

PRAGMATIC SOFTWARE MAINTENANCE:
TOOL AND TECHNIQUES TO SUPPORT SOFTWARE
MAINTENANCE AT COMMIT-TIME.

MATHIEU NAYROLLES

A RESEARCH PROPOSAL
IN
THE DEPARTMENT
OF
ELECTRICAL & COMPUTER ENGINEERING

PRESENTED IN PARTIAL FULFILLMENT OF THE REQUIREMENTS
FOR THE DEGREE OF DOCTOR OF PHILOSOPHY
CONCORDIA UNIVERSITY
MONTRÉAL, QUÉBEC, CANADA

APRIL 2016

© MATHIEU NAYROLLES, 2016

Abstract

Pragmatic Software Maintenance:

Tool and technics to support software maintenance at commit-time.

Mathieu Nayrolles, Ph.D.

Concordia University, 2016

The maintenance and evolution of complex software systems account for more than 70% software's life cycle. More than two decades of research have been conducted to improve our knowledge of these processes in terms of issue triaging, issue prediction, duplicate issue detection, issue reproduction and changes prediction. This research gave meaning to the millions of issues that can be found in project and revision management systems. Context-aware IDE and think tank in open source architecture, opened the path to approaches that support developers during their programming sessions by leveraging past knowledge and architectures.

However, these techniques are still not broadly adopted by praticionners in large software companies. The main problem is the lack of actionable intelligence. Indeed, such systems tend to be seen as black boxes that yield false positive by end-users. Moreover, actions to resolve an identified problem can be hard to identify or to apply.

In this research proposal, we present four approaches: (a) an online bug-fix search engine and API (**BUMPER**, Bug MetarePository for dEveloper and Researcher), (b) an approach to automatically reproduce bug submitted to bug-report systems (**JCHARMING**, Java CrasH Automatic Reproduction by directed Model checkING), (c) a recommendation system to propose the right auto-completion at the right time using actionnable intelligence (**RESEMBLE**, REcommendation System based on cochangeE Mining at Block LEvel) and, finally, (d) an approach to prevent bug insertion at commit-time by leveraging decade of open-source history (**BIANCA**, Bug Insertion ANticipation by Clone Analysis at commit time). We also propose a taxonomy of bugs. When combined into **pErICOPE** (Ecosystem Improve source COde during Programming session with real-time mining of common knowlEdge), these tools (i) provide the possibility to search related software artifacts using natural language, (ii) accurately reproduce field-crash in lab environment, (iii) recommend improvement or completion of block of code under edition and (iv) prevent the introduction of issues at commit time.

This proposal develops these ideas while highlighting remaining issues and the PhD schedule.

Contents

1	Introduction	1
1.1	Preliminaries	3
1.1.1	Version control systems	4
1.1.2	Project Tracking Systems	5
1.2	Thesis Contributions	7
1.2.1	Problems in the Literature	7
1.2.2	Thesis contributions	8
1.3	Outline	13
2	Related Work	14
2.1	Crash reproduction	14
2.2	Issue and source code relationships	17
2.3	Crash Prediction	17
2.4	Clone Detection	20
3	Aggregating versionning systems and project managements systems	22
3.1	BUMPER - Bug Meta-repository For Developers & Researchers	23
3.1.1	Data collection	23
3.1.2	Architecture	24
3.1.3	UML Metamodel	25
3.1.4	Features	26
3.1.5	Application Program Interface (API)	29
3.2	JCHARMING - Java CrasH Automatic Reproduction by directed Model checkING	30
3.2.1	Preliminaries	31
3.2.2	The JCHARMING Approach	34
3.2.3	Case studies	40

3.2.4	Results	41
4	Using clone detection for pragmatic software maintenance	47
4.1	Software maintenance at branching-time	49
4.2	Software maintenance at commit-time	51
4.2.1	PRECINCT: PREventing Clones INsertion at Commit Time	52
4.2.2	RESEMBLE - REcommendation System based on cochangE Mining at Block LEvel	65
4.3	Software maintenance at merge-time	68
5	A taxonomy to classify the research	73
5.0.1	Study Setup	74
5.0.2	Datasets	75
5.0.3	Study Design	76
5.0.4	Study result and discussion	79
6	Reflection on pragmatic software maintenance	88
7	Research Plan	90
7.1	Current State of Research	90
7.2	Current Contribution	90
7.3	Plan for short term work	91
7.4	Publication Plan	92
8	Conclusion	93
	Bibliography	95

List of Figures

1	Proportion of papers containing “Empirical Study” or “Mining software repository” with regards to the paper in Software quality indexed by Google Scholar	9
2	Proposed Architecture	10
3	PRECINCT sample output	11
4	Overview of the bumper database construction.	24
5	Overview of the bumper architecture.	25
6	Overview of the bumper meta-model.	26
7	Screenshot of https://bumper-app.com with “Exception” as research. . . .	30
8	A toy program under testing	32
9	A toy program under model checking	33
10	A toy program under directed model checking	33
11	Overview of JCHARMING.	34
12	Java InvalidActivityException is thrown in the Bar.Goo loop if the control variable is greater than 2.	35
13	Hypothetical example representing $bslice_{[entry \leftarrow f_0]}$ Vs. $\cup_{i=0}^n bslice_{[f_{i+1} \leftarrow f_i]} = \{f_0, f_1, f_2, z_2\}$	37
14	Data structure of a commit.	48
15	Data structure of two commits.	49
16	Two branches pointing on one commit.	49
17	Two branches pointing on two commits.	50
18	Pre-checkout hook	50
19	JCHARMING generated test imported into task branch	51
20	Overview of the PRECINCT Approach.	53
21	PRECINCT output when replaying commit 710b6b4 of the Monit system used in the case study.	59

22	PRECINCT Branching.	61
23	The BIANCA Approach	69
24	BIANCA warnings from April to August 2008 using the first normalization.	71
25	BIANCA warnings from April to August 2008 using the second normalization.	72
26	Class diagram showing the relationship between bugs and fixed	73
27	Proposed Taxonomy of Bugs	74
28	Data collection and analysis process of the study	75
29	Proportions of different types of bugs	79
30	Proportions of Types 1 and 3 versus Types 2 and 4 with respect to their severity in the Apache dataset.	81
31	Proportions of Types 1 and 3 versus Types 2 and 4 with respect to their severity in the Netbeans dataset.	82
32	Fixing time of Types 1 and 3 versus fixing time of Types 2 and 4.	84

List of Tables

1	Hypothetical BUMPER data	12
2	List of target systems in terms of Kilo line of code (KLoC), number of classes (NoC) and Bug # ID	40
3	Effectiveness of JCHARMING using directed model checking (DMC) and model checking (MC) in minutes	42
4	List of Target Systems in Terms of Files and Kilo Line of Code (KLOC) at current version and Language	60
5	Datasets	70
6	Datasets	76
7	Contingency table and Pearson's chi-squared tests	79
8	Proportion of bug types in amount and percentage	80
9	Pearson's chi squared p-values for the severity, the reopen and the duplicate factor with respect to a dataset	81
10	Proportion of each bug type with respect to severity.	83
11	Percentage and occurrences of bugs duplicated by other bugs and reopened with respect to their bug type and dataset.	84
12	Average fixing time with respect to bug type	85

Chapter 1

Introduction

Maintenance activities are known to be costly and challenging [Pre05]. Studies have shown that the cost of software maintenance can reach up to 70% of the overall cost of the software development process [HR02]. By “maintenance activity” we mean any change to software beyond its first release or iteration; i.e. development of software where there is already an existing system that is to be changed. This is the broadest possible meaning of the term. It includes adding new features, creating new iterations, as well as classic adaptive and corrective maintenance.

The difficulties encountered by maintainers are various. In large software project, difficulties are partially attributable to the fact that the software are created, distributed and maintained globally. Developers and users are scattered across the globe.

Consequently, maintainers can hardly have complete knowledge over the software at hand. Yet, such knowledge will be an asset to avoid known mistakes, programming errors and improve the software.

In the last decade, source code revision control system and project tracking systems have grown to contain hundreds of thousands of revision, bug and crash report per project. Naturally, this plethora of data pushed researchers across the world to conduct hundreds of studies in several active research fields: bug reproduction, bug triaging, duplicated bug identification, bug comprehension, bug re-production.

In order to support software maintenance, researchers and developers started to mine software repository. Software repository contains the each modification done to the software and related issues.

Mining software repository—the science of interpreting software artifacts—is perhaps one of the most active research field today. The reason is that their analysis provides useful insight that can help with many maintenance activities such as bug fixing [WPZZ07,

SKP14], bug reproduction [AKE08, JO12, Che13], fault analysis [JLL⁺12, JO13], bug prediction [Hov07], etc.

A great length of effort have been conceded by the mining software repository community to create tools and approaches ranging from simple text-pattern matching to complex prediction models. This tools are often automatically tested and verified in laboratory[LLS⁺13]. Indeed, manually testing and validating industrial-sized repository for each iteration of a given algorithm is counterproductive. This automation of experiment allows researchers to improve their approaches and compare tools to each others. The current state-of-the-art approach can always be improved by new and more complex technics[HP04].

Yet, one question remains: Is it usefull to **humans** maintainers?

This simple question has been studied by researchers and companies alike (noticeable examples [LLS⁺13, FM15, LWA07, APM⁺07, AP08, JSMHB13, Nor13, LvdH11]). Of course, results of such studies are not binary and can be difficult to interpret. Nevertheless, it has been discovered that although human maintainers agree that such tools are beneficial, false positives and the way in which the warnings are presented, among other things, are barriers to use[JSMHB13]. In addition, human maintainers agree to some of the characteristics that maintenance-oriented tools should have [HP04, LvdH11, LLS⁺13]:

- Actionable messages. Presenting a warning about bug-proneness of a given line is not enough. Clear actions to improve the source code should be provided
- Obvious reasoning. The conditions that led to a given warning should be understandable by the maintainer. If the conditions are hidden in complex statistical models, then maintainers cannot review them and will find it difficult to trust the tool.
- Scaling. Industrial sized project contains thousand files and dependencies which can be updated many times a day. Maintenance-oriented tools should not hinder the productivity of maintainers.
- Contextualization. Warnings and messages should always be contextualized with respect to the project at hand and not generic rules.
- Integration. Developers and maintainers are overwhelm by the amount of existing tools. Yet, their daily use three different kind of tools. An integrated development environment, a versionning system and a project tracking system to produce, version and manage their software, respectively. Maintenance-oriented tools should fit in the existing ecosystem rather than complexifying the deployment process.

In the following sections, we present some preliminaries to the mining software repository field. Then, we describe the thesis contribution and this proposal's outline.

1.1 Preliminaries

Software maintenance, comprehension, evolution, specifications and testing are research areas overlapping each other in terms of terminology.

In this proposal, we will use precise set of definitions. We do not claim ownership of these definitions, they have been established using various resources [ALRL04, Pra01, Bur06, RGK90, WAC12].

We limit software maintenance to the following three artifacts:

- Bug report: A bug report describes a behaviour observed on the field and considered abnormal by the reporter. Bug reports are submitted manually to bug report systems (bugzilla/jira). There is no mandatory format to report a bug, nevertheless, it should have : Version of the software / OS / Platform used, steps to reproduce, screenshots, stack trace and anything that could help a developer to assess the internal state of the software system.
- Crash report: A crash report is the last action a software system does before crashing. Crash reports are automatic (they have to be implemented into the software system by developer) and contain data (that can be proprietary) to help developer understand the crash (e.g memory dump, ...).
- Tasks: A task is a new feature, or the improvement of an existing feature, to be implemented in a future release of the software.

These artifacts are produced in response to the following phenomenons:

- Software Bug: A software bug is an error, flaw, failure, defect or fault in a computer program or system that causes it to violate at least one of its functional or nonfunctional requirements.
- Error: An error is a mistake, misconception, or misunderstanding on the part of a software developer. Fault/defect: A fault (defect) is defined as an abnormal condition or defect at the component, equipment, or subsystem level which may lead to a failure. A fault (defect) is not final (the system still works) and does not prevent a given feature to be accomplished. A fault (defect) is a deviation (anomaly) of the healthy system that can be caused by an error or external factors (hardware, third parties, ...).

- Failure: The inability of a software system or component to perform its required functions within specified requirements.
- Crash: The software system encountered a fault (defect) that triggered a fatal failure from which the system could not recover from/overcome. As a result, the system stopped.

In the remaining of this section, we introduce the two types of software repositories: version control and project tracking system.

1.1.1 Version control systems

Version control consists in maintaining the versions of files — such as source code and other software artifacts. This activity is extremely complex and cannot be done by hand on real world project. Consequently, numerous softwares have been created to help practitioners manage the version of their software artifacts. Each evolution of a software is a version¹ (or revision) and each version (revision) is linked to the one before through modifications of software artifacts. These modifications consist in updating, adding or deleting software artifacts. They can be referred as `diff`, `patch` or `commit`². Each `diff`, `patch` or `commit` have the following characteristics:

- Number of Files: The number of software artifacts that have been modified, added or deleted.
- Number of Hunks: The number of consecutive blocks of modified, added or deleted lines in textual files. Hunks are useful to determine, in each file, how many different places the developer has modified.
- Number of Churns: The number of lines modified. However, the churn value for a line change should be at least two as the line had to be deleted first and then added back with the modifications.

Modern version control systems also support branching. A branch is a derivation in the evolution that contain a duplication of the source code so that both versions can be modified in parallel. Branches can be reconciled with a merge operation that merge modification from the two branches. This operation is completely automated at the exception of merge conflicts that arise when both branches contain modification on the same line. Such conflicts

¹Software version is not to be confused with the version of a software which refer to the shipping of a final product to customers.

²These names are not to be used interchangeably as difference exists.

cannot be reconciled automatically and have to be dealt with by developer. This allows for a greater agility among developers as changes one branch do not affect other developers of other branches.

Branching have been used for more than testing hazardous refactoring or testing framework upgrades. Indeed, task branching is an agile branching strategy where a new branch is created for each task[Mar09]. It is common to see branch names `123_implement_X` where 123 is the #id of tasks X given by the project tracking systems. Project tracking systems are presented in Section1.1.2.

Providers

In this proposal, we mainly refer to three version control systems: **Svn**, **Git** and, to a lesser extent, **Mercurial**. **SVN** is distributed by the Apache foundation and is a centralized concurrent version system that can handle conflict in the different versions of different developers and it is widely used. At the opposite, **Git** is a distributed revision control system — originally developed by Linus Torvald — where revisions can be kept locally for a while and then shared with the rest of the team. Finally **Mercurial** is also a distributed revision system, but share a lot of concepts with **Svn**. Consequently, it will be easier for people used to **Svn** to switch to a distributed revision system if they use **Mercurial**.

1.1.2 Project Tracking Systems

Project tracking systems allow end-users to directly create bug reports (BRs) to report unexpected system behaviour, manager can create tasks to drive the evolution forward and crash report (CRs) can be automatically created. These systems are also used by development teams to keep track of the modification induced by bug and to crash reports, and keep track of the fixes.

The life cycle of a report is as follows: After an it is s submitted by an end-user, it is set to the **UNCONFIRMED** state until it receives enough votes or that a user with the proper permissions modifies its status to **NEW**. The report is then assigned to a developer to be fixed. When the report is in the **ASSIGNED** state, the assigned developer(s) starts working on the report. A fixed report moves to the **RESOLVED** state. Developers have five different possibilities to resolve a report: **FIXED**, **DUPLICATE**, **WONTFIX**, **WORKSFORME** and **INVALID**.

- **RESOLVED/FIXED**: A modification to the source code has been pushed, i.e., a changeset (also called a patch) has been committed to the source code management system and fixes the root problem described in the report.

- **RESOLVED/DUPLICATE:** A previously submitted report is being processed. The report is marked as duplicate of the original report.
- **RESOLVED/WONTFIX:** This is applied in the case where developers decide that a given report will not be fixed.
- **RESOLVED/WORKSFORME:** If the root problem described in the report cannot be reproduced on the reported OS / hardware.
- **RESOLVED/INVALID:** If the report is not related to the software itself.

Finally, the report is **CLOSED** after it is resolved. A report can be reopened (sent to the **REOPENED** state) and then assigned again if the initial fix was not adequate (the fix did not resolve the problem). The elapsed time between the report marked as the new one and the resolved status are known as the *fixing time*, usually in days. In case of task branching, the branch associated with the report is marked as ready to be merged. Then, the person in charge (quality assurance team, manager, ect...) will be able to merge the branch with the mainline. If the report is reopened: the days between the time the report is reopened and the time it is marked again as **RESOLVED/FIXED** are cumulated. Reports can be reopened many times.

Tasks, follow a similar life cycle at the exception of the **UNCONFIRMED** and **RESOLVED** states. Tasks are created by management and do not need to be confirmed in order to be **OPEN** and **ASSIGNED** to developers. When a task is complete, it will not go to the **RESOLVED** state, but to the **IMPLEMENTED** state. Bug and crash reports are considered as problems to eradicate in the program. Tasks are considered as new features or amelioration to include in the program.

Reports and tasks can (and must according to [BJS⁺08]) have a severity. The severity is a classification to indicate the degree of impact on the software. The possible severities are

- **blocker:** blocks development and/or testing work.
- **critical:** crashes, loss of data, severe memory leak.
- **major:** major loss of function.
- **normal:** regular report, some loss of functionality under specific circumstances.
- **minor:** minor loss of function, or other problem where easy workaround is present.
- **trivial:** cosmetic problem like misspelled words or misaligned text.

The relationship between an report or a task and the actual modification can be hard to establish and it has been a subject of various research studies (e.g., [ACC⁺02, BBR⁺10, WZKC11]). This reason is that they are in two different systems: the version control system and the project tracking system. While it is considered a good practice to link each report with the source code revision system by indicating the report *#id* on the modification message, more than half of the reports are not linked to a modification[WZKC11].

Providers

We have collected and plan to collect data from four different project tracking systems: *Bugzilla*, *Jira*, *Github* and *Sourceforge*. *Bugzilla* belongs to the Mozilla foundation and has first been released in 1998. *Jira*, provided by Altassian, has been released 14 years ago, in 2002. *Bugzilla* is 100% open source and it's difficult to estimate how many project uses it. However, we can, without any risks envision that it owns a great share of the market as major organizations such as Mozilla, Eclipse and the Apache Software Foundation uses it. *Jira*, in the other hand, is a commercial software — with a freemium business model — and Altassian claims that they have 25,000 customers over the world.

Github and *Sourceforge* are different from *Bugzilla* and *Jira* in a sense that they were created as source code revision system and evolve, later on, to add project tracking capabilities to their softwares. This common particularity have the advantage to ease the link between reports and source code.

1.2 Thesis Contributions

In this section we present the problems in the literature (Section 1.2.1) and the research challenges we face in our work (Section ??). Sections ?? and 1.2.2 present the scope and the contributions of this research.

1.2.1 Problems in the Literature

- **Problem 1:** The literature contains numerous papers about tools that improve the overall software quality with static [Dan00, Bur03, Hov07, MGD10] and dynamic [NMHIL, NMV13, Pal13] analysis. At a few exceptions (such as [LvdH11, MBFV13]), these tools do not account for the context of project. By context, we refer to past versions, reports, comments or any direct and undirect software artifact that could help provide precise step to resolve a given warning.

- **Problem 2:** Tools aiming to ease maintenance processes do not integrate themselves seamlessly in the five steps of maintenance: acknowledge a report or a task, create a branch, fix or enhance source code, test, mark branch as ready [KKI02a, NMHIL, NMV13, JFG09, Che13, Gav13, JO12, NAW⁺08, DZM, NBMV15, DMT13, JMSG07, CHA09]. Indeed most of the tools are plain-old tools, at the exception of IDE plugins [KKI02a, Hov07], that maintainers have to install and update manually. Then, for each completed task, the tools have to be started and their results analyzed: false positive discarded and warning messages decrypted. This among other things, are barriers to their broad adoption [JSMHB13, LLS⁺13].
- **Problem 3:** As shown by Figure 1, the proportion of empirical studies and studies based on mining software repositories regarding to software quality has been increasing exponentially since 2005 ([KZWG11, LNH⁺11, SLKJ11, BN11, TSL12, ZNGM12, SNH13, CNSH14, MANH14, HNH15] are some noticeable examples). Yet, the field lack of a clear taxonomy to be able to compare approaches efficiently and replicate experiments [Has08, GHH⁺09].

1.2.2 Thesis contributions

Figure 2 depicts our proposed solution for fulfilling our objectives. More specifically, this figure describes how our different contribution interface themselves in a classical maintenance process supported by feature branching.

During maintenance process supported by feature branching. First, a developer will browse the project management system and acknowledge reports or tasks assigned to him. Assigning reports or issues to developers is known as triaging. There have been research on how to perform automated triaging that determines which developer is the best suited to accomplish it [SKP14, TNAKN11b, BvdH13]. We do not propose ameliorations in the field of triaging.

The second step will be to create a feature branch for the report or the issue at hand. Creating a feature branch is a straightforward command, for example *git checkout -b 123.task.name* for git. We can identify which task or report is being tackled by the developer with the name of the branch and automatically fetch relevant informations such as steps to reproduce the reports or similar reports closed in the past. We do that using a pre-checkout and one of our contribution: BUMPER (Bug Metarepository for Research and Developers).

A hook is a process that can be implemented to receive the modification to the source code done by the developer. Hooks are custom scripts set to fire off before and after

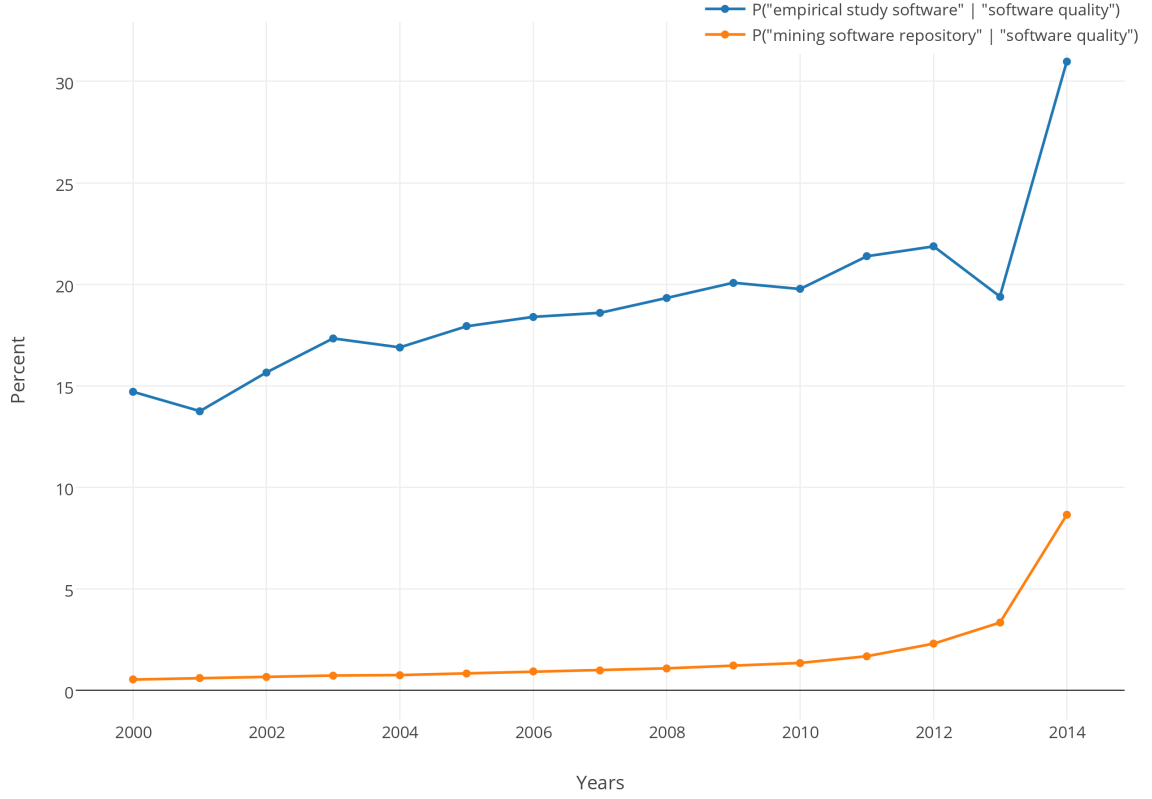


Figure 1: Proportion of papers containing “Empirical Study” or “Mining software repository” with regards to the paper in Software quality indexed by Google Scholar

versionning actions occur. There are two types of hooks: client-side and server-side. Client-side hooks are triggered by operations such as committing and merging, whereas server-side hooks run on network operations such as receiving pushed commits. These hooks can be used for all sorts of reasons such as compliance to coding rules or automatic run of unit test suites.

Leveraging the hook-feature, we can fetch information from a central repository called **BUMPER** (Bug Metarepository for Research and DevelopeRs) about the task or report at hand. **BUMPER** is a meta-repository that makes reports, tasks and related source code searchable using a structured API. **BUMPER** contains several adaptator to fetch new reports and tasks submitted to project management systems. It updates itself every night.

