. Ridge estimators in logistic regression. *Applied Statistics* 41, 1, 191–201.
- [40] LEHMAN, M. M. 1980. Programs, life cycles, and laws of software evolution. *Proceedings of the IEEE* 68, 9, 1060–1076.
- [41] LESZAK, M., PERRY, D. E., AND STOLL, D. 2002. Classification and evaluation of defects in a project retrospective. *Journal of Systems and Software* 61, 3, 173–187.
- [42] LEWIS, C., LIN, Z., SADOWSKI, C., ZHU, X., OU, R., AND WHITEHEAD JR, E. J. 2013. Does bug prediction support human developers? Findings from a Google case study. In *Proceedings of the 2013 International Conference on Software Engineering*. ICSE 2013. IEEE Press, 372–381.
- [43] LIAW, A. AND WIENER, M. 2002. Classification and regression by randomforest. *R News* 2, 3, 18–22.
- [44] MADEYSKI, L. AND KAWALEROWICZ, M. 2017. Continuous defect prediction: the idea and a related dataset. In *Proceedings of the 14th International Conference on Mining Software Repositories*. IEEE Press, 515–518.
- [45] MALHOTRA, R. 2015. A systematic review of machine learning techniques for software fault prediction. *Appl. Soft Comput.* 27, C, 504–518.
- [46] MATSON, J. E., BARRETT, B. E., AND MELLICHAMP, J. M. 1994. Software development cost estimation using function points. *IEEE Transactions on Software Engineering* 20, 4, 275–287.
- [47] MCCABE, T. J. 1976. A complexity measure. *IEEE Transactions on software Engineering* 4, 308–320.
- [48] MCINTOSH, S., KAMEI, Y., ADAMS, B., AND HASSAN, A. E. 2014. The impact of code review coverage and code review participation on software quality: A case study of the qt, vtk, and itk projects. In *Proceedings of the 11th Working Conference on Mining Software Repositories*. ACM, 192–201.
- [49] MENDE, T. AND KOSCHKE, R. 2009. Revisiting the evaluation of defect prediction models. In *Proceedings of the 5th International Conference on Predictor Models in Software Engineering*. ACM, 7.
- [50] MENDE, T. AND KOSCHKE, R. 2010. Effort-aware defect prediction models. In *Software Maintenance and Reengineering (CSMR), 2010 14th European Conference on*. IEEE, 107–116.
- [51] MENZIES, T., BUTCHER, A., COK, D., MARCUS, A., LAYMAN, L., SHULL, F., TURHAN, B., AND ZIMMERMANN, T. 2013. Local versus global lessons for defect prediction and effort estimation. *IEEE Transactions on Software Engineering* 39, 6, 822–834.
- [52] MOCKUS, A. AND WEISS, D. M. 2000. Predicting risk of software changes. *Bell Labs Technical Journal* 5, 2, 169–180.
- [53] MOONEN, L. 2001. Generating robust parsers using island grammars. In *Proceedings of the Eighth Working Conference on Reverse Engineering (WCRE'01)*. WCRE '01. IEEE Computer Society, Washington, DC, USA, 13–.
- [54] MOSER, R., PEDRYCZ, W., AND SUCCI, G. 2008. A comparative analysis of the efficiency of change metrics and static code attributes for defect prediction. In *Proceedings of the 30th international conference on Software engineering*. ACM, 181–190.
- [55] MUNSON, J. C. AND ELBAUM, S. G. 1998. Code churn: A measure for estimating the impact of code change. In *Software Maintenance, 1998. Proceedings., International Conference on*. IEEE, 24–31.
- [56] NAGAPPAN, N. AND BALL, T. 2005. Use of relative code churn measures to predict system defect density. In *Software Engineering, 2005. ICSE 2005. Proceedings. 27th International Conference on*. IEEE, 284–292.
- [57] NUCCI, D. D., PALOMBA, F., ROSA, G. D., BAVOTA, G., OLIVETO, R., AND LUCIA, A. D. 2017. A developer centered bug prediction model. *IEEE Transactions on Software Engineering PP*, 99, 1–1.