When a report is created, **JCHARMING** (Java CrasH Automatic Reproduction by directed Model checkING) will fetch the content of the issue and try to create a scenario to reproduce

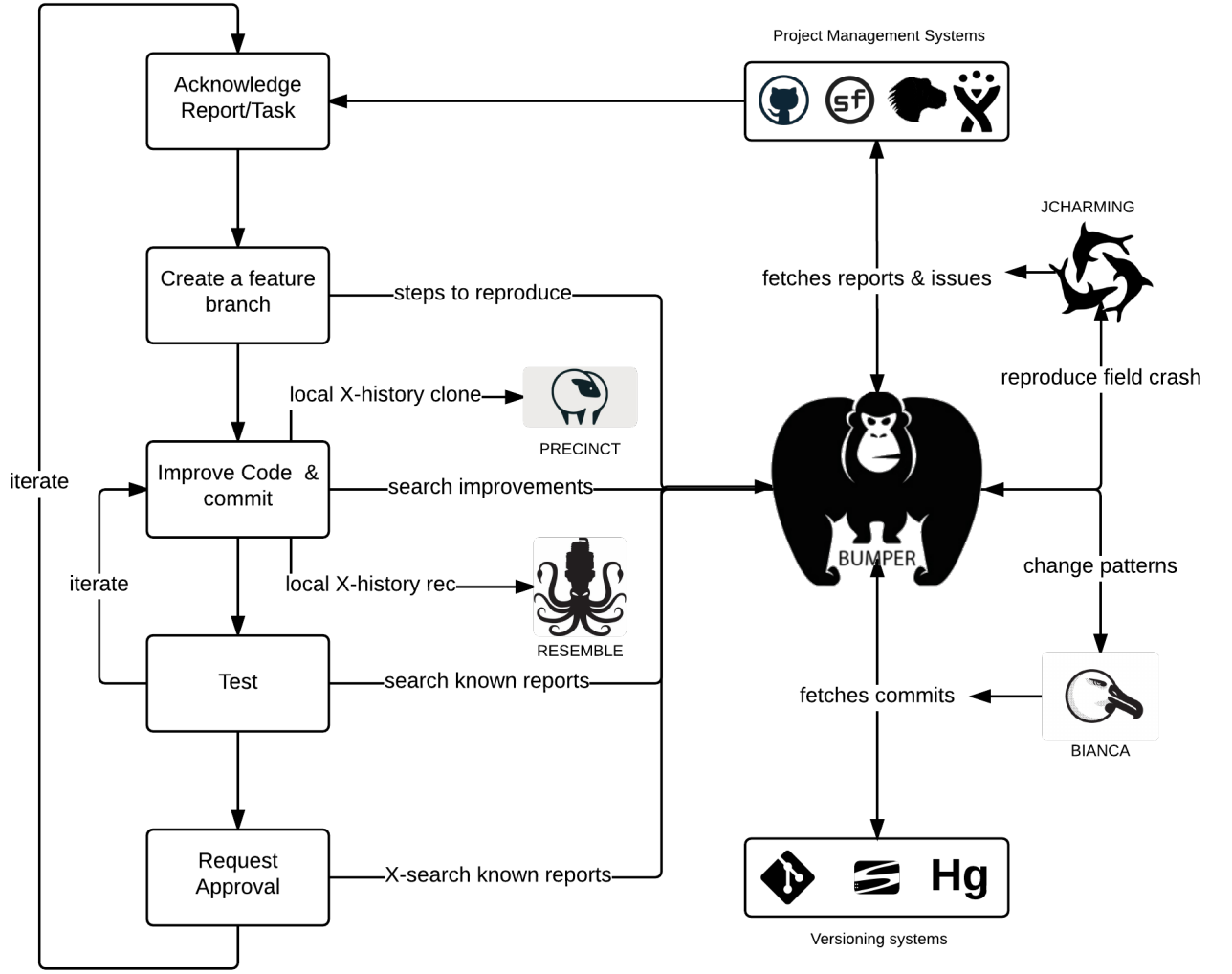


Figure 2: Proposed Architecture

the on-field crash. In case of success the scenario is stored in **BUMPER**.

Consequently, at branching-time, we can add the scenarios to reproduce the report directly inside the source code and recommend similar reports. We provide *accurate* and *contextual* information in a transparent manner for the developer. Indeed, the reproduction of a field-crash is binary, reproduced or not. In addition, this process occurs directly inside the versioning system that the maintainer already uses.

At this point, the maintainers start editing the code and commit their changes to their local branch. For each commit, we perform two operations on the local history: clones insertion and co-change recommendation using **PRECINCT** (PREventing Clones INsertion at Commit Time) and **RESEMBLE** (REcommendation System based on cochangeE Mining at Block LEvel). Once again, we leverage hooks, and more specifically pre-commits hooks to

do so. For each commit, we are able to prevent the insertion of a clone before it reaches to central repository and recommend changes based on mining similar co-changes at block level. This is possible because the local branch of the developer contains all the history of the project, from its beginning to the present. At commit-time, we provide **contextualised** advice directly into the source versionning system to improve the software before the changes reach the central repository. Figure 3 presents an output of PRECINCT, directly in the source versionning system.

```

Terminal - math@math-hp: ~/workspace/monit
File Edit View Terminal Tabs Help
math@math-hp ~/workspace/monit (git)-[master] % git comm
it -m "Use the more reliable fcntl function instead of i
octl"
*****
*
* PRECINCT (PREventing Clones INsertion at Commit Time)
*
*****
*
Following File(s) insert clones
libmonit/src/system/Net.c
(I) Inspect (D) disregard (R) remove from commit.

```

Figure 3: PRECINCT sample output

When maintainers are confident enough to share their changes, they will push them to the central repository. Using a server-side push hook, we can analyse the source code pushed by the maintainer. Our last tool, **BIANCA** (Bug Insertion ANTicipation by Clone Analysis at psuh time), analyzes the change pattern performed by the maintainer in order to identify if it this pattern is likely to introduce a default in the application. **BIANCA** does this by analyzing every changes ultimately led to a report. In case of a match, **BIANCA** displays the matching source code to the maintainer. The matching source code that led to a report and its subsequent modifications to fix that report are actionnable messages. Indeed, the maintainers are informed their modifications raised a warning and they can see how to transform their code, based on the provided example, if they decide to take action. Futhermore, the reasoning behind **BIANCA**'s warning is obvious: the submitted modification matches modification known to have introduced a deffect. At push-time, the response time of any analysis shall be as short as possible: maitainer are waiting for their modification synchronizied with the central repository. Hence, **BIANCA** only checks for matches inside the

current project history. At merge-time, however, we can compute more expensive operation. Indeed, the maintainer submits the modification for integration to the mainline and the modification will be reviewed by other developers in code reviews. Consequently, **BIANCA** checks for matching example in all known projects. **BIANCA** uses the code normalizations and text-based clone detection technics [Cor06b, ROY09, CR11].

All these approaches have been designed with the aim to reduce known barriers of adoption, namely, actionable messages, obvious reasoning, scaling and contextualization[JSMHB13, HP04, LvdH11, LLS⁺13].

As a motivating example, we draft the following scenario. Table 1 presents hypothetical data stored in **BUMPER** in terms of sequence #id, sequence of code blocks, a flag to know if a said sequence introduced an issue in a given system and step to reproduce the issue if any.

Seq #ID	Language #ID	Blocks	Root of Issue	Steps to reproduce
1	1	A-A-B-C-A-A	Yes	E-F-G
2	1	A-A-B-C	No	-
3	2	D-E-A-C	No	-

Table 1: Hypothetical **BUMPER** data

During a maintenance activity, let's assume that a developer has committed $A - B - C$ to its local history. Then, **RESEMBLE** will recommend to transform the current code to $A - A - B - C$ as it seems to be the right thing to do. If the developer follows **RESEMBLE** recommendation and then, adds another two A s, the sequence is transformed to $A - A - B - C - A - A$. If the developer commit its changes, **BIANCA** will raise a warning saying that this sequence is known to be the root of an issue and invite the developer to execute the steps $E - F - G$ — that were produced by **JCHARMING** — in order to see if s/he did introduce a defect. Moreover, **BIANCA** will take the time to compare $A - A - B - C - A - A$ and $D - E - A - C$, at merge-time, using our normalization algorithms even if they are not in the same programming language. Finally, when a new report is submitted, **BUMPER** indexes it and **JCHARMING** tries to reproduce it and update the step to reproduce part of **BUMPER**.

We can envision the potential of such a system and its complexity, knowing that it would contain millions of issues, hundreds of thousand projects, dozens of programming languages and will help developers leveraging the knowledge of other developers.

In addition to these contributions, we want to provide a way to classify the research related to defaults in software maintenance. Such taxonomy will allow researchers to specialize on classes of defect and to be able to compare accurately their approaches with

other, very much like the taxonomy that exists for clones detection [Cor].

1.3 Outline

Chapter 2

Related Work

2.1 Crash reproduction

In his Ph.D thesis [Che13], Chen proposed an approach named STAR (Stack Trace based Automatic crash Reproduction). Using only the crash stack, STAR starts from the crash point and goes backward towards the entry point of the program. During the backward process, STAR computes the required condition to reach the crash point using an SMT (Satisfiability Modulo Theories) solver named Yices [DM06]. The objects that satisfy the required conditions are generated and orchestrated inside a JUnit test case. The test is run and the resulting crash stack is compared to the original one. If both match, the bug is said to be reproduced. When applied to different systems, STAR achieved 60% accuracy.

Jaygarl et al. [JKXC10] created OCAT (Object Capture based Automated Testing). The authors' approach starts by capturing objects created by the program when it runs on-field in order to provide them in an automated test process. Indeed the coverage of automated tests is often low due to the lack of correctly constructed objects. Also, the objects can be mutated by means of evolutionary algorithms. These mutations target primitive fields in order to create even more objects and therefore improve the code coverage once more. While not targeting the reproduction of a bug, OCAT is a well-known approach and was used as the main mechanism for bug reproduction.

Narayanasamy et al. [NPC05] proposed BugNet, a tool that continuously records program execution for deterministic replay debugging. According to the authors, the size of the recorded data needed to reproduce a bug with high accuracy is around 10MB. This recording is then sent to the developers and allows the deterministic replay of a bug. The

authors argued that, with nowadays Internet bandwidth, the size of the recording is not an issue during the transmission of the recorded data, however, the instrumentation of the system is problematic since it slows down considerably the execution.

Jin et al. [JO12] proposed BugRedux for reproducing field failures for in-house debugging. The tool aims to synthesize in-house executions that mimic field failures. To do so, the authors use several types of data collected in the field such as stack traces, crash stacks, and points of failure. The data that successfully reproduced the field crash is sent to software developers to fix the bug.

Based on the success of BugRedux, the authors built F3 (Fault localization for Field Failures) [JO13]. F3 performs many executions of a program on top of BugRedux in order to cover different paths leading to the fault. It then generates many pass and fail paths which can lead to a better understanding of the bug. They also use grouping, profiling and filtering, to improve the fault localization process.

While being close to our approach, BugRedux and F3 may require the call sequence and/or the complete execution trace in order to achieve bug reproduction. When using only the crash traces (referred to as call stack at crash time in their paper), the success rate of BugRedux significantly drops to 37.5% (6/16). The call sequence and the complete execution trace required to reach 100% of bug reproduction can only be obtained through instrumentation and with an overhead ranging from 1% to 1066%.

Clause et al. [CO07] proposed a technique for enabling and supporting debugging of field failures. They record the execution of the program on the client side and propose to compress the generated data to the minimal required size to ensure that the reproduction is feasible. This compression is also performed on the client side. Moreover, the authors keep traces of all accessed documents in the operating system and also compress/reduce them to the minimal. Overall, they are able to reproduce on-field bug using a file weighting 70Kb. The minimal execution paths triggering the failure are then sent to the developers who can replay the execution on a sandbox, simulating the clients environment. While efficient, this approach suffers from severe security and privacy issues.

RECORE (REconstructing CORE dumps) is a tool proposed by Rossler et al. [RZF⁺13]. It instruments Java bytecode to wrap every method in a try and catch block while keeping a quasi-null overhead. The tool starts from the core dump and tries (with evolutionary

algorithms) to reproduce the same dump by executing the programs many times. The set of inputs responsible for the failure is generated when the generated dump matches the collected one. ReCrash [AKE08] is a tool that aims to make software failures reproducible by preserving object states. It uses an in-memory stack, which contains every argument and object clone of the real execution in order to reproduce a crash via the automatic generation of unit test cases. Unit test cases are used to provide hints to the developers on the buggy code. This approach suffers from overhead when they record everything (between 13% to 64% in some cases). The authors also propose an alternative in which they record only the methods surrounding the crash. For this to work, the crash has to occur at least once so they could use the information causing the crash to identify the methods surrounding it when (and if) it appears.

JRapture [SCFP00] is a capture/replay tool for observation-based testing. The tool captures execution of Java programs to replay it in-house. To capture the execution of a Java program, the authors used their own version of the Java Virtual Machine (JVM) and employ a lightweight, transparent capture process. Using their own JVM allows one to capture any interactions between a Java program and the system, including GUI, file, and console inputs, and on replay, it presents each thread with exactly the same input sequence it saw during capture. Unfortunately, they have to make their customer use their own JVM in order to support their approach, which limits the generalization of the approach to mass-market software.

Finally, Zamfir et al. [PO11] proposed ESD, an execution synthesis approach which automatically synthesizes failure execution using only the stack trace information. However, this stack trace is extracted from the core dump and may not always contain the components that caused the crash.

Except for STAR, approaches targeting the reproduction of field crashes require the instrumentation of the code or the running platform in order to save the stack call or the objects to successfully reproduce bugs. As we discussed earlier, instrumentation can cause a massive overhead (1% to 1066%) while running the system. In addition, data generated at run-time using instrumentation may contain sensitive information.

2.2 Issue and source code relationships

Researchers started studying the relationships between issues and source code repositories more than two decades ago. To the best of our knowledge the first ones who conduct this type of study on a significant scale were Perry and Stieg [PS93]. In these two decades, many aspects of these relationships have been studied in length. For example, researchers interested themselves in ameliorating the issues report by specifying guidelines to make a good report [BJS⁺08] and try to further simplify the existing models [HGGBR08].

Then, we can find approaches on how long it will take for an issue to get fixed [BN11, ZGV13, SKP14] and where it should be fixed [ZZL12, KTM⁺13]. With the rapidly increasing number of issues, the community also interested itself in prioritizing the issues report compared to one another [KWM⁺11] and do so by predicting the severity of an issue [LDGG10].

Finally, researchers proposed approaches to predict which issues will get reopened [ZNGM12, Lo13] which issues report is a duplicate of which other one [JW08, BPZ08, TSL12].

Another field of study consists in assigning these issues reports, if possible automatically to the right developers through triaging [AHM06, JKZ09, TNAKN11a, BvdH13] and predicting which locations are likely to yield new bugs [KZPJ06, KZWZ07a].

2.3 Crash Prediction

Predicting crash, fault and bug is very large and popular research area. The main goal behind the plethora of papers is to save on manpower—being the most expensive resource to build software—by directing their efforts on locations likely to contain a bug, fault or crash.

There are two distinct trends in crash, fault and bug prediction in the papers accepted to major venues such as MSR, ICSE, ICSME and ASE: history analysis and current version analysis.

In the history analysis, researchers extract and interpret information from the system. The idea being that the files or locations that are the most frequently changed are more likely to contain a bug. Additionally, some of these approaches also assume that locations linked to a previous bug are likely to be linked to a bug in the future.

On the other hand, approaches using only the current version to predict bugs assume that the current version, i.e. its design, call graph, quality metrics and more, will trigger the appearance of the bug in the future. Consequently, they do not require the history and only need the current source-code.

In the remaining of this section, we will describe approaches belonging to the two families.

Change logs approaches

Change logs based approaches rely on mining the historical data of the application and more particularly, the source code *diffs*. A source code *diffs* contains two versions of the same code in one file. Indeed, it contains the lines of code that have been deleted and the one that has been added. It is worth noting that, *diffs* files do not represent the concept of modified line. Indeed, a modified line will be represented by a deletion and an addition. Researchers mainly use five metrics when dealing with *diffs* files:

- Number of files: The number of modified files in a given commit
- Insertions: The number of added lines
- Deletions: The number of deleted lines
- Churns: The number of deleted lines immediately followed by an insertion which give an approximation of how many lines have been modified
- Hunks: The number of consecutive blocks of lines. This gives an approximation of how many distinct locations have been edited to accomplish a unit of work.

Naggapan *et al.* studied the churns metric and how it can be connected to the apparition of new defect in a complex software systems. They established that relative churns are, in fact, a better metric than classical churn [NB05a] while studying Windows Server 2003.

Hassan, interested himself with the entropy of change, i.e. how complex the change is [Has09]. Then, the complexity of the change, or entropy, can be used to predict bugs. The more complex a change is, the more likely it is to bring the defect with it. Hassan used its entropy metric, with success, on six different systems. Prior to this work, Hassan, in collaboration with Holt proposed an approach that highlights the top ten most susceptible locations to have a bug using heuristics based on *diffs* file metrics [HH05]. Moreover, their heuristics also leverage the data of the bug tracking system. Indeed, they use the past defect location to predict new ones. The conclusion of these two approaches has been that recently modified and fixed locations were the most defect-prone compared to frequently modified ones.

Similarly to Hassan and Hold, Ostrand *et al.* predict future crash location by combining the data from changed and past defect locations [OWB05]. The main difference between

Hassan and Hold and Ostrand *et al.* is that Ostrand *et al.* validate their approach on industrial systems as they are members of the AT&T lab while Hassan and Hold validated their approach on open-source systems. This proved that these metrics are relevant for open-source and industrial systems.

Kim *et al.* applied the same recipe and mined recent changes and defects with their approach named bug cache [KZWZ07b]. However, they are more accurate than the previous approaches at detecting defect location by taking into account that is more likely for a developer to make a change that introduces a defect when being under pressure. Such changes can be pushed to revision-control system when deadlines and releases date are approaching.

Single-version approaches

Approaches belonging to the single-version family will only consider the current version of the software at hand. Simply put, they don't leverage the history of changes or bug reports. Despite this fact, that one can see as a disadvantage compared to approaches that do leverage history, these approaches yield interesting results using code-based metrics.

Chidamber and Kemerer published the well-known CK metrics suite [CK94] for object oriented designs and inspired Moha *et al.* to publish similar metrics for service oriented programs [MPNC12]. Another famous metric suite for assessing the quality of a given software design is Briand's coupling metrics [BDW99].

The CK and Briand's metrics suites have been used, for example, by Basili *et al.* [BBM96], El Emam *et al.* [EMM01], Subramanyam *et al.* [SK03] and Gyimothy *et al.* [GFS05] for object oriented designs. Service oriented designs have been far less studied than object oriented design as they are relatively new, but, Nayrolles *et al.* [NMHIL, Nay13], Demange *et al.* [DMT13] and Palma *et al.* [Pal13] used Moha *et al.* metric suites to detect software defects.

All these approaches, proved software metrics to be useful at detecting software fault for object oriented and service oriented designs, respectively.

Finally, Nagappan *et al.* [NB05b, NBZ06] and Zimmerman [ZPZ07, ZN08] further refined metrics-based detection by using statical analysis and call-graph analysis.

While hundreds of bug prediction papers have been published by academia over the last decade, the developed tools and approaches fail to change developer behavior while deployed in industrial environment [LLS⁺13]. This is mainly due to the lack of actionable message, i.e. messages that provide concrete steps to resolve the problem at hand.

pErICOPE (Ecosystem Improve source COde during Programming session with real-time

mining of common knowlEdge), our proposed ecosystem, will be different in a sense that we will provide actionable messages as presented in Section 1.2.2.

2.4 Clone Detection

Clone detection is an important and difficult task. Throughout the years, researchers and practitioners have developed a considerable number of methods and tools in order to detect efficiently source code clones.

Text-based techniques use the code — often raw (e.g. with comments) — and compare sequences of code (blocks) to each other in order to identify potential clones. Johnson was perhaps the first one to use fingerprints to detect clones[Joh93, Joh94]. Blocks of code are hashed, producing fingerprints that can be compared. If two blocks share the same fingerprint, they are considered as clones. Manber et al. [Man94] and Ducasse et al.[DRD99] refined the fingerprint technique by using leading keywords and dot-plots.

Tree-matching and metric-based are two sub-categories of syntactic analysis for clone detection. Syntactic analysis consists of building abstract syntax trees (AST) and analyse them with a set of dedicated metrics or searching for identical sub-trees. Many approaches using AST have been published using sub-tree comparison including the work of Baxter et al.[BYM⁺98], Wahleret et al. [WSWF], or more recently, the work of Jian et al. with Deckard [JMSG07]. An AST-based approach compares metrics computed on the AST, rather than the code itself, to identify clones [PMDL99, BMD⁺].

Another approach to detect clones is to use static analysis and to leverage the semantics of the program to improve the detection. These techniques rely on program dependency graphs where nodes are statements and edges are dependencies. Then, the problem of finding clones is reduced to the problem of finding identical sub-groups in the program dependency graph. Examples of recent techniques that fall into this category are the ones presented by Krinke et al.[Kri01] and Gabel et al. [GJS08].

Many clone detection tools have been created using a lexical approach for clone detection. Here, the code is transformed into a series of tokens. If sub-series repeat themselves, it means that a potential clone is in the code. Some popular tools that use this technique include, but not limited to, Dup[Bak], CCFinder[KKI02b], and CP-Miner[LLMZ06].

Furthermore, a large number of taxonomies have been published in an attempt to classify clones and ease the research on clone detection[MLM96, BMD⁺99, KFF06, BKA⁺07, Kon, KG].

Other active research activities in clone detection focus on clone removal and management. Once detected, an obvious step is to provide approaches to remove clones in an

automatic way or (at least) keep track of them if removing them is not an option. Most modern IDEs provide the *extract method* feature that transforms a potentially copy-pasted block of code into a method and a call to the newly generated method[KH, HKKI04]. More advanced techniques (see Codelink[TBG] and[DER07]) involve analysing the output of CCFinder[KKI02a, LHMI07] or program dependencies graphs[HKKI04] to automatically suggest a method that would go through the *extract method* process.

The aforementioned techniques, however, focus on detecting clones after they are inserted in the code. Only a few studies focus on preventing the insertion of clones. Lague et al. [LPM⁺] conducted a very large empirical study with 10,000 developers over 3 years, where developers were asked to use clone detection tools during the development process of a very large telecom system. The authors found that while clones are being removed over time, using clone detection tools help improving the quality of the system as it prevents defects to reach the customers. Duala et al. [DER07, DER10] proposed to create clone region descriptors (CRDs), which describe clone regions within methods in a robust way that is independent from the exact text of the clone region or its location in a file. Then, using CRDs, clone insertion can be prevented.

PRECINCT aims to prevent clone insertion while integrating the clone detection process in a transparent manner in the day-to-day development process. This way, software developers do not have to resort to external tools to remove clones after they are inserted. Our approach operates at commit time, notifying software developers of possible clones as they commit their code.

Chapter 3

Aggregating versionning systems and project managements systems

In this chapter, we present two approaches: **BUMPER** and **JCHARMING**. **BUMPER** is the keystone of our proposed framework for programmatic software maintenance. The role of **BUMPER** is to aggregate information belonging to versionning systems and project management systems. **BUMPER** acts as our consolidated dataset.

JCHARMING is an approach to reproduce field-failure. Every time a bug report is submitted to a project management system and aggregated by **BUMPER**, **JCHARMING** will try to reproduce it. In case of success, the steps to reproduce the bug are saved in **BUMPER**.

The materials presented in this chapter are based on the following publications:

- Nayrolles, M. , Hamou-Lhadj, W., Tahar, S. Larsson, A. (2016). A Bug Reproduction Approach Based on Directed Model Checking and Crash Traces. *Journal of Software: Evolution and Process*. Wiley. (Accepted)
- Nayrolles, M. & Hamou-Lhadj, W. (2016). **BUMPER**: A Tool to Cope with Natural Language Search of Millions Bugs and Fixes. *International Conference on Software Analysis, Evolution, and Reengineering - Tool Track (SANER'16)*. IEEE.
- Mathieu Nayrolles, Abdelwahab Hamou-Lhadj, Tahar Sofiene, and Alf Larsson. **JCHARMING** : A Bug Reproduction Approach Using Crash Traces and Directed Model Checking. In *SANER'15*, pages 101-110, 2015. (Best Paper Award).

3.1 BUMPER - Bug Meta-repository For Developers & Researchers

With the goal to support the research towards analyzing relationships between bugs and their fixes we constructed a dataset of 380 projects, more than 100,000 resolved/fixed and with 60,000 changesets that were involved in fixing them from Netbeans and The Apache Software foundation’s software that is (1) searchable in natural language at research.mathieu-nayrolles.com/bumper, (2) contains clear relationships between the bug report and the code involved to fix it, (3) supports complex queries such as parent-child relationships, unions or disjunctions and (4) provide easy exports in json, csv and xml format.

In what follows, we will present the projects we selected. Then, we present the features related to the bugs and their fixes we integrate in BUMPER (BUg Metarepository for dEvelopers and Researchers) and how we construct our dataset. Then, we present the API, based on Apache Solr [Nay14], which allows the NLP search with practical examples before providing research opportunities based on our dataset.

However, to the best of our knowledge, no attempt has been made towards building a unified and online dataset where all the information related to a bug, or a fix can be easily accessed by researchers and engineers.

3.1.1 Data collection

Figure 6 illustrates our data collection and analysis process that we present here and discuss in more detail in the following subsections. First, we extract the raw data from the two bug report management systems used in this study (Bugzilla¹ and Jira²). The extracted data is consolidated in one database called BUMPER where we associate each bug report with its fix. The fixes are mined from different type of source versioning system. Indeed, Netbeans is based on mercurial³ while we used the git⁴ mirrors⁵ for the Apache Foundation software.

We used the same datasets presented earlier in Table 6. We choose to use the same ones as these two datasets because they exposed a great diversity in programming languages, teams, localization, utility and maturity. Moreover, the used different tools, i.e. Bugzilla, JIRA, Git and Mercurial, and therefore, BUMPER is ready to host any other datasets that used any composition of these tools.

¹<https://netbeans.org/bugzilla/>

²<https://issues.apache.org/jira/issues/?jql=>

³ <http://mercurial.selenic.com/>

⁴<http://git-scm.com/>

⁵<https://github.com/apache>

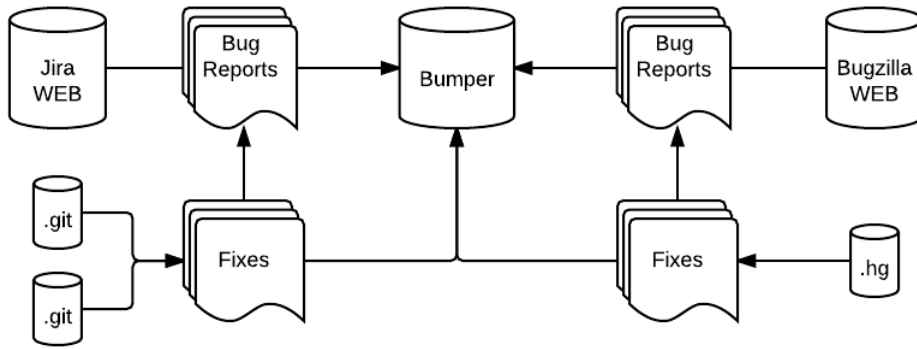


Figure 4: Overview of the bumper database construction.

3.1.2 Architecture

BUMPER rely on a highly scalable architecture composed of two distinct servers as depicted in Figure 5. The first server, on the left, handles the web requests and runs three distinct components:

- Pound is a lightweight open source reverse proxy program and application firewall. It is also served us to decode to request to http. Translating an request to http and then, use this HTTP request instead of the one allow us to save the http's decryption time required at each step. Pound also acts as a load-balancing service for the lower levels.
- Translated requests are then handled to Varnish. Varnish is an HTTP accelerator designed for content-heavy and dynamic websites. What it does is caching request that come in and serve the answer from the cache is the cache is still valid.
- NginX (pronounced engine-x) is a web-server that has been developed with a particular focus on high concurrency, high performances and low memory usage.

On the second server, that concretely handles our data, we have the following items:

- Pound. Once again, we use pound here, for the exact same reasons.
- SolrCloud is the scalable version of Apache Solr where the data can be separated into shards (e.g chunk of manageable size). Each shard can be hosted on a different server, but it's still indexed in a central repository. Hence, we can guarantee a low query time while exponentially increasing the data.

- Lucene is the full text search engine powering Solr. Each Solr server has its own embedded engine.

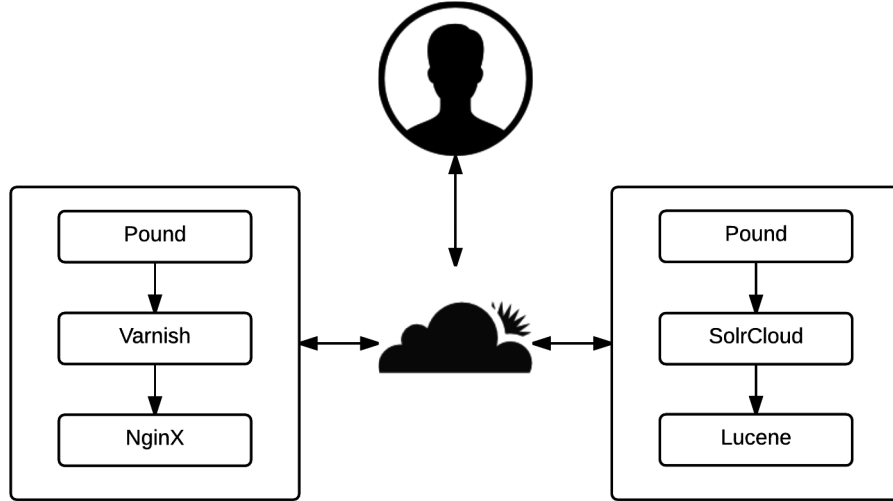


Figure 5: Overview of the bumper architecture.

Request from users to the servers and the communication between our servers are going through the CloudFlare network. CloudFlare acts as a content delivery network sitting between the users and the webserver. They also provide an extra level of caching and security.

To give the reader a glimpse about the performances that this unusual architecture can yield; we are able to request and display the result of a specific request in less than 100 ms while our two servers are, in fact, two virtual machines sharing an AMD Opteron (tm) Processor 6386 SE (1 core @ 2,000 MHz) and 1 GB of RAM.

3.1.3 UML Metamodel

Figure 6 presents the simplified BUMPER metamodel that we designed according to our bug taxonomy presented in section 27 and according to our future needs for JCHARMING, RESSEMBLE and BIANCA.

An *issue (task)* is characterized by a *date*, *title*, *description*, and a *fixing time*. They are reported (created) by and assigned to *users*. Also, *issues (tasks)* belong to *project* that are in *repository* and might be composed of *sub-projects*. *Users* can modify an *issue (task)* during *life cycle events* which impact the *type*, the *resolution*, the *platform*, the *OS* and the

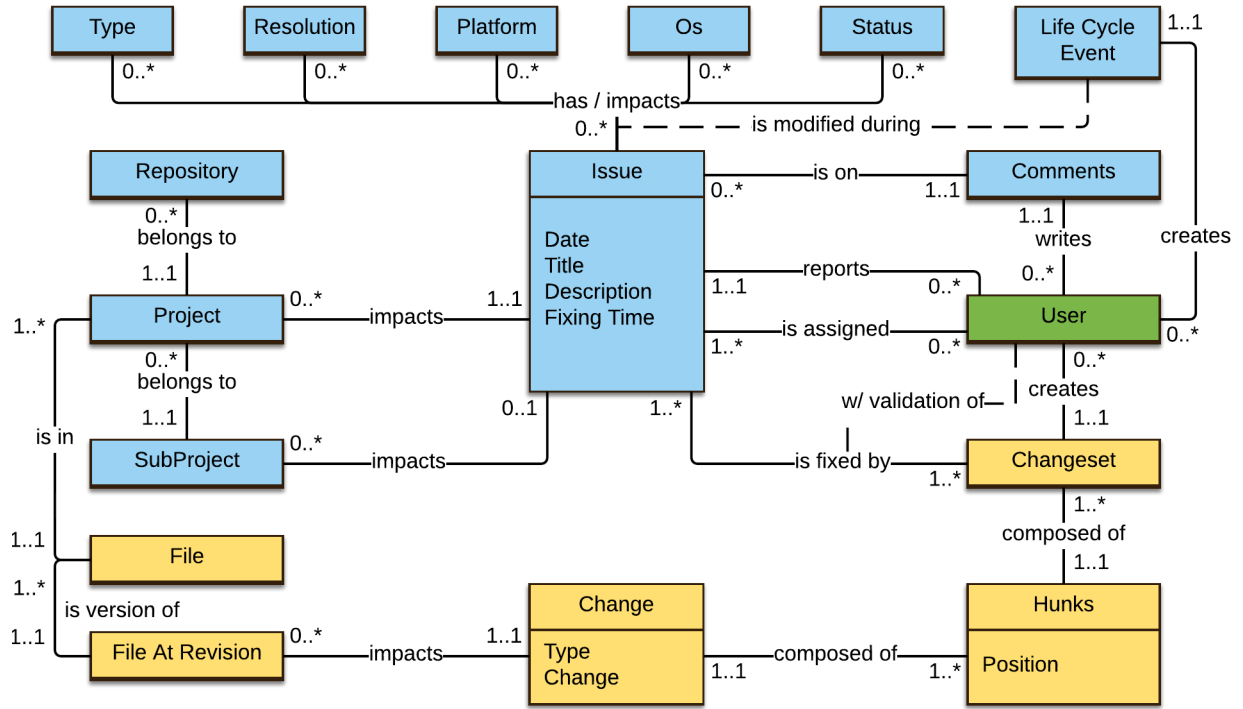


Figure 6: Overview of the bumper meta-model.

status. *Issues* (*tasks*) are resolved (implemented) by *changeset* that are composed of *hunks*. *Hunks* contain the actual changes to a *file* at a given revision, which are versions of the *file* entity that belongs to a *Project*.

3.1.4 Features

In this section, we present the features of bug report and their fixes in details.

Bug Report

A bug report is characterized by the following features:

- ID: unique string id of the form bug_dataset_project_bug_id
- Dataset: the dataset of which the bug is extracted from.
- Type: The type help us to distinguish different type of entities in BUMPER, i.e the bugs, changesets and hunks. For bug report, the type is always set to BUG
- Date: The date at which the bug report has been submitted.
- Title: The title of the bug report.

- Project: The project that this bug affects.
- Sub_project: The sub-project that this bug affects.
- Full_name_project: The combination of the project and the sub-project.
- Version: the version of the project that this bug affects
- Impacted_platform: the platform that this bug affects
- Impacted_os: the operating system that this bug affects
- Bug_status: The status of the bug. As in bumper, our main concern is on the relationship between of fix and a bug, we only have RESOLVED bugs
- Resolution: How the bug was resolved. Once again, as we are interested in investigating the fixes and the bugs, we only have FIXED bugs.
- Reporter_pseudo: the pseudonym of the person who report the bug.
- Reporter_name: the name of the person who reported the bug
- Assigned_to_pseudo: the pseudonym of the person who have
- been assigned to fix this bug
- Assigned_to_name: the name of the person who have been assigned to fix this bug
- Bug_severity: the severity of a bug
- Description: the description of the bug the reporter gave
- Fixing_time: The time it took to fix the bug, i.e the elapsed time between the creation of the BR and its modification to resolve/fixed, in minutes
- Comment_nb: How many comments have been posted on the bug report system for that bug
- Comment: Contains one comment. A bug can have 0 or many comments
- File: A file qualified name that has been modified in order to fix a bug. A bug can have 0 (in case we did not find its related commit) or many files.

We selected this set of features for bug report as they are the ones that are analyzed in many past and recent studies. In addition, bugs can contain 0 or many .

Changesets

In this section, we present the features that characterize changeset entities in BUMPER.

- ID: the SHA1 hash
- User: the name and email of the person who submitted that commit
- Date: the date at which this commit has been fixed
- Summary: the commit message entered by the user
- File: The fully qualified name of a file modified on that commits. A changeset can have 1 or many files.
- Number_files: How many files have been modified in that commit
- Insertions: the number of inserted lines
- Deletions: the number of deleted lines
- Churns: the number of modified lines
- Hunks: the number of sets of consecutive changed lines
- Parent_bug: the id of the bug this changeset belongs to.

In addition, changesets contain one or many hunks.

Hunks

A hunks are a set of consecutive lines changed in a file in order. A set of hunks form a fix that can be scattered across one or many files. Knowing how many hunks a fixed required and what are the changes in each of them is useful, as explained by [2] to understand how many places developers have to go to fix a bug.

Hunks are composed of:

- ID: unique id based on the files, the insertion and the SHA1 of the commits
- Parent_changeset: the SHA1 of the Changeset this hunk belongs to
- Parent_bug: the id of the bug this hunk belongs to.
- Negative_churns: how many lines have been removed in that hunk

- Positive_churns: how many lines have been added in that hunk
- Insertion: the position in a file at which this hunk takes place.
- Change: One line that have been added or removed. A Hunk can contain one or many changes.

3.1.5 Application Program Interface (API)

BUMPER is available for engineers and researchers at <https://bumper-app.com> and take the form of a regular search engine. Bumper supports (1) natural language query, (2) parent-child relationships, query, (3) disjunctions and union between complex queries and (4) a straight forward export of query results in XML, CSV or JSON format.

Browsing BUMPER, the basic query mode, perform the following operation:

$$\begin{aligned}
 (type : BUG \text{ AND } report_t : ("YOUR TERMS")) \text{ OR } (!parent \text{ which} = type "BUG") \\
 fix_t : "YOUR TERMS"
 \end{aligned}
 \tag{1}$$

The first part of the query component of the query retrieves all the bugs that contains the “YOUR TERMS” query in at least one its features by selecting type: BUG and report_t, which is an index composed of all the features of the bug, set to “YOUR TERMS”. Then, we merge this query with another one that reads

(!parent which = type “BUG”)fix_t : “YOUR TERMS”). In this one, we retrieve the parent documents, i.e the bugs, of fixes that contains “YOUR TERMS” in their fix_t index. The fix_t index is, as for the BUG, an index based on all the fields of changeset and hunk both. As a result, we search seamlessly in the bug report and their fixes in natural language.

As a more practical example, Figure 7 illustrate a query on <https://bumper-app.com>. The search term is “Exception” and we can see that 20,285 issues / tasks have been found in 25 ms This particular set of issues, displayed on the left side, match because they contain “Exception” in the issue report or in the source code modified to fix this issue (implement this task). Then on the right side of the screen, the selected issue (task) is displayed. We can see the basic characteristic of the issue (task) followed by comments and finally, the source code.

Moreover, BUMPER supports AND, OR, NOR operators and provide results in order of seconds.

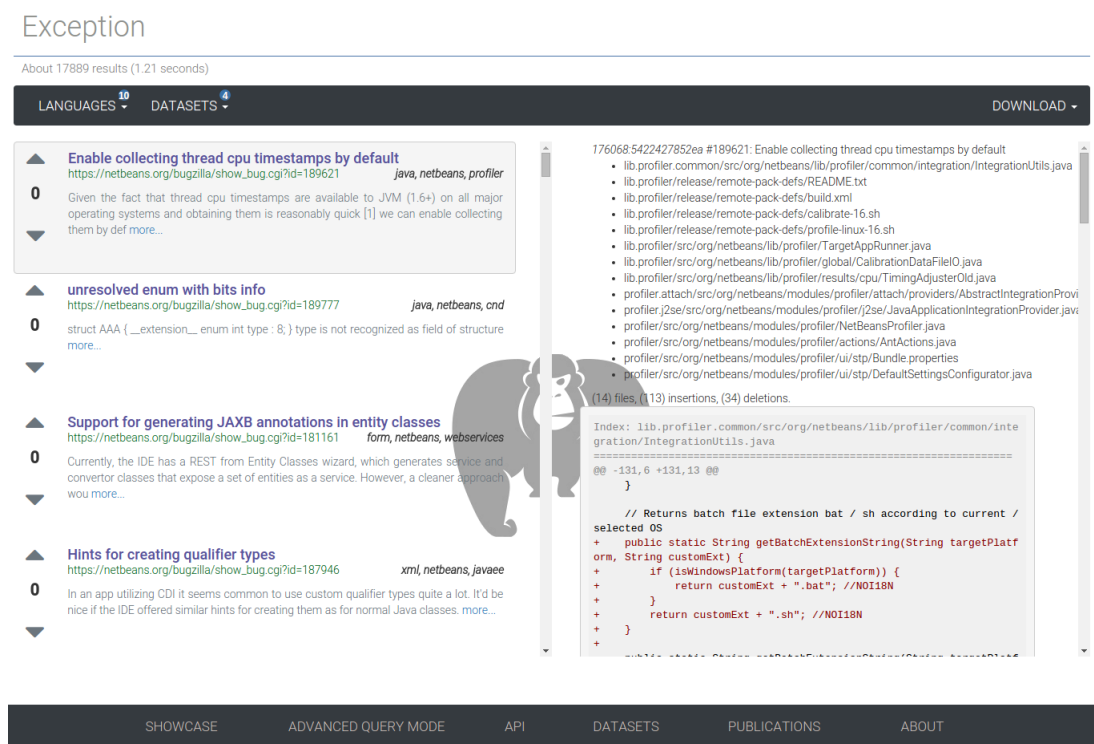


Figure 7: Screenshot of <https://bumper-app.com> with “Exception” as research.

As we said before, BUMPER is based on Apache Solr which have an incredibly rich API that is available online⁶.

BUMPER serves as data repositories for the upcoming approaches presented in the next sections.

3.2 JCHARMING - Java Crash Automatic Reproduction by directed Model checkING

Field failures are challenging to reproduce because the data provided by the end users is often scarce. A survey conducted with developers of major open source software systems such as Apache, Mozilla and Eclipse revealed that one of the most valuable piece of information that can help locate and fix the cause of a crash is the one that can help reproduce it [BJS⁺08]. It is therefore important to invest in techniques and tools for automatic bug reproduction to ease the maintenance process and accelerate the rate of bug fixes and patches.

⁶ <http://lucene.apache.org/solr/resources.html>

In this section, we present an approach, called **JCHARMING** (Java CrasH Automatic Reproduction by directed Model checkING) that uses a combination of crash traces and model checking to automatically reproduce bugs that caused field failures. Unlike existing techniques, such as on-field record and in-house replay [NPC05, AKE08, JKXC10] or crash explanation [MSA04, CFS09] **JCHARMING** does not require instrumentation of the code. It does not need access to the content of the heap either. Instead, **JCHARMING** uses a list of functions output when an uncaught exception in Java occurs (i.e., the crash trace) to guide a model checking engine to uncover the statements that caused the crash. While we do not filter any personal information that may appear in the crash trace, **JCHARMING** raises less privacy concerns than a tool recording every call or dump the content of the memory.

To assess the efficiency of **JCHARMING** we try to reproduce bug reports contained in **BUMPER**.

3.2.1 Preliminaries

Model checking (also known as property checking) will, given a system (that could be software [VHB⁺03] or hardware based [Kro99]), check if the system meets a specification Spec by testing exhaustively all the states of the system under test (SUT), which can be represented by a Kripke [Kri63] structure:

$$SUT = \langle S, T, P \rangle \quad (2)$$

where S is the set of states, $T \subseteq S * S$ represents the transitions between the states and P is the set of properties that each state satisfies. The SUT is said to satisfy a set of properties p when there exists a sequence of states transition x leading towards these properties. This can be written as:

$$(SUT, x) \models p \quad (3)$$

However, this only ensures that $\exists x$ such that p is reached at some point in the execution of the program and not that p holds nor that $\forall x, p$ is satisfiable. In **JCHARMING**, SUTs are bound to a simple specification: they must not crash under a fair environment. In the framework of this study, we consider a fair environment as any environment where the transitions between the states represent the functionalities offered by the program. For example, in a fair environment, the program heap or other memory spaces cannot be modified. Without this fairness constraint, all programs could be tagged as buggy since we could, for example, destroy objects in memory while the program continues its execution. As we are interested in verifying the absence of unhandled exceptions in the SUT, we aim

to verify that for all possible combinations of states and transitions there is no path leading towards a crash. That is:

$$\forall x.(SUT, x) \models \neg c \quad (4)$$

If such a path exists (i.e., $\exists x$ such that $(SUT, x) \models c$) then the model checker engine will output the path x (known as the counter-example) which can then be executed. The resulting Java exception crash trace is compared with the original crash trace to assess if the bug is reproduced. While being accurate and exhaustive in finding counter-examples, model checking suffers from the state explosion problem, which hinders its applicability to large software systems.

To show the contrast between testing and model checking, we use the hypothetical example of Figures 8, 9 and 10 to sketch the possible results of each approach. These figures depicts a toy program where from the entry point, unknown calls are made (dotted points) and, at some points, two methods are called. These methods, called `Foo.Bar` and `Bar.Foo`, implement a for loop from 0 to `loopCount`. The only difference between these two methods is that the `Bar.Foo` method throws an exception if `i` becomes larger than two. Hereafter, we denote this property as $c_{i>2}$.

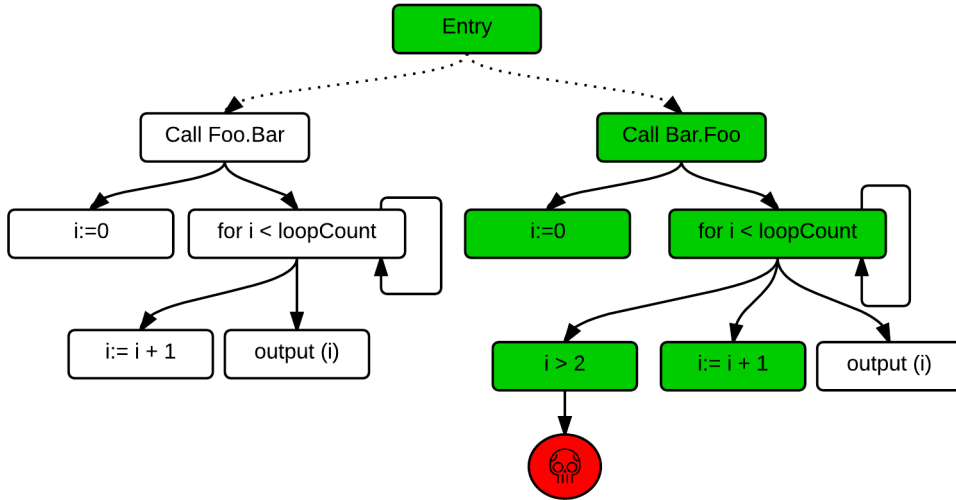


Figure 8: A toy program under testing

Figure 8 shows the program statements that could be covered using testing approaches. Testing software is a demanding task where a set of techniques is used to test the SUT according to some input.

Software testing depends on how well the tester understands the SUT in order to write

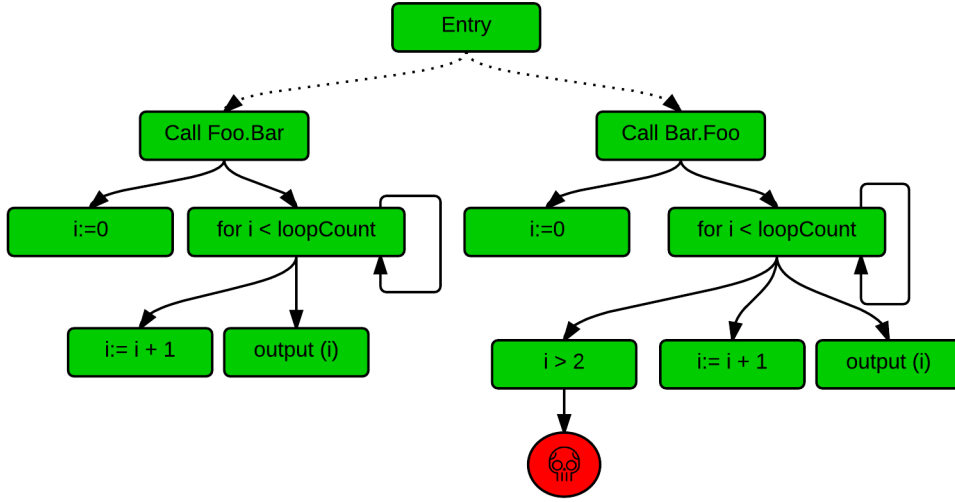


Figure 9: A toy program under model checking

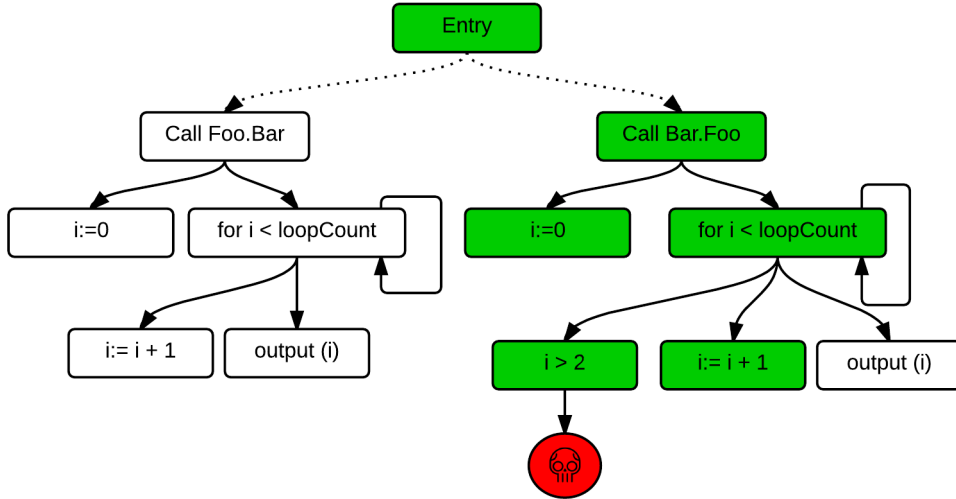


Figure 10: A toy program under directed model checking

relevant test cases that are likely to find errors in the program. Program testing is usually insufficient because it is not exhaustive. In our case, using testing would mean that the tester knows what to look for in order to detect the causes of the failure. We do not assume this knowledge in JCHARMING.

Model checking, on the other hand, explores each and every state of the program (Figure 9), which makes it complete, but impractical for real-world and large systems. To overcome the state explosion problem of model checking, directed (or guided) model checking has

been introduced [RM09]. Directed model checking use insights generally heuristics about the SUT in order to reduce the number of states that need to be examined. Figure 10 explores only the states that may lead to a specific location, in our case, the location of the fault. The challenge, however, is to design techniques that can guide the model checking engine. As we will describe in the next section, we use crash traces and program slicing to overcome this challenge.

3.2.2 The JCHARMING Approach

Figure 11 shows an overview of JCHARMING. The first step consists of collecting crash traces, which contain raw lines displayed to the standard output when an uncaught exception in Java occurs. In the second step, the crash traces are preprocessed by removing noise (mainly calls to Java standard library methods). The next step is to apply backward slicing using static analysis to expand the information contained in the crash trace while reducing the search space. The resulting slice along with the crash trace are given as input to the model checking engine. The model checker executes statements along the paths from the main function to the first line of the crash trace (i.e., the last method executed at crash time, also called the crash location point). Once the model checker finds inconsistencies in the program leading to a crash, we take the crash stack generated by the model checker and compare it to the original crash trace (after preprocessing). The last step is to build a JUnit test, to be used by software engineers to reproduce the bug in a deterministic way.

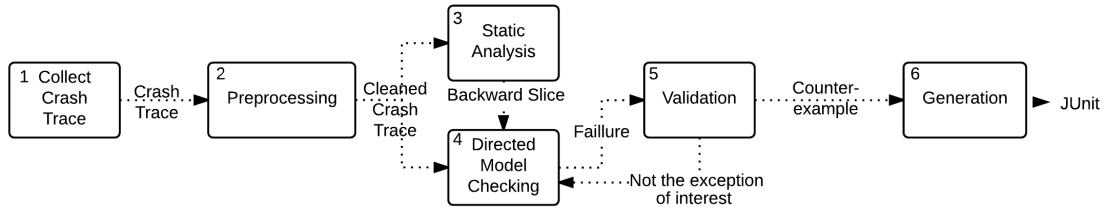


Figure 11: Overview of JCHARMING.