- [58] O'BRIEN, R. 2007. A caution regarding rules of thumb for variance inflation factors. *Quality & Quantity* 41, 5, 673.
- [59] OSTRAND, T. J., WEYUKER, E. J., AND BELL, R. M. 2005. Predicting the location and number of faults in large software systems. *IEEE Transactions on Software Engineering* 31, 4, 340–355.
- [60] PALOMBA, F., BAVOTA, G., DI PENTA, M., FASANO, F., OLIVETO, R., AND DE LUCIA, A. 2018. On the diffuseness and the impact on maintainability of code smells: a large scale empirical investigation. *Empirical Software Engineering* 23, 3, 1188–1221.
- [61] PALOMBA, F., PANICHELLA, A., ZAIDMAN, A., OLIVETO, R., AND DE LUCIA, A. 2017. The scent of a smell: An extensive comparison between textual and structural smells. *IEEE Transactions on Software Engineering*.
- [62] PALOMBA, F., ZANONI, M., FONTANA, F. A., DE LUCIA, A., AND OLIVETO, R. 2016. Smells like teen spirit: Improving bug prediction performance using the intensity of code smells. In *Proceedings of the 32nd International Conference on Software Maintenance and Evolution (ICSME 2016)*. Raleigh, USA: IEEE.
- [63] PALOMBA, F., ZANONI, M., FONTANA, F. A., DE LUCIA, A., AND OLIVETO, R. 2017. Toward a smell-aware bug prediction model. *IEEE Transactions on Software Engineering*.
- [64] PANICHELLA, A., ALEXANDRU, C. V., PANICHELLA, S., BACCHELLI, A., AND GALL, H. C. 2016. A search-based training algorithm for cost-aware defect prediction. In *Proceedings of the Genetic and Evolutionary Computation Conference 2016*. GECCO '16. ACM, New York, NY, USA, 1077–1084.
- [65] PARNAS, D. L. AND LAWFOR, M. 2003. The role of inspection in software quality assurance. *IEEE Transactions on Software Engineering* 29, 8, 674–676.
- [66] PASCARELLA, L., PALOMBA, F., AND BACCHELLI, A. 2018. Fine-Grained Just-in-Time Defect Prediction: Online Appendix. Url = <https://doi.org/10.5281/zenodo.1886063>.
- [67] PERNER, T. V. 1998. What's wrong with bonferroni adjustments. *BMJ: British Medical Journal* 316, 7139, 1236.
- [68] POSNETT, D., D'SOUZA, R., DEVANBU, P., AND FILKOV, V. 2013. Dual ecological measures of focus in software development. In *Proceedings of the 2013 International Conference on Software Engineering*. ICSE '13. IEEE Press, Piscataway, NJ, USA, 452–461.
- [69] QUINLAN, J. R. 1986. Induction of decision trees. *Mach. Learn.* 1, 1, 81–106.
- [70] RAHMAN, F. AND DEVANBU, P. 2013. How, and why, process metrics are better. In *Proceedings of the 2013 International Conference on Software Engineering*. IEEE Press, 432–441.

- [71] RAHMAN, F., POSNETT, D., HINDLE, A., BARR, E., AND DEVANBU, P. 2011. Bugcache for inspections: Hit or miss? In *Proceedings of the 19th ACM SIGSOFT Symposium and the 13th European Conference on Foundations of Software Engineering. ESEC/FSE '11*. ACM, New York, NY, USA, 322–331.
- [72] RAY, B., HELLENDORF, V., GODHANE, S., TU, Z., BACCHELLI, A., AND DEVANBU, P. 2016. On the "naturalness" of buggy code. In *Proceedings of the 38th International Conference on Software Engineering. ICSE '16*. ACM, New York, NY, USA, 428–439.
- [73] ROBNIK-SIKONJA, M. 2004. Improving random forests. In *ECML*. Vol. 3201. 359–370.