Collecting Crash Traces

The first step of JCHARMING is to collect the crash trace caused by an uncaught exception. Crash traces are usually included in crash reports and can therefore be automatically retrieved using a simple regular expression. Figure 12 shows an example of a crash trace that contains the exception thrown when executing the program depicted in Figures 8 to 10. The crash trace contains a call to the `Bar.foo()` method the crash location point and calls

to Java standard library functions (in this case, GUI methods because the program was launched using a GUI).

```

1.java.awt.AWTErrorReporter:loopTimes
should be <= 3
2. at Foo.bar(Foo.java:10)
3. at GUI.buttonActionPerformed(GUI.java:88)
4. at GUI.access0(GUI.java : 85)
5.atGUI1.actionPerformed(GUI.java:57)
6. caused by java.lang.IndexOutOfBoundsException : 3
7. at scam.Foo.buggy(Foo.java:17)
8. and 4 more ...

```

Figure 12: Java `InvalidActivityException` is thrown in the `Bar.Goo` loop if the control variable is greater than 2.

As shown in Figure 12, we can see that the first line (referred to as frame f_0 , subsequently the next line is called frame f_1 , etc.) does not represent the real crash point but it is only the last exception of a chain of exceptions. Indeed, the `InvalidActivity` has been triggered by an `IndexOutOfBoundsException` in `scam.Foo.buggy`. This kind of crash traces reflects several nested try/catch blocks.

In addition, it is common in Java to have incomplete crash traces. According to the Java documentation [Ora11], line 8 of Figure 12 should be interpreted as follows: *“This line indicates that the remainder of the stack trace for this exception matches the indicated number of frames from the bottom of the stack trace of the exception that was caused by this exception (the “enclosing exception”). This shorthand can greatly reduce the length of the output in the common case where a wrapped exception is thrown from the same method as the “causative exception” is caught.”*

We are likely to find shortened traces in bug repositories as they are what the user sees without any possibility to expand their content.

Preprocessing

In the preprocessing step, we first reconstruct and reorganize the crash trace in order to address the problem of nested exceptions. Then, with the aim to obtain an optimal guidance for our directed model checking engine, we remove frames that are out of our control. Frames out of our controls refer usually, but are not limited to, Java library methods and

third party libraries. In Figure 12, we can see that Java GUI and event management components appear in the crash trace. We assume that these methods are not the cause of the crash; otherwise it means that there is something wrong with the on-field JDK. If this is the case, we will not be able to reproduce the crash. Note that removing these unneeded frames will also reduce the search space of the model checker.

Building the Backward Static Slice

For large systems, a crash trace does not necessarily contain all the methods that have been executed starting from the entry point of the program (i.e., the main function) to the crash location point. We need to complete the content of the crash trace by identifying all the statements that have been executed starting from the main function until the last line of the preprocessed crash trace. In Figure 12, this will be the function call `Bar.foo()`, which happens to be also the crash location point. To achieve this, we turn to static analysis by extracting a backward slice from the main function of the program to the `Bar.foo()` method.

A backward slice contains all possible branches that may lead to a point n from a point m as well as the definition of the variables that control these branches [De 01]. In other words, the slice of a program point n is the program subset that may influence the reachability of point n starting from point m . The backward slice containing the branches and the definition of the variables leading to n from m is noted as $bslice_{[m \leftarrow n]}$.

We perform a static backward slice between each frame to compensate for possible missing information in the crash trace. More formally, the final static backward slice is represented as follows:

$$bslice_{[entry \leftarrow f_0]} = bslice_{[f_1 \leftarrow f_0]} \cup bslice_{[f_2 \leftarrow f_1]} \cup \dots \cup bslice_{[f_n \leftarrow f_{n-1}]} \cup bslice_{[entry \leftarrow f_n]} \quad (5)$$

Note that the union of the slices computed between each pair of frames must be a subset of the final slice between f_0 and the entry point of the program. More formally:

$$\bigcup_{i=0}^{entry} bslice_{[f_{i+1} \leftarrow f_i]} \subseteq bslice_{[entry \leftarrow f_0]} \quad (6)$$

Indeed, in Figure 13, the set of states allowing to reach f_0 from f_2 is greater than the set of states to reach f_1 from f_2 plus set of states to reach f_0 from f_1 . In this hypothetical example and assuming that z_2 is a prerequisite to f_2 then $bslice_{[entry \leftarrow f_0]} = \{f_0, f_1, f_2, z_0, z_1, z_2, z_3\}$ while $\bigcup_{i=0}^n bslice_{[f_{i+1} \leftarrow f_i]}$.

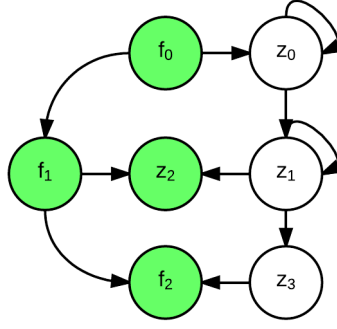


Figure 13: Hypothetical example representing $bslice_{[entry \leftarrow f_0]}$ Vs. $\cup_{i=0}^n bslice_{[f_{i+1} \leftarrow f_i]} = \{f_0, f_1, f_2, z_2\}$

In the worst case scenerio where there exists one and only one transition between each frame, which is very unlikely for real and complex systems, then $bslice_{[entry \leftarrow f_0]}$ and $\cup_{i=0}^n bslice_{[f_{i+1} \leftarrow f_i]}$ yield the same set of states with a comparable computational cost since the number of branches to explore will be the same in both cases.

Algorithm 1 is a high level representation of how we compute the backward slice between each frame. The algorithm takes as input the pre-processed call trace, the byte code of the SUT, and the entry point. From line 1 to line 5, we initialize the different variables used by the algorithm. The main loop of the algorithm begins at line 6 and ends at line 15. In this loop, we compute the static slice between the current frame and the next one. If the computed static slice is not empty then we update the final backward slice with the newly

computed slice.

Data: Crash Stack, BCode, Entry Point

Result: BSolve

Frames *frames* \leftarrow *extract frames from crash stack*;

Int *n* \leftarrow size of frame;

Int *offset* \leftarrow 1;

Bslice *bSlice* $\leftarrow \emptyset$;

for *i* \leftarrow 0 to *i* < *n* **do** *offset* < *n* - 1

BSlice *currentBSlice* \leftarrow backward slice from frames[*i*] to *i* + *offset*;

if *currentBSlice* $\neq \emptyset$ **then**

bSlice \leftarrow *bSlice* \cup *currentBSlice*;

offset \leftarrow 1;

else

offset \leftarrow *offset* + 1;

end

end

Algorithm 1: High level algorithm computing the union of the slices

Using backward slicing, the search space of the model checker that processes the example of Figures 8 to 10 is given by the following expression:

$$\exists x. \left(\bigcup_{i=0}^{entry} bslice_{[f_{i+1} \leftarrow f_i]} \subset SUT \right) \models c_{i>2} \quad (7)$$

That is, there exists a sequence of states transitions *x* that satisfies $c_{i>2}$ where both the transitions and the states are entry elements of $\bigcup_{i=0}^{entry} bslice_{[f_{i+1} \leftarrow f_i]}$. Obviously, $c_{i>2}$ also needs to be included for the final static slice to be usable by the model checking engine. Consequently, the only frame that need to be untouched for the backward static slice to be meaningful is f_0 .

Directed Model Checking

The model checking engine we use in this paper is called JPF (Java PathFinder) [VPK04], which is an extensible JVM for Java bytecode verification. This tool was first created as a front-end for the SPIN model checker [Hol97] in 1999 before being open-sourced in 2005. JPF is organized around five simple operations: (i) *generate states*, (ii) *forward*, (iii) *backtrack*, (iv) *restore state* and (v) *check*. In the forward operation, the model checking engine generates the next state s_{t+1} . If s_{t+1} has successors then it is saved in a backtrack table to be restored later. The backtrack operation consists of restoring the last state in the

backtrack table. The restore operation allows restoring any state and can be used to restore the entire program as it was the last time we choose between two branches. After each, forward, backtrack and restore state operation the check properties operation is triggered.

In order to direct JPF, we have to modify the *generate states* and the *forward* steps. The *generate states* is populated with entry the states in $\bigcup_{i=0}^{entry} bslice_{[f_{i+1} \leftarrow f_i]} \subset SUT$ and we adjust the *forward step* to explore a state if the target state $s_i + 1$ and the transition x to pass from the current state s_i to s_{i+1} are in $\bigcup_{i=0}^{entry} bslice_{[f_{i+1} \leftarrow f_i]} \subset SUT$ and $x. \bigcup_{i=0}^{entry} bslice_{[f_{i+1} \leftarrow f_i]} \subset x.SUT$.

Validation

To validate the result of directed model checking, we modify the *check properties* step that checks if the current sequence of states transitions x satisfies a set a property. If the current states transitions x can throw an exception, we execute x and compare the exception thrown to the original crash trace (after preprocessing). If the two exceptions match, we conclude that the conditions needed to trigger the failure have been met and the bug is reproduced.

However, as argued by Kim et al. in [KZN13], the same failure can be reached from different paths of the program. Although the states executed to reach the defect are not exactly the same, they might be useful to enhance the understanding of the bug by software developers, and speed up the deployment of a fix. Therefore, in this paper, we consider a defect to be partially reproduced if the crash trace generated from the model checker matches the original crash trace by a factor of t , where t is a threshold specified by the user. t is the percentage of identical frames between both crash traces.

Generating Test Cases for Bug Reproduction

To help software developers reproduce the crash in a lab environment we automatically produce the JUnit test cases necessary to run the SUT to cause the exercise of the bug.

To build a test suite that reproduces a defect, we need to create a set of objects used as arguments for the methods that will enable us to travel from the entry point of the program to the defect location. JPF has the ability to keep track of what happens during model checking in the form of traces containing the visited states and the value of the variables. We leverage this capability to create the required objects and call the methods leading to the failure location. Although we can track back the internal state of objects at a specific time using JPF, it can be too computationally taxing to recreate only the objects needed to generate the bug. To overcome this, we use serialization techniques [OP99]. We take advantage of features offered by the XStream [Xst11] library which enables the

serialization and deserialization of any Java object even objects that do not implement the Java Serializable interface. We use the serialization when the model checker engine performs too many operations modifying the property of a given object. In such case, we serialize the last state of the object.

3.2.3 Case studies

In this section, we show the effectiveness of JCHARMING to reproduce bugs in seven open source systems⁷. The aim of the case study is to answer the following question: *Can we use crash traces and directed model checking to reproduce on- field bugs in a reasonable amount of time?*

Targeted Systems

Table 2 shows the systems and their characteristics in terms of Kilo Line of Code (KLoC) and Number of Classes (NoC).

SUT	KLOC	NoC	Bug #ID
Ant	265	1233	38622, 41422
ArgoUML	58	1922	2603, 2558, 311, 1786
dnsjava	33	182	38
jfreechart	310	990	434, 664, 916
Log4j	70	363	11570, 40212, 41186, 45335, 46271, 47912, 47957
MCT	203	1267	440ed48
pdfbox	201	957	1412, 1359

Table 2: List of target systems in terms of Kilo line of code (KLoC), number of classes (NoC) and Bug # ID

Apache Ant [Apa] is a popular command-line tool to build make files. While it is mainly known for Java applications, Apache Ant also allows building C and C++ applications. We choose to analyze Apache Ant because it has been used by other researchers in similar studies.

ArgoUML [Col] is one of the major players in the open source UML modeling tools. It has many years of bug management and, similar to Apache Ant, it has been extensively used as a test subject in many studies.

⁷The bug reports used in this study and the result of the model checker are made available for download from research.mathieu-nayrolles.com/jcharming/

Dnsjava [Wel13] is a tool for the implementation of the DNS mechanisms in Java. This tool can be used for queries, zone transfers, and dynamic updates. It is not as large as the other two, but it still makes an interesting case subject because it has been well maintained for the past decade. Also, this tool is used in many other popular tools such as Aspirin, Muffin and Scarab.

JfreeChart [Obj05] is a well-known library that enables the creation of professional charts. Similar to dnsjava, it has been maintained over a very long period of time JfreeChart was created in 2005 and it is a relatively large application.

Apache Log4j [The99] is a logging library for Java. This is not a very large library, but it is extensively used by thousands of programs. As other Apache projects, this tool is well maintained by a strong open source community and allows developers to submit bugs. The bugs which are in the bug report system of Log4j are, generally speaking, well documented and almost every bug contains a related crash trace and, therefore, it is a tool of interest to us.

MCT [NAS09] stands for Mission Control technologies and was developed by the NASA Ames Research Center (the creators of JPF) for use in spaceflight mission operation. This tool benefits from two years of history and targets a very critical domain, Spacial Mission Control. Therefore, this tool has to be particularly and carefully tested and, consequently, the remaining bugs should be hard to discover and reproduce.

PDFBox [Apa14] is another tool supported by the Apache Software Foundation since 2009 and was created in 2008. PDFBox allows the creation of new PDF documents and the manipulation of existing documents.

Bug Selection and Crash Traces

In this study, we have selected the reproduced bugs randomly in order to avoid the introduction of any bias. We selected a random number of bugs ranging from 1 to 10 for each SUT containing the word “exception” and where the description of the bug contains a match a regular expression designed to find the pattern of a Java exception.

3.2.4 Results

Table 3 shows the results of JCHARMING in terms of Bug #ID, reproduction status, and execution time (in minutes) of directed model checking (DMC) and Model Checking (MC). The experiments have been conducted on a Linux machine (8 GB of RAM and using Java 1.7.0.51).

- The result is noted as “Yes” if the bug has been fully reproduced, meaning that the

crash trace generated by the model checker is identical to the crash trace collected during the failure of the system.

- The result is “Partial” if the similarity between the crash trace generated by the model checker and the original crash trace is above $t=80\%$. Given an 80% similarity threshold, we consider partial reproduction as successful. A different threshold could be used.
- Finally, the result of the approach is reported as “No” if either the similarity is below $t \geq 80\%$ or the model checker failed to crash the system given the input we provided.

SUT	Bug #ID	Reprod.	Time DMC	Time MC
Ant	38622	Yes	25.4	-
	41422	No	-	-
ArgoUML	2558	Partial	10.6	-
	2603	Partial	9.4	-
	311	Yes	11.3	-
	1786	Partial	9.9	-
DnsJava	38	Yes	4	23
jFreeChart	434	Yes	27.3	-
	664	Partial	31.2	-
	916	Yes	26.4	-
Log4j	11570	Yes	12.1	-
	40212	Yes	15.8	-
	41186	Partial	16.7	-
	45335	No	-	-
	46271	Yes	13.9	-
	47912	Yes	12.3	-
	47957	No	-	-
MCT	440ed48	Yes	18.6	-
PDFBox	1412	Partial	19.7	-
	1359	No	-	-

Table 3: Effectiveness of JCHARMING using directed model checking (DMC) and model checking (MC) in minutes

As we can see in Table 3, we were able to reproduce 17 bugs out of 20 bugs either completely or partially (85ratio). The average time to reproduce a bug is 16 minutes. This

result demonstrates the effectiveness of our approach, more particularly, the use of backward slicing to create a manageable search space that guides adequately the model checking engine. We also believe that our approach is usable in practice since it is also time efficient. Among the 20 different bugs we have tested, we will describe one bug (chosen randomly) for each category (successfully reproduced, partially reproduced, and not reproduced) for further analysis.

Successfully reproduced

The first bug we describe in this discussion is the bug #311 belonging to ArgoUML. This bug was submitted in an earlier version of ArgoUML. This bug is very simple to manually reproduce thanks to the extensive description provided by the reporter, which reads: *“I open my first project (Untitled Model by default). I choose to draw a Class Diagram. I add a class to the diagram. The class name appears in the left browser panel. I can select the class by clicking on its name. I add an instance variable to the class. The attribute name appears in the left browser panel. I can’t select the attribute by clicking on its name. Exception occurred during event dispatching:”*

The reporter also attached the following crash trace that we used as input for JCHARMING:

```
1. java.lang.NullPointerException:
2. at
3. uci.uml.ui.props.PropPanelAttribute
   .setTargetInternal (PropPanelAttribute.java)
4. at uci.uml.ui.props.PropPanel.
   setTarget (PropPanel.java)
5. at uci.uml.ui.TabProps.setTarget (TabProps.java)
6. at uci.uml.ui.DetailsPane.setTarget
   (DetailsPane.java)
7. at uci.uml.ui.ProjectBrowser.select
   (ProjectBrowser.java)
8. at uci.uml.ui.NavigatorPane.mySingleClick
   (NavigatorPane.java)
9. at uci.uml.ui.NavigatorPane\$_Navigator
   MouseListener.mouse Clicked (NavigatorPane.java)
10. at
```

```

java.awt.AWTEventMulticaster.mouseClicked
(AWTEventMulticaster.java:211)
11. at java.awt.AWTEventMulticaster.mouseClicked
(AWTEventMulticaster.java:210)
12. at java.awt.Component.processMouseEvent
(Component.java:3168)
...
19. java.awt.LightweightDispatcher
.retargetMouseEvent (Container.java:2068)
22. at java.awt.Container
.dispatchEventImpl (Container.java:1046)
23. at java.awt.Window
.dispatchEventImpl (Window.java:749)
24. at java.awt.Component
.dispatchEvent (Component.java:2312)
25. at java.awt.EventQueue
.dispatchEvent (EventQueue.java:301)
28. at java.awt.EventDispatchThread.pumpEvents
(EventDispatchThread.java:90)
29. at java.awt.EventDispatchThread.run (EventDispatch
Thread.java:82)

```

The cause of this bug is that the reference to the attribute of the class was lost after being displayed on the left panel of ArgoUML and therefore, selecting it through a mouse click throws a null pointer exception. In the subsequent version, ArgoUML developers added a TargetManager to keep the reference of such object in the program. Using the crash trace, JCHARMING's preprocessing step removed the lines between lines 11 and 29 because they belong to the Java standard library and we do not want neither the static slice nor the model checking engine to verify the Java standard library but only the SUT. Then, the third step performs the static analysis following the process described in Section IV.C. The fourth step performs the model checking on the static slice to produce the same crash trace. More specifically, the model checker identifies that the method `setTargetInternal(Object o)` could receive a null object that will result in a Null pointer exception.

Partially reproduced

As an example of a partially reproduced bug, we explore the bug #664 of the Jfreechart program. The description provided by the reporter is: “*In ChartPanel.mouseMoved there’s a line of code which creates a new ChartMouseEvent using as first parameter the object returned by getChart(). For getChart() is legal to return null if the chart is null, but ChartMouseEvent’s constructor calls the parent constructor which throws an IllegalArgumentException if the object passed in is null.*”

The reporter provided the crash trace containing 42 lines and the replaced an unknown number of lines by the following statement “`␣deleted entry␣`”. While JCHARMING successfully reproduced a crash yielding almost the same trace as the original trace, the “`␣deleted entry␣`” statement – which was surrounded by calls to the standard java library – was not suppressed and stayed in the crash trace. That is, JCHARMING produced only the 6 (out of 7) first lines and reached 83% similarity, and thus a partial reproduction.

```
1. java.lang.IllegalArgumentException: null source
2. at java.util.EventObject.<init>(
EventObject.java:38)
3. at
4 org.jfree.chart.ChartMouseEvent.<init>
(ChartMouseEvent.java:83)
5. at org.jfree.chart.ChartPanel
.mouseMoved(ChartPanel.java:1692)
6. ␣<␣deleted entry␣>␣
```

In all bugs that were partially reproduced, we found that the differences between the crash trace generated from the model checker and the original crash trace (after preprocessing) consists of few lines only.

Not Reproduced

To conclude the discussion on the case study, we present a case where JCHARMING was unable to reproduce the failure. For the bug #47957 belonging to Log4j and reported in late 2009 the reporter wrote: “*Configure SyslogAppender with a Layout class that does not exist; it throws a NullPointerException. Following is the exception trace:*” and attached the following crash trace:

```

1. 10052009 01:36:46 ERROR [Default: 1]
   struts.CPEExceptionHandler.execute
   RID[(null;25KbxlK0voima4h00ZLBQFC;236A18E60000045C3A
   7D74272C4B4A61)]
2. Wrapping Exception in ModuleException
3. java.lang.NullPointerException
4. at org.apache.log4j.net.SyslogAppender
   .append(SyslogAppender.java:250)
5. at org.apache.log4j.AppenderSkeleton
   .doAppend(AppenderSkeleton.java:230)
6. at org.apache.log4j.helper.AppenderAttachableImpl
   .appendLoopOnAppenders(AppenderAttachableImpl
   .java:65)
7. at org.apache.log4j.Category.callAppenders
   (Category.java:203)
8. at org.apache.log4j.Category
   .forcedLog(Category.java:388)
9. at org.apache.log4j.Category.info
   (Category.java:663)

```

The first three lines are not produced by the standard execution of the SUT but by an `ExceptionHandler` belonging to Struts [Apa00]. Struts is an open source MVC (Model View Controller) framework for building Java Web Application. JCHARMING examined the source code of Log4J for the crash location `struts.CPEExceptionHandler.execute` and did not find it since this method belongs to the source base of Struts – which uses log4j as a logging mechanism. As a result, the backward slice was not produced, and we failed to perform the next steps. It is noteworthy that the bug is marked as duplicate of the bug #46271 which contains a proper crash trace. We believe that JCHARMING could have successfully reproduced the crash, if it was applied to the original bug.

While JCHARMING is effective at reproducing on-field failures in lab environment, we want to reduce their number in the coming year. To do so, we built **RESSEMBLE** and **BIANCA** that we present in the next two sections.

Chapter 4

Using clone detection for pragmatic software maintenance

The adoption of tools and processes that aim to support human maintainers during maintenance is relatively low [LLS⁺13, FM15, LWA07, APM⁺07, AP08, JSMHB13, Nor13, LvdH11]. Human maintainers agree that such tools are beneficial but, in addition of disrupting their workflow, these tools tend to have a high false positives rate. Also, the way in which the warnings are presented, among other things, are barriers to use[JSMHB13]. When asked for their opinions, human maintainers agree to the following characteristics for maintenance-oriented tools[HP04, LvdH11, LLS⁺13]:

- Actionable messages. Presenting a warning about bug-proneness of a given line is not enough. Clear actions to improve the source code should be provided
- Obvious reasoning. The conditions that led to a given warning should be understandable by the maintainer. If the conditions are hidden in complex statistical models, then maintainers cannot review them and will find it difficult to trust the tool.
- Scaling. Industrial sized project contains thousand files and dependencies which can be updated many times a day. Maintenance-oriented tools should not hinder the productivity of maintainers.
- Contextualization. Warnings and messages should always be contextualized with respect to the project at hand and not generic rules.
- Integration. Developers and maintainers are overwhelm by the amount of existing tools. Yet, their daily use three different kind of tools. An integrated development environment, a versionning system and a project tracking system to produce, version

and manage their software, respectively. Maintenance-oriented tools should fit in the existing ecosystem rather than complexifying the deployment process.

In modern versioning systems, when maintainers make modifications to the source code and their modification versionned, they have to do commit. The commit operation will versionne the modifications applied to one or many files.

Figure 14 presents the data structure used to store a modification or commit. Each commit is represented as a tree. The root leaf (green) contains the commit, tree and parent hashes as same as the author and the description associated to the commit. The second leaf (blue) contains the leaf hash and the hashes of the files of the project.

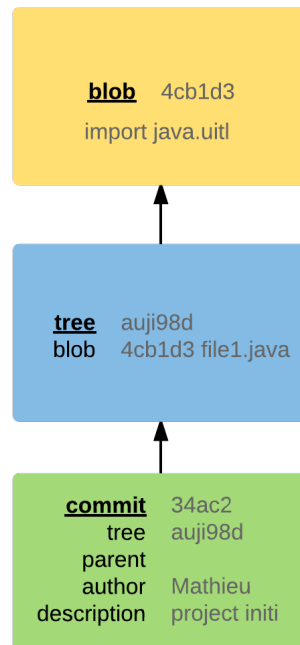


Figure 14: Data structure of a commit.

In this example, we can see that author “Mathieu” has created the file *file1.java* with the message “project init”. Figure 15 represents an ulterior modification. In this second example, *file1.java* is modified while *file2.java* is created. The second commit 98ca9 have 34ac2 as parent.

Branches point to a commit. In a task-branching environment, a branch is created via a checkout operation for each task. Task can be to fix the root cause of crash or bug report or features to implement. In figure 16, the *master* branch and the *1_fix_overflow* point on commit 98ca9.

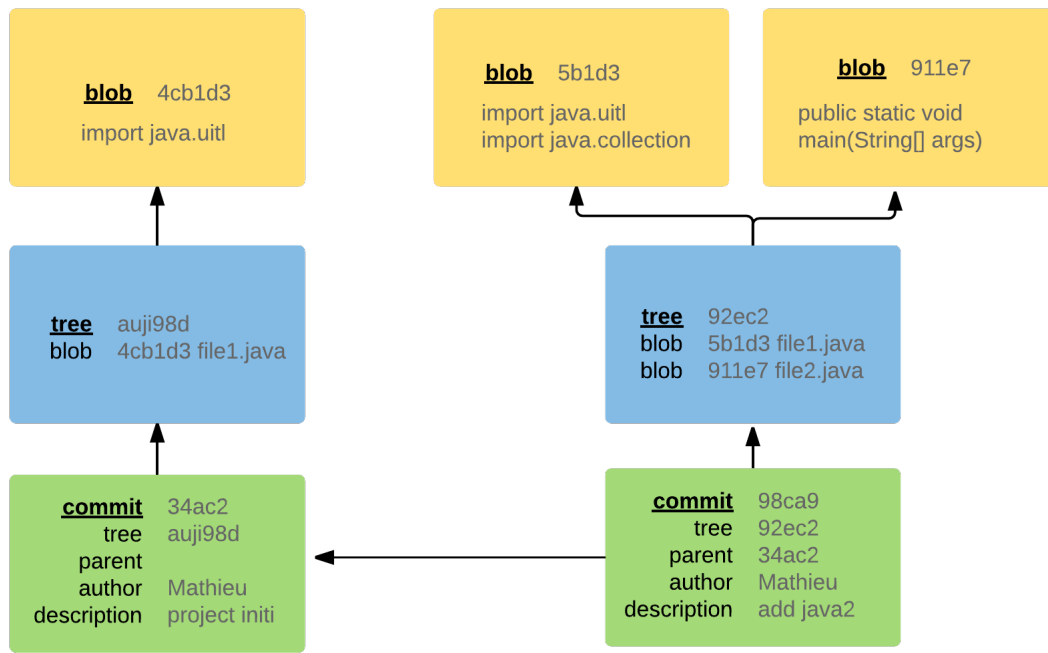


Figure 15: Data structure of two commits.

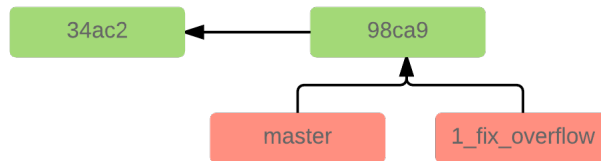


Figure 16: Two branches pointing on one commit.

Both branches can evolve separately and be merged together when the task branch is ready. In figure 17, the *master* branch points on *a13ab2* while the *1_fix_overflow* points on *ahj23k*.

In this chapter, we present approaches that have the characteristics described by human maintenaire and integrate themselves seamlessly in a task branching environment.

4.1 Software maintenance at branching-time

Software maintenance at branching-time allows humain maintainers to obtain, in an automatic fashion, the steps to reproduce the root cause of a bug or crash report.

As an example, if a maintainer is assigned to fix the root cause of bug report #38622 of the Ant project, related to parsing exception, he will create a branch named *38622_fix_parsing*

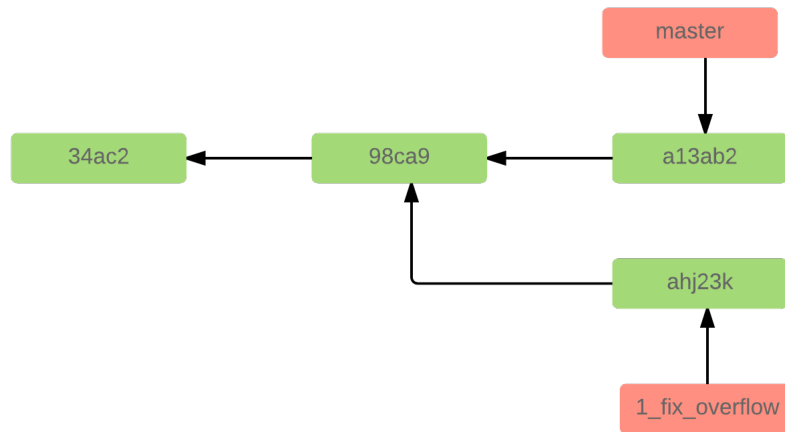


Figure 17: Two branches pointing on two commits.

```
#!/bin/sh

#Parsing the branch name
arr=$(git branch | tr "_" "\n")

#save reproduction steps to the ./tests directory
wget https://bumper-app.com/api/testsuites/$arr[0] -P $(pwd)/
  ↪ tests
```

Figure 18: Pre-checkout hook

with the command: `git checkout -b 38622_fix_parsing` for Git.

Before the checkout, our pre-checkout hook triggers. As a reminder, a hook is a bash script that can be executed before or after any versionning operation. Figure 18 presents our pre-checkout hook. The first operation is to parse the branch name in order to extract the task *#id*. With the task *#id* we query the BUMPER’s API for the test suite created by JCHARMING. Then, we write the response to the `/tests/` directory of the project.

Figure 19 shows the Java test suite that would have been downloaded and added to the `/tests/` directory in this case.

The combination of BUMPER, JCHARMING and this pre-checkout hook provide maintainers actionable and contextualized information in a non-intrusive manner. This information can be used to ease the maintenance process.

```

public class Ant-38622 extends TestCase {

    /**
     * JCHARMING GENERATED TEST — USE WITH CAUTION
     */
    @Test
    public testAnt() {

        org.apache.tools.ant.helper.ProjectHelper2 v0 = org.
            ↪ apache.tools.ant.ProjectHelper.getProjectHelper
            ↪ ();
        org.apache.tools.ant.Project v1 = new org.apache.
            ↪ tools.ant.Project();
        v0.parse(v1, null);
    }
}

```

Figure 19: JCHARMING generated test imported into task branch

4.2 Software maintenance at commit-time

Software maintenance at commit-time allows developers to receive recommendation and improvement at commit-time using code clones detection techniques.

Code clones appear when developers reuse code with little to no modification to the original code. Studies have shown that clones can account for about 7% to 50% of code in a given software system[Bak, ?]. Developers often reuse code (and create clones) in their software on purpose[KSN05]. Nevertheless, clones are considered a bad practice in software development since they can introduce new bugs in the code[KG06, JDHW09, LLMZ06]. If a bug is discovered in one segment of the code that has been copied and pasted several times, then the developers will have to remember the places where this segment has been reused in order to fix the bug in each place.

In the last two decades, there have been many studies and tools that aim at detecting clones. They can be grouped into three categories. The first category includes techniques that treat the source code as text and use transformation and normalization methods to compare various code fragments[Joh94, Joh93, CR11, RC08]. The second category includes methods that use lexical analysis, where the source code is sliced into sequences of tokens,

similar to the way a compiler operates[Bak, Bak92, BG02, KKI02b, LLMZ06]. The tokens are used to compare code fragments. Finally, syntactic analysis has also been performed where the source code is converted into trees, more particularly abstract syntax tree (AST), and then the clone detection is performed using tree matching algorithms[BYM⁺98, KH00, TG06, FFK08].

Many taxonomies have been published in an attempt to classify clones into types. [MLM96, BMD⁺99, KFF06, BKA⁺07, Kon, KG]. Despite the particularities of each proposed taxonomy, researchers agree on the following classification. Type 1 clones are copy-pasted blocks of code that only differ from each other in terms of non-code artefacts such as indentation, whitespaces, comments and so on. Type 2 clones are blocks of code that are syntactically identical at the exception of literals, identifiers and types that can be modified. In addition, Type 2 clones share the particularities of Type 1 about indentation, whitespaces and comments. Type 3 clones are similar to Type 2 clones in terms of modification of literals, identifiers, types, indentation, whitespaces and comments but also contain added or deleted code statements. Finally, Type 4 are code blocks that perform the same tasks, but using a completely different implementation.

4.2.1 PRECINCT: PREventing Clones INsertion at Commit Time

In this section, we present PRECINCT (PREventing Clones INsertion at Commit Time) that focuses on preventing the insertion of clones at commit time, i.e., before they reach the central code repository. PRECINCT is an online clone detection technique that relies on the use of pre-commit hooks capabilities of modern source code version control systems. PRECINCT intercepts this modification and analyses its content to see whether a suspicious clone has been introduced or not. A flag is raised if a code fragment is suspected to be a clone of an existing code segment. In fact, PRECINCT, itself, can be seen as a pre-commit hook that detects clones that might have been inserted in the latest changes with regard to the rest of the source code. This said, only a fraction of the code is analysed, making PRECINCT efficient compared to leading clone detection techniques such as NICAD (Accurate Detection of Near-miss Intentional Clones) [CR11]. Moreover, the detected clones are presented using a classical ‘diff’ output that developers are familiar with. PRECINCT is also well integrated with the workflow of the developers since it is used in conjunction with a source code version control systems such as Git¹.

In this study, we focus on Type 3 clones as they are more challenging to detect. Since Type 3 clones include Type 1 and 2 clones, then these types could be detected separately

¹<https://git-scm.com/>

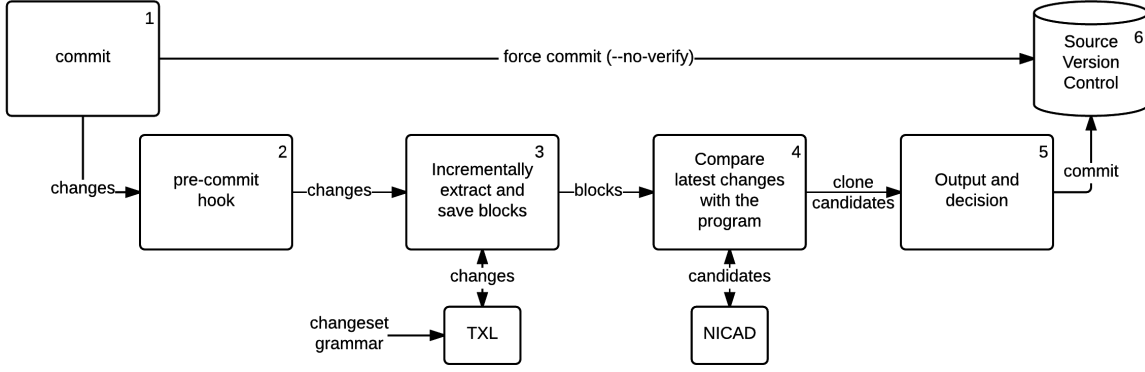


Figure 20: Overview of the PRECINCT Approach.

by PRECINCT as well.

We evaluated the effectiveness of PRECINCT using precision and recall on three systems, developed independently and written in both C and Java. The results show that PRECINCT prevents Type 3 clones to reach the final source code repository with an average accuracy of 97.7%.

The rest of this paper is organized as follows: In Section ??, we present the studies related to PRECINCT. Then, in Section 4.2.1, we present the PRECINCT approach. The evaluation of PRECINCT is the subject of Section 4.2.1, followed by threats to validity. Finally, we conclude the paper in Section 4.2.1.

The PRECINCT Approach

The PRECINCT approach is composed of six steps. The first and last steps are typical steps that a developer would do when committing code. Indeed, the first step is the commit step where developers send their latest changes to the central repository and the last step is the reception of the commit by the central repository. The second step is the pre-commit hook, which kicks in as the first operation when one wants to commit. The pre-commit hook has access to the changes in terms of files that have been modified, more specifically, the lines that have been modified. The modified lines of the files are sent to TXL[Cor06a] for block extraction. Then, the blocks are compared to previously extracted blocks in order to identify candidate clones using the comparison engine of NICAD[CR11]. We chose NICAD engine because it has been shown to provide high accuracy [CR11]. The tool is also readily available, easy to use, customizable, and works with TXL. Note, however, that PRECINCT can also work with other engines for comparing code fragments. Finally, the output of NICAD is further refined and presented to the user for decision. These steps are discussed in more detail in the following subsections.

PRECINCT Pre-Commit Hook

Depending on the exit status of the hook, the commit will be aborted and not pushed to the central repository. Also, developers can choose to ignore the pre-hook. In Git, for example, they will need to use the command `git commit --no-verify` instead of `git commit`. This can be useful in case of an urgent need for fixing a bug where the code has to reach the central repository as quickly as possible. Developers can do things like check for code style, check for trailing white spaces (the default hook does exactly this), or check for appropriate documentation on new methods.

PRECINCT is a set of bash scripts where the entry point of these scripts lies in the pre-commit hooks. Pre-commit hooks are easy to create and implement as depicted in Listing 4.1. This pre-hook is shipped with Git, a popular version control system. Note that even though we use Git as the main version control to present PRECINCT, we believe that the techniques presented in this paper are readily applicable to other version control systems. In Listing 4.1, from lines 3 to 11, the script identifies if the commit is the first one in order to select the revision to work against. Then, in Lines 18 and 19, the script checks for trailing whitespace and fails if any are found.

For PRECINCT to work, we just have to add the call to our script suite instead or in addition of the whitespace check.

Listing 4.1: Git Pre-Commit Hook Sample

<code>#!/bin/sh</code>	1
	2
<code>if git rev-parse --verify HEAD > \</code>	3
<code> /dev/null 2>&1</code>	4
<code>then</code>	5
<code> against=HEAD</code>	6
<code>else</code>	7
<code> # Initial commit: diff against</code>	8
<code> # an empty tree object</code>	9
<code> against=4b825dc642....</code>	10
<code>fi</code>	11
	12
<code># Redirect output to stderr.</code>	13
<code>exec 1>&2</code>	14
	15
<code># If there are whitespace errors,</code>	16
<code># print the offending file names and fail.</code>	17
<code>exec git diff-index --check \</code>	18
<code> --cached \$against --</code>	19

Extract and Save Blocks

A block is a set of consecutive lines of code that will be compared to all other blocks in order to identify clones. To achieve this critical part of PRECINCT, we rely on TXL[Cor06a], which is a first-order functional programming over linear term rewriting, developed by Cordy et al.[Cor06a]. For TXL to work, one has to write a grammar describing the syntax of the source language and the transformations needed. TXL has three main phases: *parse*, *transform*, *unparse*. In the parse phase, the grammar controls not only the input but also the output form. Listing 4.2 — extracted from the official documentation² — shows a grammar matching an *if-then-else* statement in C with some special keywords: [IN] (indent), [EX] (exdent) and [NL] (newline) that will be used for the output form.

²<http://txl.ca>

Listing 4.2: Txl Sample Sample

define if_statement	1
if ([expr]) [IN][NL]	2
[statement] [EX]	3
[opt else_statement]	4
end define	5
	6
define else_statement	7
else [IN][NL]	8
[statement] [EX]	9
end define	10

Then, the *transform* phase will, as the name suggests, apply transformation rules that can, for example, normalize or abstract the source code. Finally, the third phase of TXL, called *unparse*, unparses the transformed parsed input in order to output it. Also, TXL supports what the creators call Agile Parsing[DCMS], which allow developers to redefine the rules of the grammar and, therefore, apply different rules than the original ones.

PRECINCT takes advantage of that by redefining the blocks that should be extracted for the purpose of clone comparison, leaving out the blocks that are out of scope. More precisely, before each commit, we only extract the blocks belonging to the modified parts of the source code. Hence, we only process, in an incremental manner, the latest modification of the source code instead of the source code as a whole.

We have selected TXL for several reasons. First, TXL is easy to install and to integrate with the normal workflow of a developer. Second, it was relatively easy to create a grammar that accepts commits as input. This is because TXL is shipped with C, Java, Csharp, Python and WSDL grammars that define all the particularities of these languages, with the ability to customize these grammars to accept changesets (chunks of the modified source code that include the added, modified, and deleted lines) instead of the whole code.

Algorithm ?? presents an overview of the “extract” and “save” blocks operations of PRECINCT. This algorithm receives as arguments, the changesets, the blocks that have been previously extracted and a boolean named `compare_history`. Then, from Lines 1 to 9 lie the *for* loop that iterates over the changesets. For each changeset (Line 2), we extract the blocks by calling the `extract_blocks(Changeset cs)` function. In this function, we expand our changeset to the left and to the right in order to have a complete block.

As depicted by Listing ??, changesets contain only the modified chunk of code and not necessarily complete blocks. Indeed, we have a block from Line 3 to Line 6 and deleted

lines from Line 8 to 14. However, in Line 7 we can see the end of a block, but we do not have its beginning. Therefore, we need to expand the changeset to the left in order to have syntactically correct blocks. We do so by checking the block’s beginning and ending (using { and }) in C for example. Then, we send these expanded changesets to TXL for block extraction and formalization.

For each extracted block, we check if the current block overrides (replaces) a previous block (Line 4). In such a case, we delete the previous block as it does not represent the current version of the program anymore (Line 5). Also, we have an optional step in PRECINCT defined in Line 4. The `compare_history` is a condition to delete overridden blocks.

We believe that deleted blocks have been deleted for a good reason (bug, default, removed features, ...) and if a newly inserted block matches an old one, it could be worth knowing in order to improve the quality of the system at hand. This feature is deactivated by default.

In summary, this step receives the files and lines, modified by the latest changes made by the developer and produces an up to date block representation of the system at hand in an incremental way. The blocks are analysed in the next step to discover potential clones.

Compare Extracted Blocks

In order to compare the extracted blocks and detect potential clones, we can only resort to text-based techniques. This is because lexical and syntactic analysis approaches (alternatives to text-based comparisons) would require a complete program to work, a program that compiles. In the relatively wide-range of tools and techniques that exist to detect clones by considering code as text[Joh93, Joh94, MM, Man94, ?, WM05], we selected NICAD as the main text-based method for comparing clones [CR11] for several reasons. First, NICAD is built on top of TXL, which we also used in the previous step. Second, NICAD is able to detect all Types 1, 2 and 3 software clones.

NICAD works in three phases: *Extraction*, *Comparison* and *Reporting*. During the *Extraction* phase all potential clones are identified, pretty-printed, and extracted. We do not use the *Extraction* phase of NICAD as it has been built to work on programs that are syntactically correct, which is not the case for changesets. We replaced NICAD’s *Extraction* phase with our own scripts, described in the previous section.

In the *Comparison* phase, extracted blocks are transformed, clustered and compared in order to find potential clones. Using TXL sub-programs, blocks go through a process called pretty-printing where they are stripped of formatting and comments. When code

fragments are cloned, some comments, indentation or spacing are changed according to the new context where the new code is used. This pretty-printing process ensures that all code will have the same spacing and formatting, which renders the comparison of code fragments easier. Furthermore, in the pretty-printing process, statements can be broken down into several lines. Table ?? shows how this can improve the accuracy of clone detection with three `for` statements, `for (i=0; i<10; i++)`, `for (i=1; i<10; i++)` and `for (j=2; j<100; j++)`. The pretty-printing allows NICAD to detect Segments 1 and 2 as a clone pair because only the initialization of *i* changed. This specific example would not have been marked as a clone by other tools we tested such as Duploc[DRD99]. In addition to the pretty-printing, code can be normalized and filtered to detect different classes of clones and match user preferences.

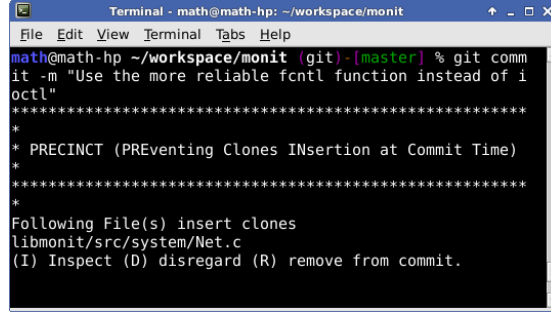
Finally, the extracted, pretty-printed, normalized and filtered blocks are marked as potential clones using a Longest Common Subsequence (LCS) algorithm[HS77]. Then, a percentage of unique statements can be computed and, depending on a given threshold (see Section 4.2.1), the blocks are marked as clones.

The last step of NICAD, which acts as our clone comparison engine, is the *reporting*. However, to prevent PRECINCT from outputting a large amount of data (an issue that many clone detection techniques face), we implemented our own reporting system, which is also well embedded with the workflow of developers. This reporting system is the subject of the next section.

As a summary, this step receives potentially expanded and balanced blocks from the extraction step. Then, the blocks are pretty-printed, normalized, filtered and fed to an LCS algorithm in order to detect potential clones. Moreover, the clone detection in PRECINCT is less intensive than NICAD because we only compare the latest changes with the rest of the program instead of comparing all the blocks with each other.

Output and Decision

In this final step, we report the result of the clone detection at commit time with respect to the latest changes made by the developer. The process is straightforward. Every change made by the developer goes through the previous steps and is checked for the introduction of potential clones. For each file that is suspected to contain a clone, one line is printed to the command line with the following options: (I) Inspect, (D) Disregard, (R) Remove from the commit as shown by Figure 21. In comparison to this simple and interactive output, NICAD outputs each and every detail of the detection result such as the total number of potential clones, the total number of lines, the total number of unique line text chars, the



```
Terminal - math@math-hp: ~/workspace/monit
File Edit View Terminal Tabs Help
math@math-hp ~/workspace/monit (git) - [master] % git comm
it -m "Use the more reliable fcntl function instead of i
octl"
*****
* PRECINCT (PREventing Clones INsertion at Commit Time)
*
*****
*
Following File(s) insert clones
libmonit/src/system/Net.c
(I) Inspect (D) disregard (R) remove from commit.
```

Figure 21: PRECINCT output when replaying commit 710b6b4 of the Monit system used in the case study.

total number of unique lines, and so on. We think that so many details might make it hard for developers to react to these results. A problem that was also raised by Johnson et al. [JSMHB13] when examining bug detection tools. Then the potential clones are stored in XML files that can be viewed using an Internet browser or a text editor.

(I) Inspect will cause a diff-like visualization of the suspected clones while (D) disregard will simply ignore the finding. To integrate PRECINCT in the workflow of the developer we also propose the remove option (R). This option will simply remove the suspected file from the commit that is about to be sent to the central repository. Also, if the user types an option key twice, e.g., II, DD or RR, then the option will be applied to all files. For instance, if the developer types DD at any point, the PRECINCT's results will be disregarded and the commit will be allowed to go through. We believe that this simple mechanism will encourage developers to use PRECINCT like they would use any other feature of Git (or any other version control system).

Case Study

In this section, we show the effectiveness of PRECINCT for detecting clones at commit time in three open source systems³.

The aim of the case study is to answer the following question: *Can we detect clones at commit time, i.e., before they are inserted in the final code, if so, what would be the accuracy compared to a traditional clone detection tool such as NICAD?*

³The programs used and instructions to reproduce the experiments are made available for download from <https://research.mathieu-nayrolles.com/precinct/>

Table 4: List of Target Systems in Terms of Files and Kilo Line of Code (KLOC) at current version and Language

SUT	Revisions	Files	KLoC	Language
Monit	826	264	107	C
Jhotdraw	735	1984	44	Java
dnsjava	1637	233	47	Java

Target Systems

Table 4 shows the systems used in this study and their characteristics in terms of the number files they contain and the size in KLoC (Kilo Lines of Code). We also include the number of revisions used for each system and the programming language in which the system is written.

Monit⁴ is a small open source utility for managing and monitoring Unix systems. Monit is used to conduct automatic maintenance and repair and supports the ability to identify causal actions to detect errors. This system is written in C and composed of 826 revisions, 264 files, and the latest version has 107 KLoC. We have chosen Monit as a target system because it was one of the systems NICAD was tested on.

JHotDraw⁵ is a Java GUI framework for technical and structured graphics. It has been developed as a “design exercise”. Its design relies heavily on the use of design patterns. JHotDraw is composed of 735 revisions, 1984 files, and the latest revision has 44 KLoC. It is written in Java and it is often used by researchers as a test bench. JHotDraw was also used by NICAD’s developers to evaluate their approach.

Dnsjava⁶ is a tool for implementing the DNS (Domain Name Service) mechanisms in Java. This tool can be used for queries, zone transfers, and dynamic updates. It is not as large as the other two, but it still makes an interesting case subject because it has been well maintained for the past decade. Also, this tool is used in many other popular tools such as Aspirin, Muffin and Scarab. Dnsjava is composed of 1637 revisions, 233 files, the latest revision contains 47 KLoC. We have chosen this system because we are familiar with it as we used it before[?].

Process

Figure 22 shows the process we followed to validate the effectiveness of PRECINCT.

⁴<https://mmonit.com/monit/>

⁵<http://www.jhotdraw.org/>

⁶<http://www.dnsjava.org/>

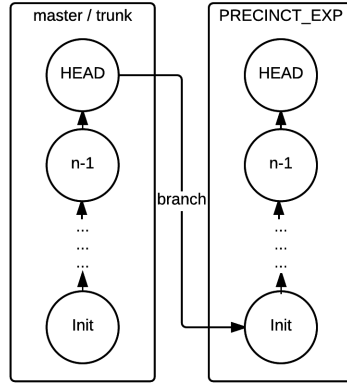


Figure 22: PRECINCT Branching.

As our approach relies on commit pre-hooks to detect possible clones during the development process (more particularly at commit time), we had to find a way to *replay* past commits. To do so, we *cloned* our test subjects, and then created a new branch called *PRECINCT_EXT*. When created, this branch is reinitialized at the initial state of the project (the first commit) and each commit can be replayed as they have originally been. At each commit, we store the time taken for PRECINCT to run as well as the number of detected clone pairs. We also compute the size of the output in terms of the number of lines of text output by our method. The aim is to reduce the output size to help software developers interpret the results.

To validate the results obtained by PRECINCT, we needed to use a reliable clone detection approach to extract clones from the target systems and use these clones as a baseline for comparison. For this, we turned to NICAD because of its popularity, high accuracy, and availability [CR11], as discussed before. This means, we run NICAD on the revisions of the system to obtain the clones then we used NICAD clones as a baseline for comparing the results obtained by PRECINCT.

It may appear strange that we are using NICAD to validate our approach, knowing that our approach uses NICAD’s code comparison engine. In fact, what we are assessing here is the ability for PRECINCT to detect clones at commit time using changsets. The major part of PRECINCT is the ability to intercept code changes and build working code blocks that are fed to a code fragment engine (in our case NICAD’s engine). PRECINCT can be built on the top of any other code comparison engine.

We show the result of detecting Type 3 clones with a maximum line difference of 30% as discussed in Table ???. As discussed in the introductory section, we chose to report on Type 3 clones because they are more challenging to detect than Type 1 and 2. PRECINCT detects Type 1 and 2 too so does NICAD. For the time being, PRECINCT is not designed

to detect Type 4 clones. These clones use different implementations. Detecting Type 4 clones is part of future work.

We assess the performance of PRECINCT in terms of precision (Equation 1) and recall (Equation 2). Both precision and recall are computed by considering NICAD’s results as a baseline. We also compute F_1 -measure (Equation 3), i.e., the weighted average of precision and recall, to measure the accuracy of PRECINCT.

$$precision = \frac{|\{NICAD_{detection}\} \cap \{PRECINCT_{detection}\}|}{|\{PRECINCT_{detection}\}|} \quad (8)$$

$$recall = \frac{|\{NICAD_{detection}\} \cap \{PRECINCT_{detection}\}|}{|\{NICAD_{detection}\}|} \quad (9)$$

$$F_1 - measure = 2 * \frac{precision * recall}{precision + recall} \quad (10)$$

Results

Figures ??, ??, ?? show the results of our study in terms of clone pairs that are detected per revision for our three subject systems: Monit, JHotDraw and Dnsjava. We used as baseline for comparison the clone pairs detected by NICAD. The blue line shows the clone detection performed by NICAD. The red line shows the clone pairs detected by PRECINCT. The brown line shows the clone pairs that have been missed by PRECINCT. As we can quickly see, the blue and red lines almost overlap, which indicates a good accuracy of the PRECINCT approach.