- [74] SAMMUT, C. AND WEBB, G. I., Eds. 2010. *Leave-One-Out Cross-Validation*. Springer US, Boston, MA, 600–601.
- [75] ŚLIWERSKI, J., ZIMMERMANN, T., AND ZELLER, A. 2005. When do changes induce fixes? In *ACM sigsoft software engineering notes*. Vol. 30. ACM, 1–5.
- [76] ŚLIWERSKI, J., ZIMMERMANN, T., AND ZELLER, A. 2005. When do changes induce fixes? *SIGSOFT Softw. Eng. Notes* 30, 4, 1–5.
- [77] SOMMERVILLE, I. 2006. *Software Engineering: (Update) (8th Edition) (International Computer Science)*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA.
- [78] SPADINI, D., ANICHE, M., STOREY, M.-A., BRUNTINK, M., AND BACCHELLI, A. 2018. When testing meets code review: Why and how developers review tests. to appear.
- [79] SPADINI, D., PALOMBA, F., ZAIDMAN, A., BRUNTINK, M., AND BACCHELLI, A. 2018. On the relation of test smells to software code quality. In *Proceedings of the International Conference on Software Maintenance and Evolution (ICSME)*. IEEE.
- [80] TAN, M., TAN, L., DARA, S., AND MAYEUX, C. 2015. Online defect prediction for imbalanced data. In *Proceedings of the 37th International Conference on Software Engineering-Volume 2*. IEEE Press, 99–108.
- [81] TUFANO, M., BAVOTA, G., POSHYVANYK, D., DI PENTA, M., OLIVETO, R., AND DE LUCIA, A. 2017. An empirical study on developer-related factors characterizing fix-inducing commits. *Journal of Software: Evolution and Process* 29, 1, e1797–n/a. e1797 JSME-15-0185.R2.
- [82] VAHABZADEH, A., FARD, A. M., AND MESBAH, A. 2015. An empirical study of bugs in test code. In *Software Maintenance and Evolution (ICSME), 2015 IEEE International Conference on*. IEEE, 101–110.
- [83] VAN DER MALSBERG, C. 1986. *Frank Rosenblatt: Principles of Neurodynamics: Perceptrons and the Theory of Brain Mechanisms*. Springer Berlin Heidelberg, 245–248.
- [84] WILLIAMS, C. AND SPACCO, J. 2008. Szz revisited: Verifying when changes induce fixes. In *Proceedings of the 2008 Workshop on Defects in Large Software Systems. DEFECTS '08*. ACM, New York, NY, USA, 32–36.
- [85] WU, R., ZHANG, H., KIM, S., AND CHEUNG, S.-C. 2011. Re-link: Recovering links between bugs and changes. In *Proceedings of the 19th ACM SIGSOFT Symposium and the 13th European Conference on Foundations of Software Engineering. ESEC/FSE '11*. ACM, New York, NY, USA, 15–25.
- [86] XIA, X., LO, D., WANG, X., AND YANG, X. 2016. Collective personalized change classification with multiobjective search. *IEEE Transactions on Reliability* 65, 4, 1810–1829.
- [87] YANG, X., LO, D., XIA, X., AND SUN, J. 2017. Tlel: A two-layer ensemble learning approach for just-in-time defect prediction. *Information and Software Technology* 87, 206–220.
- [88] YANG, X., LO, D., XIA, X., ZHANG, Y., AND SUN, J. 2015. Deep learning for just-in-time defect prediction. In *Software Quality, Reliability and Security (QRS), 2015 IEEE International Conference on*. IEEE, 17–26.
- [89] YANG, Y., ZHOU, Y., LIU, J., ZHAO, Y., LU, H., XU, L., XU, B., AND LEUNG, H. 2016. Effort-aware just-in-time defect prediction: Simple unsupervised models could be better than supervised models. In *Proceedings of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering. FSE 2016*. ACM, New York, NY, USA, 157–168.