Table ?? summarizes PRECINCT’s results in terms of precision, recall, F_1 -measure, execution time and output reduction. The first version of Monit contains 85 clone pairs and this number stays stable until Revision 100. From Revision 100 to 472 the detected clone pairs vary between 68 and 88 before reaching 219 at Revision 473. The number of clone pairs goes down to 122 at Revision 491 and decreases to 128 in the last revision. PRECINCT was able to detect 96.1% (123/128) of the clone pairs that are detected by NICAD with a 100% recall. It took in average around 1 second for PRECINCT to execute on a Debian 8 system with Intel(R) Core(TM) i5-2400 CPU @ 3.10GHz, 8Gb of DDR3 memory. It is also worth mentioning that the computer we used is equipped with SSD (Solid State Drive). This impacts the running time as clone detection is a file intensive operation. Finally, the PRECINCT was able to output 88.3% less lines than NICAD.

JHotDraw starts with 196 clone pairs at Revision 1 and reaches a pick of 2048 at Revision 180. The number of clones continues to go up until Revisions 685 and 686 where the number of pairs is 1229 before picking at 6538 and more from Revisions 687 to 721.

PRECINCT was able to detect 98.3% of the clone pairs detected by NICAD (6490/6599) with 100% recall while executing on average in 1.7 second (compared to 5.1 seconds for NICAD). With JHotDraw, we can clearly see the advantages of incremental approaches. Indeed, the execution time of PRECINCT is loosely impacted by the number of files inside the system as the blocks are constructed incrementally. Also, we only compare the latest change to the remaining of the program and not all the blocks to each other as NICAD. We also were able to reduce by 70.1% the number of lines output by NICAD.

Finally, for Dnsjava, the number of clone pairs starts high with 258 clones and goes up until Revision 70 where it reaches 165. Another quick drop is observed at Revision 239 where we found only 25 clone pairs. The number of clone pairs stays stable until Revision 1030 where it reaches 273. PRECINCT was able to detect 82.8% of the clone pairs detected by NICAD (226/273) with 100% recall, while executing on average in 1.1 second while NICAD took 3 seconds in average. PRECINCT outputs 83.4% less lines of code than NICAD.

Overall, PRECINCT prevented 97.7% of the 7000 clones (in all systems) to reach the central source code repository while executing more than twice as fast as NICAD (1.2 seconds compared to 3 seconds in average) while reducing the output in terms of lines of text output the developers by 83.4% in average. Note here that we have not evaluated the quality of the output of PRECINCT compared to NICAD's output. We need to conduct user studies for this. We are, however, confident, based on our own experience trying many clone detection tools, that a simpler and more interactive way to present the results of a clone detection tool is warranted. PRECINCT aims to do just that.

The difference in execution time between NICAD and PRECINCT stems from the fact that, unlike PRICINCT, NICAD is not an incremental approach. For each revision, NICAD has to extract all the code blocks and then compares all the pairs with each other. On the other hand, PRECINCT only extracts blocks when they are modified and only compares what has been modified with the rest of the program.

The difference in precision between NICAD and PRECINCT (2.3%) can be explained by the fact that sometimes developers commit code that does not compile. Such commits will still count as a revision, but TXL fails to extract blocks that do not comply with the target language syntax. While NICAD also fails in such a case, the disadvantage of PRECINCT comes from the fact that the failed block is saved and used as reference until it is changed by a correct one in another commit.

Threats to Validity

The selection of target systems is one of the common threats to validity for approaches aiming to improve the analysis of software systems. It is possible that the selected programs share common properties that we are not aware of and therefore, invalidate our results. However, the systems analysed by PRECINCT are the same as the ones used in similar studies. Moreover, the systems vary in terms of purpose, size and history.

In addition, we see a threat to validity that stems from the fact that we only used open source systems. The results may not be generalizable to industrial systems. We intend to undertake these studies in future work.

The programs we used in this study are all based on the Java, C and Python programming languages. This can limit the generalization of the results. However, similar to Java, C, Python, if one writes a TXL grammar for a new language — which can be a relatively hard work — then PRECINCT can work since PRECINCT relies on TXL.

Finally, we use NICAD as the code comparison engine. The accuracy of NICAD affects the accuracy of PRECINCT. This said, since NICAD has been tested on large systems, we are confident that it is a suitable engine for comparing code using TXL. Also, there is nothing that prevents us from using other code comparisons engines, if need be.

In conclusion, internal and external validity have both been minimized by choosing a set of three different systems, using input data that can be found in any programming languages and version systems (commit and changesets).

Conclusion

PRECINCT (PREventing Clones INSERTion at Commit Time) is an incremental approach for preventing clone insertion at commit time that combines efficient block extraction and clone detection and integrate itself seamlessly in the day-to-day workflow of developers. PRECINCT takes advantage of TXL and NICAD to create a clone detection tool approach that runs automatically before each commit in 1.2 second with a 97.7% precision and a 100% recall (when using NICAD results as a baseline).

Our approach also assesses two major factors that contribute to the slow adoption of clone detection tools: large number of data output by clone detection methods, and smooth integration with the task flow of the developers. PRECINCT is able to reduce the number of lines output by a classical clone detection tool such as NICAD by 83.4% while keeping all the necessary information that allow developers to decide whether the detected clone is in fact a clone. Also, our approach is seamlessly integrated with the developers' workflow by means of pre-commit hooks, which are part any version control systems.

To build on this work, we need to experiment with additional (and larger) systems with the dual aim to (a) improve and fine-tune the approach, and (b) assess the scalability of our approach when applied to even larger (and proprietary) systems. Also, we want to improve PRECINCT to support Type 4 clones.

4.2.2 RESEMBLE - REcommendation System based on cochangeE Mining at Block LLevel

RESEMBLE (REcommendation System based on cochangeE Mining at Block LLevel) is a contextual recommendation system that will take the form of another pre-commit hook. RESEMBLE will leverage the branches' history and the blocks extracted by PRECINCT to identify, on the fly, sub-optimum or hazardous code and display an actual solution to the problem.

Mining sequences of changes

In the data mining field, ARM is a well-established method for discovering co-occurrences between attributes in the objects of a large data set [G. 91, HHA97]. Plain associations have the form $X \rightarrow Y$, where X and Y , called the *antecedent* and the *consequent*, respectively, are sets of descriptors (purchases by a customer, network alarms, or any other general kind of events). Even though plain association rules could serve some relevant information, we are interested here in the sequences of changes that we believe will yield more precise result. Indeed, we think that similar modifications are often done in the same order by the same developer (e.g top to bottom or bottom to top). We, therefore, adopt a variant called sequential association rules in which both X and Y become sequences of descriptors. Moreover, our sequences follow a temporal order with the antecedent preceding the consequent. Rules of this type mined from changes reveal crucial information about the likelihood of blocks of code to be modified together in a programming session and, more importantly, in a specific order. For instance, a strong rule $Block_A, Block_B$ implies $Block_C$ would mean that after modifying $Block_A$ and then $Block_B$, there are good chances that the developer needs to modify $Block_C$. The conciseness of this example should not confuse the reader as in practical cases the sequences appearing in a rule can be of an arbitrary length. Furthermore, the strength of the rule is measured by the *confidence* metric. In probabilistic terms, it measures the conditional probability of C appearing down the line. Beside that, the significance of a rule, i.e. how many times it appears in the data, is provided by its *support* measure. To ensure only rules of potentially high interestingness are mined, the mining task is tuned by minimal thresholds to output only the sufficiently high scores for both metrics.

To extract the association rules from changes, two choices were possible. On one hand, sequential pattern mining and rule mining algorithms have been designed for structures that are slightly more general than the ones used here. In fact, sequential patterns are defined on transactions that represent sequences of sets. Efficient sequential pattern miners have been published, e.g. the PrefixSpan method [PHM⁺04]. On the other hand, sequence of changes do not compile to fully-blown sequential transactions as the underlying structures are mere sequences of individual elements. Such data has been known since at least the mid-90s but received less attention by the data mining community, arguably because it is less challenging to mine. In the general data mining literature, mining from pure sequences, as opposed to sequences made of sets, has been addressed under the name of episode mining [HHA97]. Episodes are made of *events* and in a sense, code changes are events. Arguably the largest body of knowledge on the subject belongs to the web usage mining field: The input data is again a system log, yet this time the log of requests sent to a web server [PHMa00]. It is noteworthy that sequential patterns are more general than the pure sequence ones, hence mining algorithms designed for the former might prove to be less efficient when applied to the latter (as additional steps might be required for listing all significant set). Nevertheless, to jump-start our experimental study, we used a sequential pattern/rule miner that has the advantage to be freely available on the web⁷. Although it has not been optimized for pure sequences its performances are more than satisfactory.

Code Normalization

Normalizing code is the action of making its structure consistent throughout the program according to a defined model. In our case, we only normalize the block of code that are in the current sequential association rules. To do so, we improve and combine several technologies such as source code transformation [Cor06b, Cor06a], source code pretty-printing [RC08], flexible source code normalization [CR11].

More specifically, the code first goes to a pretty printer. A pretty-printer is a component that will slightly transform the code in order to obtain consistent control structure. Concretely, spaces will be added, accolade moved and tabulation added for an *if*, a *while* and others structures to always appear the same way, regardless of the programming language. Then, the code is normalized several times. Each normalization is different and targets specific feature of the code. For example, we have one normalization that removes completely the variable names and replace the types by the highest known object in the object oriented hierarchy before `Object` itself⁸. Another normalization only keeps the structures

⁷<http://www.philippe-fournier-viger.com/spmf/>

⁸For Java programs

of the source code by normalizing both variables names and types.

Comparing normalization

Comparing different normalizations is the easiest step and can be done efficiently using the longest common sequence (LCS) algorithm [Hir77]. If the LCS is above an user-defined threshold, then a two different behaviors can be observed:

- If the modified blocks' — and potentially the ones that are likely to be modified after according to our sequential association rules — normalizations match the normalizations of blocks of code that have been removed in past history. **RESEMBLE** recommends the replacing code as a better solution. In other words, if blocks A , B and C have been replaced by blocks A' , B' and C' and blocks D , E and F normalizations match A , B and C , then, **RESEMBLE** recommends to transform D , E and F to look like A' , B' and C' . This could lead to the introduction of software clones but we argue that (i) software clones are not always harmful [JDHW09] and (ii) informing the developer that A' , B' and C' exist could lead him to re-use these blocks.
- If the modified blocks' — and potentially of the ones that are likely to be modified after according to our sequential association rules — normalizations match the normalization of blocks of code that are present in the history. **RESEMBLE** computes the differences between the developer code and the history code in order to suggest what the developer have to do next.

Planned experiments

We did not started the experiments for **RESEMBLE** yet. Nevertheless, we plan to conduct the following experiments:

- Full history test with Normalization 1.
- Full history test with Normalization 2.
- An human study where:
 - Developers use **RESEMBLE** in order to determine whether or not developers take into account our recommendations to avoid inserting defects in the code.
 - Developers use **RESEMBLE** in order to determine whether or not developers take into account our recommendations to complete their modification according to what we found in the history.

- Rate the suggested solution in a scale from 1 to 10 in order to determine if the proposed change pattern does resolve the current problem.

We believe that **RESEMBLE** will be a real asset in a developer tool belt in order to ship better code in terms of quality, performances and security. However, as **RESEMBLE** aims to provide recommendations in at commit-time, using the local history and ressources, it will not be able to be as exhaustive as an offline process. To fill this gap, we built **BIANCA** that we present in the next section.

In summary, using **PRECINCT** and **RESEMBLE** we are able to provide contextualized and accurate information in a non-intrusive manner.

4.3 Software maintenance at merge-time

BIANCA (Bug Insertion ANTicipation by Clone Analysis at commit time) is the final piece of the proposed ecosystem and, as such, the final failsafe that prevents developer to ship code that we know to be sub-optimum or to be at the very root of issues.

Many tools exist to prevent a developer to ship *bad* code [Dan00, Hov07, MGD10] or to identify *bad* code after executions (e.g in test or production environment) [NMHIL, NMV13]. However, these tools rely on metrics and rules to statically and/or dynamically identify sub-optimum code.

The **BIANCA** approach

BIANCA is different than the approaches presented in the previous sections because it mines and analyzes the change patterns in commits and matches it against past commits known to have introduced a defect in the code (or that have just been replaced by better implementation). Also, **BIANCA** is an offline approach that is triggered by a merge request. When a maintainers estimate that their work are ready to be integrated with the main branch, they open a merge request⁹. Merging a task branch is not an instantaneous process as the code need to pass code review. **BIANCA** leverages this *down* time to perform a complete history check on all projects contained in **BUMPER**.

Figure 23 presents an overview of our approach.

BIANCA builds a model where each issue is represented by three versions of the same file. These three versions are stored in **BUMPER**. The first version n is called the *stable state* because the code of this version was used to fix an issue. The $n - 1$ version, however, is called the *unstable state* as it was marked as containing an issue. Finally, the third version is

⁹Also known as pull request

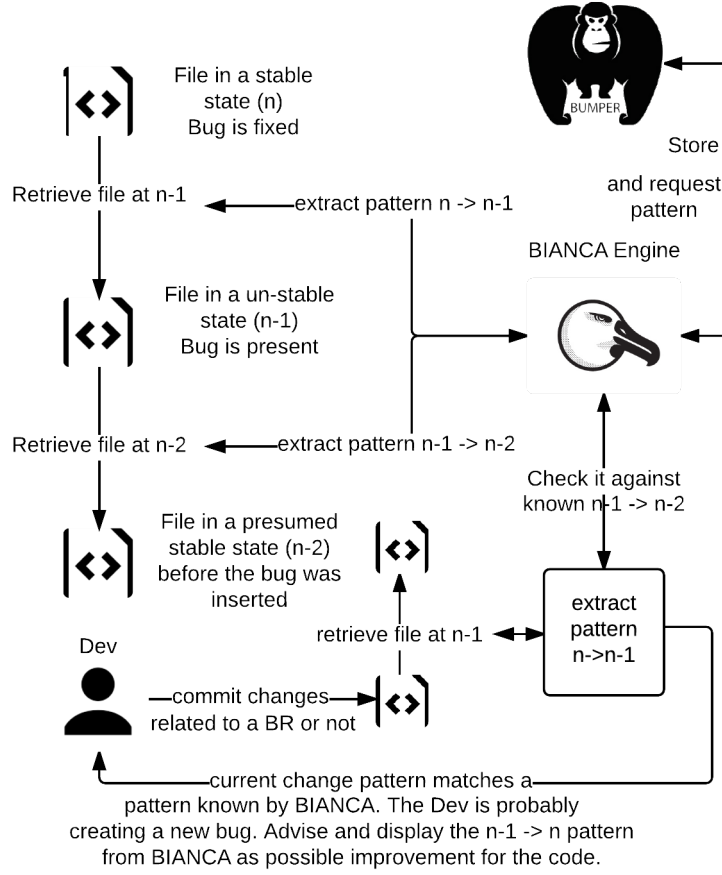


Figure 23: The BIANCA Approach

called the *before state* and represents the file before the introduction of the bug. Hereafter, we refer to the *before state* as $n - 2$. BIANCA extracts the change patterns from $n - 2$ to $n - 1$ and from $n - 1$ to n . It also generates the changes to go from $n - 2$ to n .

When a developer commits new modifications, BIANCA extracts the change pattern from the version n_{dev} (current version) and $n-1_{dev}$ (version before modification) of the developer's source code and compares this change pattern to known $n - 2$ to $n - 1$ patterns. If n_{dev} to $n - 1_{dev}$ matches a $n - 2$ to $n - 1$ then it means that the developer is inserting a known defect in the source code. In such a case, BIANCA will propose the related $n - 1$ to n pattern to the developer, so s/he could improve the source code and will show the related $n - 2$ for the n pattern so the developer will learn how to s/he should have modified the code in the first place.

Moreover, if the issue was previously reproduced by JCHARMING, then BIANCA will display the step to reproduce it.

To extract the change patterns and compare them, we used the same technique as the

one presented in sections 4.2.2 and 4.2.2 with new normalizations. The third and the fourth normalizations are removing all *less* important calls in the normalization one and two of RESEMBLE. We classify a call as less-important if, for example, it only does display-related functionalities such as generating HTML or printing something to the console. Finally, the fourth normalization will transform the code to an intermediate language of our own that will allow us to compare source code implemented in different programming languages.

Then, as in RESEMBLE, if the LCS is above a user-defined threshold, then a warning is raised by BIANCA alerting the developer that the commit is suspected to insert a defect. The given defect is shown to the developer and can either force the commit if s/he don't find the warning relevant or abort the commit.

We believe that the warning, alongside the previously mined change patterns and steps to reproduce the suspected default — provided by JCHARMING, if available — will satisfy developers in terms of actionable intelligence. Thus, BIANCA could succeed, where other tools failed, at being used in industrial environment [LLS⁺13].

Early experiments

We have experimented the efficiency of BIANCA with the same datasets we used to build our bug taxonomy proposed in section ??.

Dataset	Fixed Issues	Commit	Files	Projects
Netbeans	53,258	122,632	30,595	39
Apache	49,449	106,366	38,111	349
Total	102,707	229,153	68,809	388

Table 5: Datasets

We choosed to use the same datasets for several reasons. First of all, we spare the time needed to collect new datasets. Then, because these datasets contain a very large system mainly implemented in Java: Netbeans; and 349 independent Apache projects implemented in a very wide range of programming language. Consequently, these datasets allow us to test the efficiency of our different code normalizations.

We ran two different experiments using the two first normalizations we described in section 4.3. Both experiments consider only a few months of history, from April to August 2008. While this could hinder the pertinence of our results, these five months of history contain 167,597 commits related to bug fix. Consequently, we believe our results to be representative.

The first experiment yields the result presented by Figure 24.

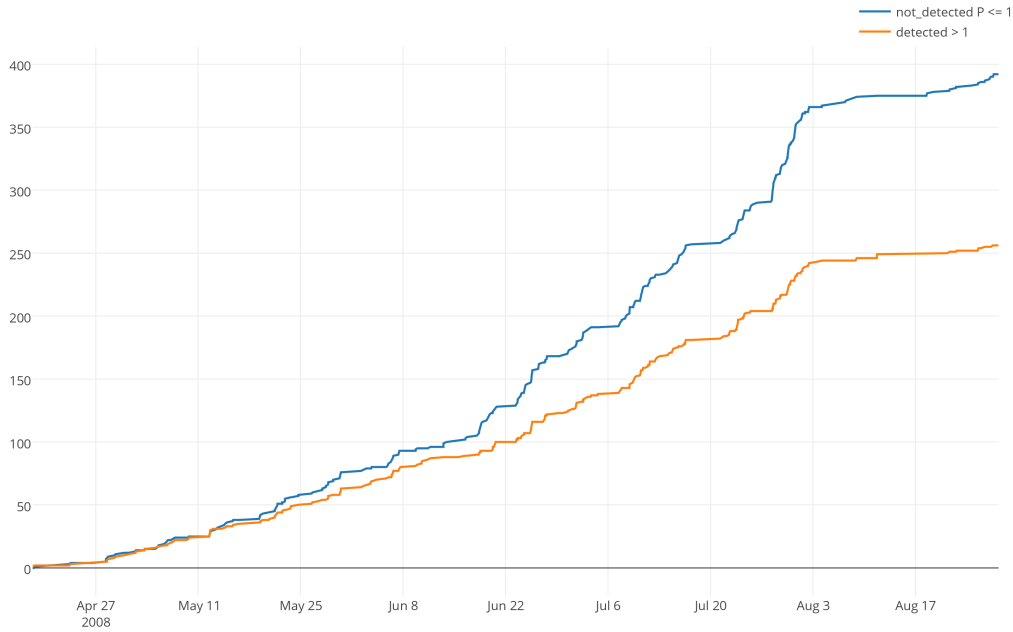


Figure 24: BIANCA warnings from April to August 2008 using the first normalization.

With the first normalization, BIANCA raised 69,519 warnings out of 167,597 (41.5%) analyzed commits. Out of these 69,519, 13.4% turned out to be false positives. A false positive is a commit that have been tagged as introducing a bug by BIANCA but did not according to the history. However, false positives have to be dealt with carefully in this study as the commit might have introduced a bug but the bug could have not been reported yet.

In our second experiment, we used the second normalization and BIANCA raised 83,627 warnings out of 167,597 (48.89%) commit we analyze. However, the false positive rate increases to 21%. Figure 25 shows the results.

Planned experiments

BIANCA experiments are still in their early stage and we are still trying to improve our normalizations in order to reduce the false positive rate. In addition to these improvements, we want to conduct the following additional experiments:

- Full history test with Normalization 1.

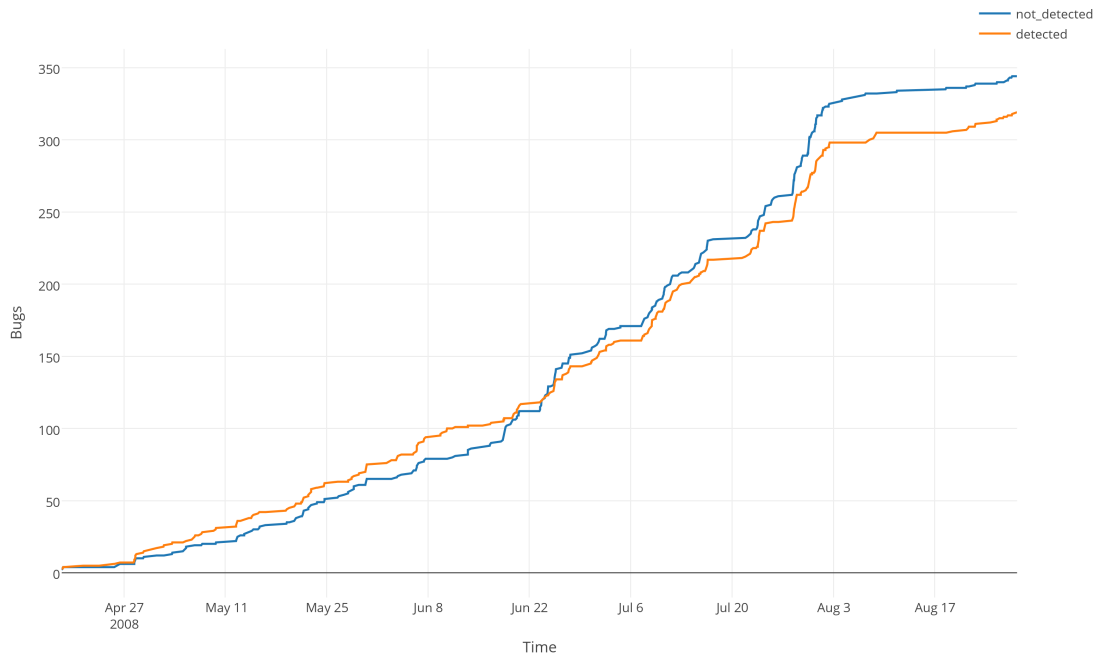


Figure 25: BIANCA warnings from April to August 2008 using the second normalization.

- Full history test with Normalization 2.
- Full history test with Normalization 3.
- Full history test with Normalization 4.
- An human study where:
 - Developers use BIANCA in order to determine whether or not developers take into account our warnings or override them.
 - Rate the proposed solution in a scale from 1 to 10 in order to determine if the proposed change pattern does resolve the current problem.

Chapter 5

A taxonomy to classify the research

In order to classify the research on the different field related to software maintenance, we can reason about types of bugs at different levels. For example, we can group bugs based on the developers that fix them or using information about the bugs such as crash traces. In this paper, we are interested in bugs that share similar fixes. By a fix, we mean a modification (adding or deleting lines of code) to an exiting file that is used to solve the bug. With this in mind, the relationship between bugs and fixes can be modeled using the UML diagram in Figure 26. The diagram only includes bugs that are fixed. From this figure, we can think of four instances of this diagram, which we refer to as bug taxonomy or simply bug types (see Figure 27).



Figure 26: Class diagram showing the relationship between bugs and fixed

The first and second types are the ones we intuitively know about. Type 1 refers to a bug being fixed in one single location (i.e., one file), while Type 2 refers to bugs being fixed in more than one location. In Figure 2, only two locations are shown for the sake of clarity, but many more locations could be involved in the fix of a bug. Type 3 refers to multiple bugs that are fixed in the exact same location. Type 4 is an extension of Type 3, where multiple bugs are resolved by modifying the same set of locations. Note that Type 3 and Type 4 bugs are not duplicates, they may occur when different features of the system fail due to the same root causes (faults). We conjecture that knowing the proportions of each type of bugs in a system may provide insight into the quality of the system. Knowing, for

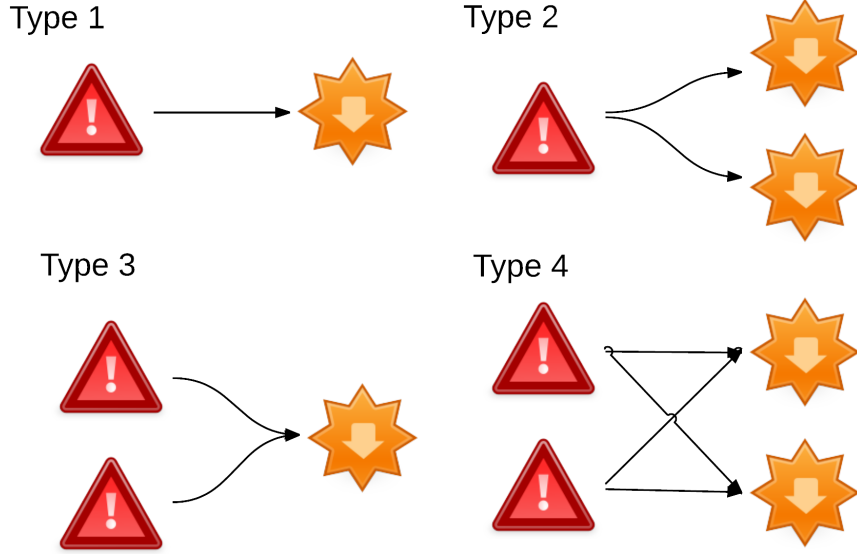


Figure 27: Proposed Taxonomy of Bugs

example, that in a given system the proportion of Type 2 and 4 bugs is high may be an indication of poor system quality since many fixes are needed to address these bugs. In addition, the existence of a high number of Types 3 and 4 bugs calls for techniques that can effectively find bug reports related to an incoming bug during triaging. This is similar to the many studies that exist on detection of duplicates (e.g., [RAN07, SLW⁺10, NNN⁺12]), except that we are not looking for duplicates but for related bugs (bugs that are due to failures of different features of the system, caused by the same faults). To our knowledge, there is no study that empirically examines bug data with these types in mind, which is the main objective of this paper. More particularly, we are interested in the following research questions:

- RQ1: What are the proportions of different types of bugs?
- RQ2: How complex is each type of bugs?
- RQ3: How fast are these types of bugs fixed?

5.0.1 Study Setup

Figure 28 illustrates our data collection and analysis process that we present here and discuss in more detail in the following subsections. First, we extract the raw data from the two bug report management systems used in this study (Bugzilla and Jira). Second, we extract the fix to the bugs from the source code version control system of Netbeans and Apache (Maven and Git).

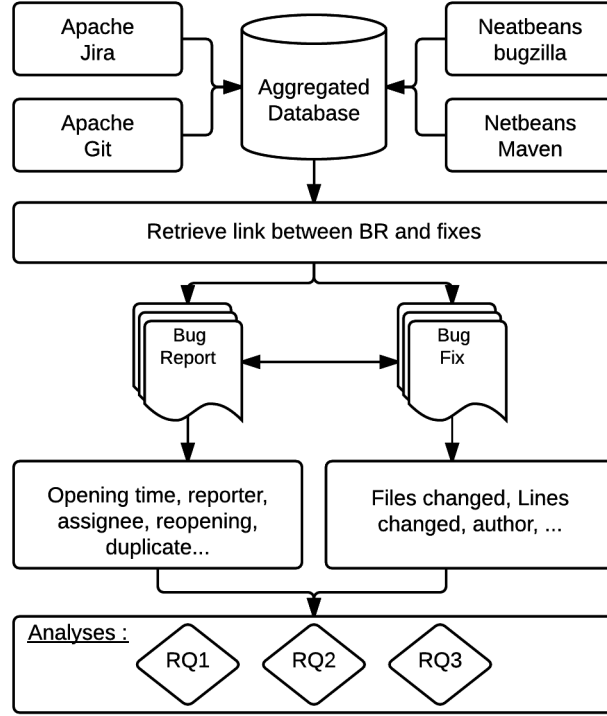


Figure 28: Data collection and analysis process of the study

The extracted data is consolidated in one database where we associate each bug report to its fix. We mine relevant characteristics of BRs and their fixes such as opening time, number of comments, number of times the BR is reopened, number of changesets for BR and the number of files changed and lines modified for fixes or patch. Finally, we analyze these characteristics to answer the aforementioned research questions (RQ).

5.0.2 Datasets

In this study, we used two distinct datasets: Netbeans and the Apache Software Foundation projects. Netbeans is an integrated development environment (IDE) for developing with many languages including Java, PHP, and C/C++. The very first version of Netbeans, then known as Xelfi, appeared in 1996. The Apache Software Foundation is a U.S non-profit organization supporting Apache software projects such as the popular Apache web server since 1999. The characteristics of the Netbeans and Apache Software Foundation datasets used in this paper are presented in Table 6.

Cumulatively, these datasets span from 2001 to 2014. In summary, our consolidated dataset contains 102,707 bugs, 229,153 changesets, 68,809 files that have been modified to fix the bugs, 462,848 comments, and 388 distinct systems. We also collected 221 million

Dataset	R/F BR	CS	Files	Projects
Netbeans	53,258	122,632	30,595	39
Apache	49,449	106,366	38,111	349
Total	102,707	229,153	68,809	388

Table 6: Datasets

lines of code modified to fix the bugs, identified 3,284 sub-projects, and 17,984 unique contributors to these bug report and source code version management systems. Finally, the cumulated opening time for all the bugs reaches 10,661 working years (3,891,618 working days).

5.0.3 Study Design

We describe the design of our study by first stating the research questions, and then explaining the variables, and analysis methods we used to answer these questions. We formulate three research questions (RQs) with the ultimate goal to improve our understanding of each bug type. We focus, however, on Types 2 and 4. This is because these bugs require multiple fixes. They are therefore expected to be more complex. The objective of the first research question is to analyze the proportion of each type of bugs. The remaining two questions address the complexity of the bugs and the bug fixing duration according to the type of bugs. 1)

RQ 1: What are the proportions of different types of bugs?

The answer to this question provides insight into the distribution of bugs according to their type with a focus on Type 2 and 4 bugs. As discussed earlier, knowing, for example, that bugs of Type 2 and 4 are the most predominant ones suggests that we need to investigate techniques to help detect whether an incoming bug is of Types 2 and 4 by examining historical data. Similarly, if we can automatically identify a bug that is related to another one that has been fixed then we can reuse the results of reproducing the first bug in reproducing the second one.

Hypothesis: To answer this question, we analyze whether Type 2 and 4 bugs are predominant in the studied systems, by testing the null hypothesis:

- H_{01A} : The proportion of Types 2 and 4 does not change significantly across the studied systems

We test this hypothesis by observing both a global (across systems) and a local predominance (per system) of the different types of bugs. We must observe these two aspects to ensure that the predominance of a particular type of bug is not circumstantial (in few given systems only) but is also not due to some other, unknown factors (in all systems but not in a particular system).

Variables: We use as variables the amount of resolved/fixed bugs of each type (1, 2, 3 and 4) that are linked to a fix (commit). As mentioned earlier, duplicate bugs are excluded. These are marked as resolved/duplicate in our dataset.

Analysis Method: We answer RQ1 in two steps. The first step is to use descriptive statistics; we compute the ratio of Types 2 and 4 bugs and the ratio of Types 1 and 3 bugs to the total number of bugs in the dataset. This shows the importance of Types 2 and 4 bugs compared to Types 1 and 3 bugs.

In the second step, we compare the proportions of the different types of bugs with respect to the system where the bugs were found. We build the contingency table with these two qualitative variables (the type and studied system) and test the null hypothesis H_{01A} to assess whether the proportion of a particular type of bugs is related to a specific system or not.

We use the Pearson's chi-squared test to reject the null hypothesis H_{01A} . Pearson's chi-squared independence test is used to analyze the relationship between two qualitative data, in our study the type bugs and the studied system. The results of Pearson's chi-squared independence test are considered statistically significant at $\alpha = 0.05$. If $p\text{-value} \leq 0.05$, we reject the null hypothesis H_{01A} and conclude that the proportion of types 3 and 4 bugs is different from the proportion of type 1 and 2 bugs for each system.

RQ 2: How complex is each type of bugs?

We address the relation between Types 2 and 4 bugs and the complexity of the bugs in terms of severity, duplicate and reopened. We analyze whether Types 2 and 4 bugs are more complex to handle than Types 1 and 3 bugs, by testing the null hypotheses:

- H_{02S} : The severity of Types 2 and 4 bugs is not significantly different from the severity of Types 1 and 3
- H_{02D} : Types 2 and 4 bugs are not significantly more likely to get duplicated than Types 1 and 3.
- H_{02R} : Type 2 and 4 bugs are not significantly more likely to get reopened than Types 1 and 3.

Variables: We use as independent variables for the hypotheses H_{02S} , H_{02D} , H_{02R} the bug type (if the bug is from Types 2 and 4 or if it is from Types 1 and 3). For H_{02S} we use the severity as dependent variable to assess the relationship between the bug severity and the bug type. For H_{02D} (respectively H_{02R}) we use a dummy variable duplicated (reopened) to assess if a bug has been duplicated (reopened) at least once or not. This will be used to assess the relationship between the type of the bugs and the fact that the bug is more likely to be reopened or duplicated.

Analysis Method: For each hypothesis, we build a contingency table with the qualitative variables type of bugs (2 and 4 or 1 and 3) and the dependent variable duplicated (respectively reopened) and the severity variable.

We use the Pearsons chi-squared test to reject the null hypothesis H_{02D} (respectively H_{02R}) and H_{02S} . The results of Pearsons chi-squared independence test are considered statistically significant at $\alpha = 0.05$. If a p-value ≤ 0.05 , we reject the null hypothesis H_{02D} (respectively H_{02R}) and conclude the fact that the bug is more likely to be duplicated (respectively reopened) is related to the type of the bug and we reject H_{02S} and conclude that the severity level of the bug is related to the bug type.

RQ 3 : How fast are these types of bugs fixed ?

In this question, we study the relation between the different types of bugs and the fixing time. We are interested in evaluating whether developers take more time to fix Types 2 and 4 bugs than Type 1 and 3, by testing the null hypothesis:

- H_{03} : There is no statistically-significant difference between the duration of fixing periods for Types 2 and 4 bugs and that of Types 1 and 3 bugs.

Variables: To compare the bug fixing time with respect to their type, we use as independent variable the type T_i of a bug B_i , to distinguish between Types 1 and 3 bugs and Types 2 and 4 bugs. We consider as dependent variable the fixing time, FT_i , of the bug B_i . We compute the fixing time FT_i of a bug B_i . The fixing time FT_i is the time between when the bug is submitted to when it is closed/fixed.

Analysis Method: We compute the (non-parametric) Mann- Whitney test to compare the BR fixing time with respect to the BR type and analyze whether the difference in the average fixing time is statistically significant. We use the Mann- Whitney test because, as a non-parametric test, it does not make any assumption on the underlying distributions. We analyze the results of the test to assess the null hypothesis H_{03} . The result is considered as statistically significant at $\alpha = 0.05$. Therefore, if p-value ≤ 0.05 , we reject the null hypothesis H_{03} and conclude that the average fixing time of Types 1 and 3 bugs is significantly

different from the average fixing time of Types 2 and 4 bugs.

5.0.4 Study result and discussion

In this section, we report on the results of the analyses we performed to answer our research questions. We then dedicate a section to discussing the results.

RQ 1 : What are the proportions of different types of bugs?

Figure 29 shows the percentage of the different types of bugs. As shown in the figure, we found that 65% of the bugs are from Types 2 and 4. This shows the predominance of this type of bugs in all the studied systems. Figure 5 shows the repartition per dataset. We can see that Netbeans and Apache have 66% and 64% bugs of Type 1 and 3, respectively. To ensure that this observation is not related to a particular system, we perform Pearsons chi-squared test across the studied systems. Table 7 shows the contingency table for the studied systems and the result of Pearsons chi-squared test. The results show that there is statistically significant difference between the proportions of the different types of bugs.

System	Type 1 and 3	Type 2 and 4	Pearsons chisquared p-value
Apache	4910	8626	p-value <0,0001
Netbeans	9050	17586	

Table 7: Contingency table and Pearson’s chi-squared tests

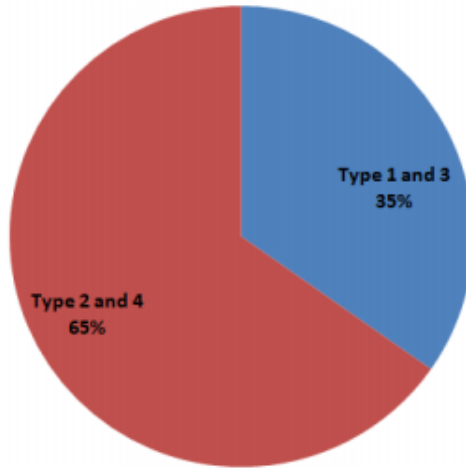


Figure 29: Proportions of different types of bugs

Table 8 shows the number of bugs for each type of bugs and the percentage of each type

of bugs. We can see that Types 3 and 4 bugs represent 28.33% and 61.21% of the total of bugs, respectively. Types 1 and 2 represent only 6.78% and 3.74%. Together, Types 3 and 4 bugs represent almost 90% of the total number of bugs linked to a commit.

Datasets	T1	T2	T3	T4	Total
Netbeans	776 (2.90%)	240 (0.90%)	8372 (31.29%)	17366 (64.91%)	26754
Apache	1968 (14.32%)	1248 (9.08%)	3101 (22.57%)	7422 (54.02%)	13739
Total	2744 (6.78%)	1488 (3.74%)	11473 (28.33%)	24788 (61.21%)	40493

Table 8: Proportion of bug types in amount and percentage

We can thus reject the null hypothesis H_{01A} and conclude that there is a predominance of Types 2 and 4 bugs in all studied systems and this observation is not related to a specific system.

RQ 2 : How complex is each type of bugs?

Figure 30 and 31 show the proportion of each bug type with respect to their severity for each dataset. Table V shows the proportion of each bug type with respect to their severity and dataset. For Netbeans, the bugs we examined in our dataset are either labeled as Blocker or Normal (despite the fact that Netbeans uses Bugzilla that supports all the severity levels presented in the previous section).

For the Apache dataset, the severity levels range from Blocker to Trivial as shown in Figure 30. Figure 31 shows that in Netbeans around 67% of Types 2 and 4 bugs are normal. The same holds for Types 1 and 3 bugs (66% are considered of normal severity). This indicates that most Types 2 and 4 bugs and Types 1 and 3 bugs are not critical in the Netbeans dataset. For the Apache dataset, the results indicate that the majority of the bugs are considered of major severity (66% for Types 1 and 3 and 72% for Types 2 and 4). It is challenging to understand the discrepancy between the two datasets partly because of the way the severity is assigned to BRs.

Table 9 shows the result of the Pearson chi-squared tests for the H_{02S} , H_{02D} and H_{02R} hypotheses.

According to the results of the test (p-value ≤ 0.005), we reject the null hypothesis H_{02S} and conclude that there is a significant difference between the severity of Types 1 and 3 bugs and the severity of Types 2 and 4 bugs.

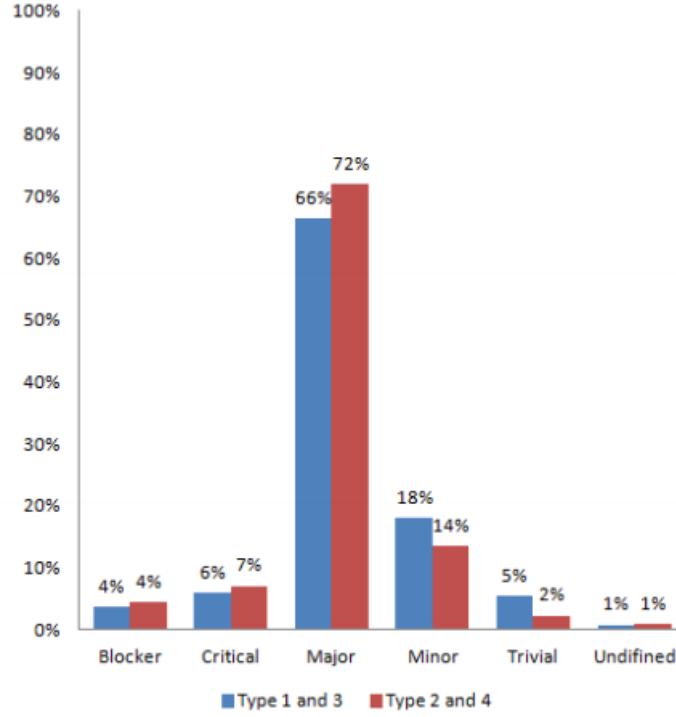


Figure 30: Proportions of Types 1 and 3 versus Types 2 and 4 with respect to their severity in the Apache dataset.

System	Factor	Pearsons chisquared p-value
Apache	Severity	p-value <0.005
	Reopened	p-value <0.005
	Duplicated	p-value <0.005
Netbeans	Severity	p-value <0.005
	Reopened	p-value <0.005
	Duplicated	p-value <0.005

Table 9: Pearson’s chi squared p-values for the severity, the reopen and the duplicate factor with respect to a dataset

Table 11 shows the occurrences of duplicate and reopened bugs with respect to their bug type in each dataset. In Netbeans, the proportion of Type 1 bugs that are marked as source of duplicate is 6.06%, 4.59% for Type 2 bugs, 5.09% for Type 3 bugs and 5.87% for Type 4 bugs with a total of 1503 bugs over 26754 (5.62%). In Apache, the proportion of Type 1 bug marked a source of a duplicate is 2.59% and 2.24%, 1.61% and 2.91% for Types

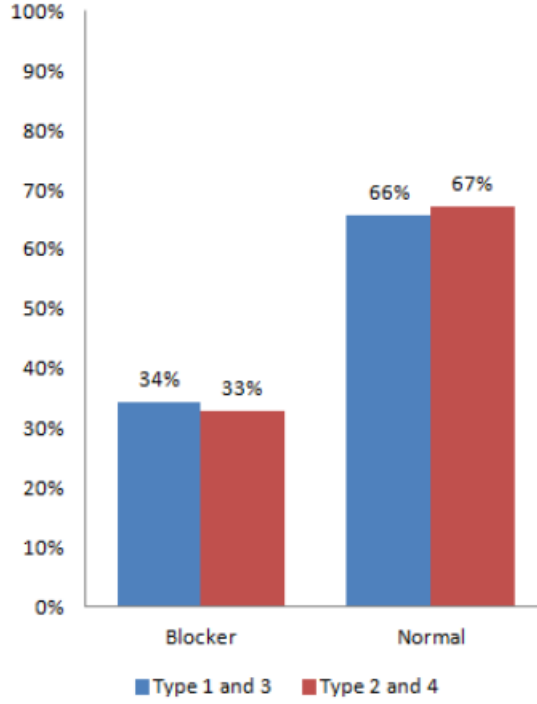


Figure 31: Proportions of Types 1 and 3 versus Types 2 and 4 with respect to their severity in the Netbeans dataset.

2, 3 and 4, respectively.

Second, we analyze the reopened bugs to see the link between the reopening and the type of bugs. We perform Pearsons chi-squared test to reject the null hypothesis H_{02R} .

According to the results of the test (p-value ; 0.005), we reject the null hypothesis H_{02R} and conclude that there is a significant relationship between the reopening of a bug and its type.

Third, we analyze the duplicated bugs to see if there is a link between the bug type and the fact duplication. We perform Pearsons chi-squared test to reject the null hypothesis H_{02D} .

According to the results of the test (p-value ; 0.005), we reject the null hypothesis H_{02D} and conclude that there is a significant relationship between the duplication of a bug and its type.

RQ 3 : How fast are these types of bugs fixed ?

Figure 32 shows the fixing time for Types 1 and 3 versus Types 2 and 4 for Netbeans and the Apache Software Foundation. In Netbeans, 98.96 and 137.05 days are required to fix Types 1 and 3 and Types 2 and 4, respectively. In Apache, 55.76 and 85.48 days are required to

Severity	T1	T2	T3	T4
Netbeans				
Blocker	340 43.81%	109 45.42%	2850 34.04%	5687 32.75%
Normal	436 56.19%	131 54.58%	5522 65.96%	11678 67.25%
Total	776 100%	240 100%	8372 100%	17365 100%
Apache				
Blocker	68 3.46%	53 4.25%	115 3.71%	329 4.43%
Critical	84 4.27%	44 3.53%	213 6.87%	565 7.61%
Major	1245 63.26%	811 64.98%	2096 67.59%	5427 73.12%
Minor	408 20.73%	276 22.12%	501 16.16%	899 12.11%
Trivial	113 5.74%	31 2.48%	159 5.13%	161 2.17%
Total	1918 100%	1215 100%	3084 100%	7381 100%

Table 10: Proportion of each bug type with respect to severity.

fix Types 1 and 3 and Types 2 and 4, respectively.

Table 12 shows the average fixing time of bugs with respect to their bug type in each dataset.

We analyze the difference in the fixing time of bugs with respect to their bug type by conducting a Mann-Whitney test to assess $H03$. The results show that the difference between the fixing time of Types 2 and 4 and Types 1 and 3 is statistically significant (p-value ≤ 0.005).

Therefore, we can reject the null hypothesis $H03$ and conclude that the fixing of Types 2 and 4 bugs takes more time than the fixing of Types 1 and 3 bugs.

Dicussion

Repartition of bug types: One important finding of this study is that there is significantly

Type	T1	T2	T3	T4	Total
Netbeans					
Dup.	6.06% (47)	4.59% (11)	5.09% (426)	5.87% (1019)	5.62% (1503)
Reop.	4.38% (34)	7.08% (17)	4.81% (403)	7.09% (1231)	6.30% (1685)
Apache					
Dup	2.59% (51)	2.24% (28)	1.61% (50)	2.91% (216)	2.51% (345)
Reop	5.59% (110)	6.49% (81)	3.10% (96)	6.90% (512)	5.82% (799)
Total					
Dup	3.57% (98)	2.62% (39)	4.15% (476)	4.98% (1235)	4.56% (1848)
Reop	5.25% (144)	6.59% (98)	4.35% (499)	7.03% (1743)	6.13% (2484)

Table 11: Percentage and occurrences of bugs duplicated by other bufs and reopened with respect to their bug type and dataset.

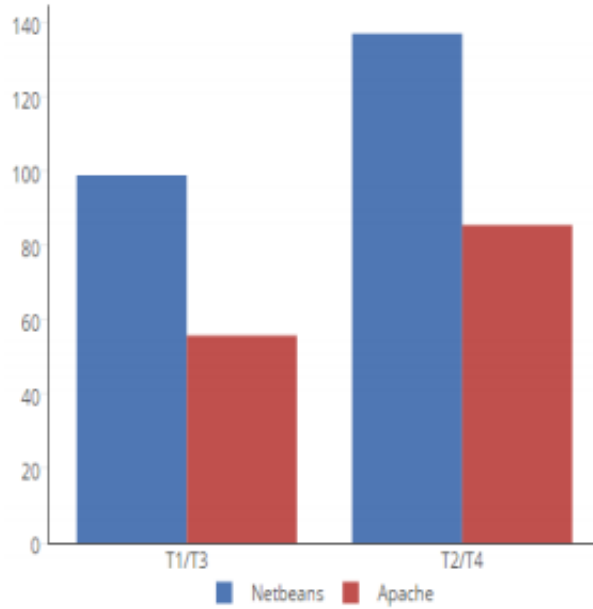


Figure 32: Fixing time of Types 1 and 3 versus fixing time of Types 2 and 4.

Dataset	T1	T2	T3	T4	Average
Netbeans	97.66	117.42	100.26	156.67	118.00
Apache	73.48	118.12	38.04	52.83	70.62
Total	85.57	117.77	69.15	104.75	94.31

Table 12: Average fixing time with respect to bug type

more Types 2 and 4 bugs than Types 1 and 3 in all studied systems. Moreover, this observation is not system-specific. The traditional one-bug/ one-fault way of thinking about bugs only accounts for 35% of the bugs. We believe that, recent triaging algorithms [JW08, JKZ09, KCZH11, TNAKN11a] can benefit from these findings by developing techniques that can detect Type 2 and 4 bugs. This would result in better performance in terms of reducing the cost, time and efforts required by the developers in the bug fixing process.

Severity of bugs: We discussed the severity and the complexity of a bug in terms of its likelihood to be reopened or marked as duplicate (RQ2). Although clear guidelines exist on how to assign the severity of a bug, it remains a manual process done by the bug reporter. In addition, previous studies, notably those by Khomh et al. [KCZH11], showed that severity is not a consistent/trustworthy characteristic of a BR, which lead to the emergence of studies for predicting the severity of bugs (e.g., [LDGG10, LDSV11, TLS12]). Nevertheless, we discovered that there is a significant difference between the severities of Types 1 and 3 compared to Types 2 and 4.

Complexity of bugs: At the complexity level, we use the number of times a bug is reopened as a measure of complexity. Indeed, if a developer is confident enough in his/her fix to close the bug and that the bug gets reopened it means that the developer missed some dependencies of the said bug or did not foresee the consequences of the fix. We found that there is a significant relationship between the number of reopenings and type of a bug. In other words, there is a significant relationship between the complexity and the type of a given bug. In our datasets, Types 1 and 3 bugs are reopened in 1.88% of the cases, while Types 2 and 4 are reopened in 5.73%. Assuming that the reopening is a representative metric for the complexity of bug, Types 2 and 4 are three times more complex than Types 1 and 3. Finally, if we consider multiple reopenings, Types 2 and 4 account for almost 80% of the bugs that reopened more than once and more than 96% of the bug opened more than twice. While current approaches aiming to predict which bug will be reopen use the amount of modified files [SIK⁺10, ZNGM12, Lo13], we believe that they can be improved by taking into account the type of a the bug. For example, if we can detect that an incoming bug is of Type 2 or 4 then it is more likely to be reopened than a bug of Type 1

or 3. Similarly, approaches aiming to predict the files in which a given bug should be fixed could be categorized and improved by knowing the bug type in advance [ZZL12, KTM⁺13].

Impact of a bug: To measure the impact of bugs in end-users and developers, we use the number of times a bug is duplicated. This is because if a bug has many duplicates, it means that a large number of users have experienced and a large number of developers are blocked the failure. We found that there is a significant relationship between the bug type and the fact that it gets duplicated. Types 1 and 3 bugs are duplicated in 1.41% of the cases while Types 2 and 4 are duplicated in 3.14%. Assuming that the amount of duplication is an accurate metric for the impact of bug, Types 2 and 4 have more than two times bigger impact than Types 1 and 3. Similarly to reopening, if we consider multiple duplication, Types 2 and 4 account for 75% of the bugs that get duplicated more than once and more than 80% of the bugs that get duplicated more than twice. We believe that approaches targeting the identification of duplicates [BPZ08, JW08, SLW⁺10, TSL12] could leverage this taxonomy to achieve even better performances in terms of recall and precision.

Fixing time: Our third research question aimed to determine if the type of a bug impacts its fixing time. Not only we found that the type of a bug does significantly impact its fixing time, but we also found that, in average Types 2 and 4, stay open 111.26 days while Types 1 and 3 last for 77.36 days. Types 2 and 4 are 1.4 time longer to fix than Types 1 and 3. We therefore believe that, approaches aiming to predict the fixing time of a bug (e.g., [Pan07, BN11, ZGV13]) can highly benefit from accurately predicting the type of a bug and therefore better plan the required man-power to fix the bug. In summary, Types 2 and 4 account for 65% of the bugs and they are more complex, have a bigger impact and take longer to be fixed than Types 1 and 3 while being equivalent in terms of severity.

Our taxonomy aimed to analyse: (1) the proportion of each type of bugs; (2) the complexity of each type in terms of severity, reopening and duplication; (3) the required time to fix a bug depending on its type. The key findings are:

- Types 2 and 4 account for 65% of the bugs.
- Types 2 and 4 have a similar severity compared to Types 1 and 3.
- Types 2 and 4 are more complex (reopening) and have a bigger impact (duplicate) than Types 1 and 3.
- It takes more time to fix Types 2 and 4 than Types 1 and 3.

Our taxonomy and results can be built upon in order to classify past and new researches in several active areas such as bug reproduction and triaging, prediction of reopening,

duplication, severity, files to fix and fixing time. Moreover, if one could predict the type of a bug at submission time, all these areas could be improved.

In order to ease all these future works we built **BUMPER** that we present in the next section.

Chapter 6

Reflection on pragmatic software maintenance

Architects, the ones that design buildings — where mistakes cost lives — spend at least five years at school and possibly their whole careers to study, understand and reproduce great designs made by great architects. Software architects, however, begin in programming 101 by displaying the famous “Hello World” statement and exponentially increase the complexity of their programs over their years of study and work. At some point, they will earn the title of software architect (or technical leader) because they have designed, maintained and evolved *enough* programs to be trustworthy on the matter. However, unlike building architects, they have to learn how to recognize, analyze and reproduce great architectural choices by themselves in addition of their day to day work. Of course, software developers do learn good practices such as design patterns [GHJV08] but in a very few occasions they will be presented with a state-of-the-art program built by great developers (Amy Brown *et al.* propose exactly that in their books [AW12, BG12, Arm13]).

While our research is not about reforming how programming classes are taught, we still want to ease the access to this knowledge for developers during their programming sessions in order to ship better programs.

We shift the focus from merely mining revision and issue management system, where knowledge of great developers lies, to integrate them in their rightful place: as the keystone of software development and evolution activities. Extracting the ground truth from repositories helped engineers and practitioners to be better at building softwares as they know, for example, *how long it will take to fix a bug* [WPZZ07], *what makes a good bug report* [BJS⁺08] or *how to fix long-lived bugs* [SKP14]. Using these discoveries, tools can be created, on a

per organization basis, to fit particular requirements such as programming languages, development processes or particular threshold. If we want to truthfully and deeply modify the software engineering landscape to have better softwares in terms of quality, maintainability and availability, we need to provide this information during the development, maintenance and evolution processes according to a specific context in an easy, reliable, actionable, free way.

If we look back at the history of software engineering, the increase of processors' speed and decrease of their price allowed one to have a compiler on its own machine rather than sending one's code to the mainframe and receive compilation errors hours (days) later. This allowed, among other factors, the democratization of software engineering as *everyone*, belonging to a major organization or not, became able to build code. We believe that, it is now time to allow developers, engineering and practitioners, regardless of their programming language and contextual environment, not only to write and build code but to write and built qualitative, robust, resilient, easy to maintain and to fix code. What better way to do so than to *stand on the shoulder of giants* by having access to all the open sources repositories, including but not limited to, issues, tasks, bug fixes, patches, comments, good practices break down to the right level and provided at the right time during day to day programming sessions?

We believe that the *right* times for distilling this information are when a developer version its code.

Chapter 7

Research Plan

In this chapter, we present the current state of research. Section 7.2 presents our current contributions and section 7.4 presents our publication plan.

7.1 Current State of Research

Chapter ?? presented the current state of our research. The completed parts are :

- To create a taxonomy of bug (section ??).
- To make software artifacts such as source code, issues, tasks, comments and more searchable using natural language and not structured query language (section 3.1).
- To efficiently reproduce field crashes in lab environment without raising any privacy concerns (section 3.2).
- To sequentially mine co-changes (section 4.2.2).
- To normalize source code (sections 4.2.2 and ??).
- To efficiently compare normalized source code (sections 4.2.2 and ??).
- To raise warning when a normalized source code matches another normalized source code known to have introduced a defect in an application (section ??).

7.2 Current Contribution

We already have one contribution referred in an international journal [NHLTL16]:

- Nayrolles, M. , Hamou-Lhadj, W., Tahar, S. Larsson, A. (2016). A Bug Reproduction Approach Based on Directed Model Checking and Crash Traces. *Journal of Software: Evolution and Process*. Wiley. (Accepted with revisions).

and three contributions referred in international conferences [NHL15, MHLN⁺15, NHL16b, NHL16a]:

- Nayrolles, M. & Hamou-Lhadj, W. (2016). BUMPER: A Tool to Cope with Natural Language Search of Millions Bugs and Fixes. *International Conference on Software Analysis, Evolution, and Reengineering - Tool Track (SANER'16)*. IEEE. (Accepted).
- Mathieu Nayrolles, Abdelwahab Hamou-Lhadj, Tahar Sofiene, and Alf Larsson. JCHARMING : A Bug Reproduction Approach Using Crash Traces and Directed Model Checking. In *SANER'15*, pages 101-110, 2015. (Best Paper Award).
- Maiga, A., Hamou-Lhadj, W., Nayrolles, M. , Sabor, K. & Larsson, A. (2015). An Empirical Study on the Handling of Crash Reports in a Large Software Company: An Experience Report. *International Conference on Software Maintenance and Evolution (ICSME)*. (pp. 342-351). IEEE.

Finally, we also presented BUMPER during the 2015 Consortium for Software Engineering Research Fall Edition [NHL15]:

- Nayrolles, M. & Hamou-Lhadj, W. (2015). BUMPER: Bug Metarepository Search Engine for Developers and Researchers. *Consortium for Software Engineering Research Fall*.

7.3 Plan for short term work

As short term work, we want to complete the following parts:

- Recommend code in order to improve a sub-optimum code.
- Recommend code in order to complete an incomplete piece of code.
- Develop IDE plugins
- Develop new parsers (Github, Sourceforge) in order to improve:
 - The coverage of BUMPER.
 - The precision of RESEMBLE and BIANCA.
 - The accuracy of our Bug Taxonomy.

7.4 Publication Plan

This section presents our planned publications over the course of the next year. Depending on time constraints and possible invitations to journal extension, publication can be removed or added during the next years.

- **Publication 6.** BIANCA will be submitted to International Symposium on Software Reliability Engineering, ISSRE 2016.
- **Publication 7.** PRECINCT will be submitted to International Working Conference on Source Code Analysis and Manipulation, SCAM 2016.
- **Publication 8.** RESEMBLE will be submitted to International Conference Software Maintenance and Evolution, ICSME 2017.
- **Publication 9.** Our proposed bug taxonomy will be submitted to Transaction of Software Engineering, TSE 2017.
- **Publication 10.** Our ecosystem as a whole, BUMPER, JCHARMING, RESEMBLE and BIANCA will be submitted to Transaction of Software Engineering, TSE 2017.
- **Thesis.** In parallel to publications 10 and 11, I plan to write my Ph.D thesis.

Chapter 8

Conclusion

The maintenance and evolution of complex software systems account for more than 70% software’s life cycle. Hundreds of papers have been published with the aim to improve our knowledge of these processes in terms of issue triaging, issue prediction, duplicate issue detection, issue reproduction and co-changes prediction. All these publications gave meaning to the millions of issues that can be found in open source issue & project and revision management systems. Context-aware IDE and think tank in open source architecture ([AW12]) open the path to approaches that support developers during their programming sessions by leveraging past indexed knowledge and past architectures.

In this research proposal, we first presented the most influential papers in the different fields our work lies on in Chapter 2. Chapter ?? presented our proposal in details while chapter 7 detailed our attempt planning.

More specifically, in Chapter ??, we presented four approaches: **BUMPER**, **JCHARMING**, **RESEMBLE**, **BIANCA**. Also, we proposed a taxonomy of bugs. When combined into **pErICOPE** (Ecosystem Improve source COde during Programming session with real-time mining of common knowlEdge), these tools (i) provide the possibility to search related software artifacts using natural language, (ii) accurately reproduce field-crash in lab environment, (iii) recommend improvement or completion of current block of code and (iv) prevent the introduction of clones / issues at commit time.

BUMPER has been designed to handle heavy traffic while **JCHARMING** can reproduce 85% of real-world issues we submitted to it. On its side **BIANCA** is able to flag 41.5% and 48.89% of commit introducing a bug as dangerous with 13.4% and 21% of false positive using two code different code normalization, respectively.

Our future works, according to our publication plan described in section 7.4, are as follows. First, we want to improve the performances of **BIANCA** in terms of false positives.

Then, create the IDE plugin that will support RESEMBLE. Finally, we want to refine our taxonomy by including as many as datasets as possible.

Bibliography

- [ACC⁺02] G. Antoniol, G. Canfora, G. Casazza, A. De Lucia, and E. Merlo. Recovering traceability links between code and documentation. *IEEE Transactions on Software Engineering*, 28(10):970–983, oct 2002.
- [AHM06] John Anvik, Lyndon Hiew, and Gail C. Murphy. Who should fix this bug? In *Proceeding of the 28th international conference on Software engineering - ICSE '06*, page 361, New York, New York, USA, may 2006. ACM Press.
- [AKE08] Shay Artzi, Sunghun Kim, and Michael D Ernst. Recrash: Making software failures reproducible by preserving object states. In *Proceedings of the 22nd European Conference on Object-Oriented Programming*, pages 542–565, 2008.
- [ALRL04] A. Avizienis, J.-C. Laprie, B. Randell, and C. Landwehr. Basic concepts and taxonomy of dependable and secure computing. *IEEE Transactions on Dependable and Secure Computing*, 1(1):11–33, jan 2004.
- [AP08] Nathaniel Ayewah and William Pugh. A report on a survey and study of static analysis users. In *Proceedings of the 2008 workshop on Defects in large software systems - DEFECTS '08*, page 1, New York, New York, USA, jul 2008. ACM Press.
- [Apa] Apache Software Foundation. Apache Ant.
- [Apa00] Apache Software Foundation. Apache Struts Project, 2000.
- [Apa14] Apache Software Foundation. Apache PDFBox — A Java PDF Library, 2014.
- [APM⁺07] Nathaniel Ayewah, William Pugh, J. David Morgenthaler, John Penix, and YuQian Zhou. Evaluating static analysis defect warnings on production software. In *Proceedings of the 7th ACM SIGPLAN-SIGSOFT workshop on*

- Program analysis for software tools and engineering - PASTE '07*, pages 1–8, New York, New York, USA, jun 2007. ACM Press.
- [Arm13] Tavish Armstrong. *The Performance of Open Source Applications*. 2013.
- [AW12] Brown Amy and Greg Wilson. *The Architecture of Open Source Applications*. CreativeCommons, 2012.
- [Bak] B.S. Baker. On finding duplication and near-duplication in large software systems. In *Proceedings of 2nd Working Conference on Reverse Engineering*, pages 86–95. IEEE Comput. Soc. Press.
- [Bak92] Brenda S. Baker. A program for identifying duplicated code. *Computing Science and Statistics*, 1992.
- [BBM96] V.R. Basili, L.C. Briand, and W.L. Melo. A validation of object-oriented design metrics as quality indicators. *IEEE Transactions on Software Engineering*, 22(10):751–761, 1996.
- [BBR⁺10] Adrian Bachmann, Christian Bird, Foyzur Rahman, Premkumar Devanbu, and Abraham Bernstein. The missing links. In *Proceedings of the eighteenth ACM SIGSOFT international symposium on Foundations of software engineering - FSE '10*, page 97, New York, New York, USA, nov 2010. ACM Press.
- [BDW99] L.C. Briand, J.W. Daly, and J.K. Wust. A unified framework for coupling measurement in object-oriented systems. *IEEE Transactions on Software Engineering*, 25(1):91–121, 1999.
- [BG02] Brenda S. Baker and Raffaele Giancarlo. Sparse Dynamic Programming for Longest Common Subsequence from Fragments. *Journal of Algorithms*, 42(2):231–254, feb 2002.
- [BG12] Amy Brown and Wilson Greg. *The Architecture of Open Source Applications, Volume II*. 2012.
- [BJS⁺08] Nicolas Bettenburg, Sascha Just, Adrian Schröter, Cathrin Weiss, Rahul Premraj, and Thomas Zimmermann. What makes a good bug report? In *Proceedings of the 16th ACM SIGSOFT International Symposium on Foundations of Software Engineering*, page 308, New York, New York, USA, 2008. ACM Press.

- [BKA⁺07] Stefan Bellon, Rainer Koschke, Giulio Antoniol, Jens Krinke, and Ettore Merlo. Comparison and Evaluation of Clone Detection Tools. *IEEE Transactions on Software Engineering*, 33(9):577–591, sep 2007.
- [BMD⁺] M. Balazinska, E. Merlo, M. Dagenais, B. Lague, and K. Kontogiannis. Partial redesign of Java software systems based on clone analysis. In *Sixth Working Conference on Reverse Engineering (Cat. No.PR00303)*, pages 326–336. IEEE Comput. Soc.
- [BMD⁺99] M. Balazinska, E. Merlo, M. Dagenais, B. Lague, and K. Kontogiannis. Measuring clone based reengineering opportunities. In *Proceedings Sixth International Software Metrics Symposium (Cat. No.PR00403)*, pages 292–303. IEEE Comput. Soc., 1999.
- [BN11] Pamela Bhattacharya and Iulian Neamtiu. Bug-fix time prediction models. In *Proceeding of the 8th working conference on Mining software repositories - MSR '11*, page 207, New York, New York, USA, may 2011. ACM Press.
- [BPZ08] Nicolas Bettenburg, Rahul Premraj, and Thomas Zimmermann. Duplicate bug reports considered harmful ... really? In *2008 IEEE International Conference on Software Maintenance*, pages 337–345. IEEE, 2008.
- [Bur03] Oliver Burn. Checkstyle, 2003.
- [Bur06] Ilene Burnstein. *Practical Software Testing: A Process-Oriented Approach*. Springer Science & Business Media, 2006.
- [BvdH13] Gerald Bortis and Andre van der Hoek. PorchLight: A tag-based approach to bug triaging. In *2013 35th International Conference on Software Engineering (ICSE)*, pages 342–351. IEEE, may 2013.
- [BYM⁺98] Ira D. Baxter, Andrew Yahin, Leonardo Moura, Marcelo Sant’Anna, and Lorraine Bier. Clone Detection Using Abstract Syntax Trees. In *Proceedings of the IEEE International Conference on Software Maintenance.*, pages 368–377. IEEE Computer Society, mar 1998.
- [CFS09] Satish Chandra, Stephen J Fink, and Manu Sridharan. Snugglebug: a powerful approach to weakest preconditions. In *ACM Sigplan Notices*, volume 44, pages 363–374. ACM, 2009.

- [CHA09] CHANCHAL K. ROY. *Detection and Analysis of Near-Miss Software Clones*. PhD thesis, Queen's University, 2009.
- [Che13] Ning Chen. *Star: stack trace based automatic crash reproduction*. PhD thesis, The Hong Kong University of Science and Technology, 2013.
- [CK94] S.R. Chidamber and C.F. Kemerer. A metrics suite for object oriented design. *IEEE Transactions on Software Engineering*, 20(6):476–493, jun 1994.
- [CNSH14] Tse-hsun Chen, Meiyappan Nagappan, Emad Shihab, and Ahmed E Hassan. An Empirical Study of Dormant Bugs Categories and Subject Descriptors. In *Mining Software Repository*, pages 82–91, 2014.
- [CO07] James Clause and Alessandro Orso. A Technique for Enabling and Supporting Debugging of Field Failures. In *Proceedings of the 29th International Conference on Software Engineering*, pages 261–270, 2007.
- [Col] CollabNet. Tigris.org: Open Source Software Engineering.
- [Cor] Michael W. Godfrey Cory Kapser. Toward a Taxonomy of Clones in Source Code: A Case Study.
- [Cor06a] James R. Cordy. Source transformation, analysis and generation in TXL. In *Proceedings of the 2006 ACM SIGPLAN symposium on Partial evaluation and semantics-based program manipulation - PEPM '06*, page 1, New York, New York, USA, jan 2006. ACM Press.
- [Cor06b] James R. Cordy. The TXL source transformation language. *Science of Computer Programming*, 61(3):190–210, aug 2006.
- [CR11] James R. Cordy and Chanchal K. Roy. The NiCad Clone Detector. In *2011 IEEE 19th International Conference on Program Comprehension*, pages 219–220. IEEE, jun 2011.
- [Dan00] Andreas Dangel. PMD, 2000.
- [DCMS] Thomas R. Dean, James R. Cordy, Andrew J. Malton, and Kevin A. Schneider. Agile Parsing in TXL. In *Proceedings of IEEE International Conference on Automated Software Engineering*, volume 10, pages 311–336. Kluwer Academic Publishers.

- [De 01] Andrea De Lucia. Program slicing: Methods and applications. In *International Working Conference on Source Code Analysis and Manipulation*, page 144. IEEE Computer Society, 2001.
- [DER07] Ekwa Duala-Ekoko and Martin P. Robillard. Tracking Code Clones in Evolving Software. In *29th International Conference on Software Engineering*, pages 158–167. IEEE, may 2007.
- [DER10] Ekwa Duala-Ekoko and Martin P. Robillard. Clone region descriptors. *ACM Transactions on Software Engineering and Methodology*, 20(1):1–31, jun 2010.
- [DM06] Bruno Dutertre and Leonardo De Moura. The yices smt solver. *Tool paper at <http://yices.csl.sri.com/tool-paper.pdf>*, 2(2), 2006.
- [DMT13] Anthony Demange, Naouel Moha, and Guy Tremblay. Detection of SOA Patterns. In *International Conference on Service-Oriented Computing*, pages 114–130, 2013.
- [DRD99] S. Ducasse, M. Rieger, and S. Demeyer. A language independent approach for detecting duplicated code. In *Proceedings IEEE International Conference on Software Maintenance - 1999 (ICSM'99). 'Software Maintenance for Business Change' (Cat. No.99CB36360)*, pages 109–118. IEEE, 1999.
- [DZM] Valentin Dallmeier, Andreas Zeller, and Bertrand Meyer. Generating Fixes from Object Behavior Anomalies.
- [EMM01] Khaled El Emam, Walcelio Melo, and Javam C. Machado. The prediction of faulty classes using object-oriented design metrics. *Journal of Systems and Software*, 56(1):63–75, feb 2001.
- [FFK08] Raimar Falke, Pierre Frenzel, and Rainer Koschke. Empirical evaluation of clone detection using syntax suffix trees. *Empirical Software Engineering*, 13(6):601–643, jul 2008.
- [FM15] Sylvie L. Foss and Gail C. Murphy. Do developers respond to code stability warnings? pages 162–170, nov 2015.
- [G. 91] G. Piatetsky-Shapiro. Discovery, analysis and presentation of strong rules. *Knowledge Discovery in Databases*, pages 229–249, jan 1991.
- [Gav13] Andrej Gavrilov. Find bugs faster using stack traces. Technical report, 2013.

- [GFS05] T. Gyimothy, R. Ferenc, and I. Siket. Empirical validation of object-oriented metrics on open source software for fault prediction. *IEEE Transactions on Software Engineering*, 31(10):897–910, oct 2005.
- [GHH⁺09] Michael W. Godfrey, Ahmed E. Hassan, James Herbsleb, Gail C. Murphy, Martin Robillard, Prem Devanbu, Audris Mockus, Dewayne E. Perry, and David Notkin. Future of Mining Software Archives: A Roundtable. *IEEE Software*, 26(1):67–70, jan 2009.
- [GHJV08] Eric Gamma, Richard Helm, Ralph Johnson, and John Vlissides. *Design patterns*. Addison Wesley, 2008.
- [GJS08] Mark Gabel, Lingxiao Jiang, and Zhendong Su. Scalable detection of semantic clones. In *Proceedings of the 13th international conference on Software engineering - ICSE '08*, page 321, New York, New York, USA, 2008. ACM Press.
- [Has08] Ahmed E. Hassan. The road ahead for Mining Software Repositories. In *2008 Frontiers of Software Maintenance*, pages 48–57. IEEE, sep 2008.
- [Has09] Ahmed E. Hassan. Predicting faults using the complexity of code changes. In *2009 IEEE 31st International Conference on Software Engineering*, pages 78–88. IEEE, may 2009.
- [HGGBR08] Israel Herraiz, Daniel M. German, Jesus M. Gonzalez-Barahona, and Gregorio Robles. Towards a simplification of the bug report form in eclipse. In *Proceedings of the 2008 international workshop on Mining software repositories - MSR '08*, page 145, New York, New York, USA, may 2008. ACM Press.
- [HH05] A.E. Hassan and R.C. Holt. The top ten list: dynamic fault prediction. In *21st IEEE International Conference on Software Maintenance (ICSM'05)*, pages 263–272. IEEE, 2005.
- [HHA97] MANNILA HEIKKI, TOIVONEN HANNU, and VERKAMO A. INKERI. Discovery of Frequent Episodes in Event Sequences. *Data Mining and Knowledge Discovery*, 289:259–289, 1997.
- [Hir77] Daniel S Hirschberg. Algorithms for the longest common subsequence problem. *Journal of the ACM (JACM)*, 24(4):664–675, 1977.

- [HKKI04] Yoshiki Higo, Toshihiro Kamiya, Shinji Kusumoto, and Katsuro Inoue. Refactoring support based on code clone analysis. In *Product Focused Software Process Improvement*, pages 220–233. Springer, 2004.
- [HNH15] H Hemmati, M Nagappan, and Ae Hassan. Investigating the effect of defect co-fix on quality assurance resource allocation: A search-based approach. *Journal of Systems and Software*, 00:1–18, 2015.
- [Hol97] Gerard J Holzmann. The model checker SPIN. *IEEE Transactions on Software Engineering*, 23(5):279–295, 1997.
- [Hov07] David Hovemeyer. FindBugs, 2007.
- [HP04] David Hovemeyer and William Pugh. Finding bugs is easy. *ACM SIGPLAN Notices*, 39(12):92, dec 2004.
- [HR02] Health, Social and Economics Research. The Economic Impacts of Inadequate Infrastructure for Software Testing. Technical report, 2002.
- [HS77] James W. Hunt and Thomas G. Szymanski. A fast algorithm for computing longest common subsequences. *Communications of the ACM*, 20(5):350–353, may 1977.
- [JDHW09] Elmar Juergens, Florian Deissenboeck, Benjamin Hummel, and Stefan Wagner. Do code clones matter? In *2009 IEEE 31st International Conference on Software Engineering*, pages 485–495. IEEE, may 2009.
- [JFG09] Dennis Jeffrey, Min Feng, and Rajiv Gupta. BugFix: A learning-based tool to assist developers in fixing bugs. In *2009 IEEE 17th International Conference on Program Comprehension*, pages 70–79. IEEE, may 2009.
- [JKXC10] Hojun Jaygarl, Sunghun Kim, Tao Xie, and Carl K Chang. OCAT: Object Capture based Automated Testing. In *Proceedings of the 19th International Symposium on Software Testing and Analysis*, pages 159–170, 2010.
- [JKZ09] Gaeul Jeong, Sunghun Kim, and Thomas Zimmermann. Improving bug triage with bug tossing graphs. In *Proceedings of the the 7th Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on The Foundations of Software Engineering*, page 111, New York, New York, USA, aug 2009. ACM Press.

- [JLL⁺12] Shujuan Jiang, Wei Li, Haiyang Li, Yanmei Zhang, Hongchang Zhang, and Yingqi Liu. Fault Localization for Null Pointer Exception Based on Stack Trace and Program Slicing. *2012 12th International Conference on Quality Software*, pages 9–12, aug 2012.
- [JMSG07] Lingxiao Jiang, Ghassan Mishherghi, Zhendong Su, and Stephane Glondu. DECKARD: Scalable and Accurate Tree-Based Detection of Code Clones. In *29th International Conference on Software Engineering*, pages 96–105. IEEE, may 2007.
- [JO12] Wei Jin and Alessandro Orso. BugRedux: Reproducing field failures for in-house debugging. In *Proceedings of the 34th International Conference on Software Engineering, IEEE*, pages 474–484. Ieee, jun 2012.
- [JO13] Wei Jin and Alessandro Orso. F3: fault localization for field failures. In *Proceedings of the 2013 International Symposium on Software Testing and Analysis*, pages 213–223, New York, New York, USA, 2013. ACM Press.
- [Joh93] J Howard Johnson. Identifying redundancy in source code using fingerprints. In *CASCON '93 Proceedings of the 1993 conference of the Centre for Advanced Studies on Collaborative research: software engineering - Volume 1*, pages 171–183. IBM Press, oct 1993.
- [Joh94] J. Howard Johnson. Visualizing textual redundancy in legacy source. In *Proceedings of the 1994 conference of the Centre for Advanced Studies on Collaborative research*, page 32. IBM Press, oct 1994.
- [JSMHB13] Brittany Johnson, Yoonki Song, Emerson Murphy-Hill, and Robert Bowdidge. Why don't software developers use static analysis tools to find bugs? In *35th International Conference on Software Engineering (ICSE)*, pages 672–681. IEEE Press, may 2013.
- [JW08] Nicholas Jalbert and Westley Weimer. Automated duplicate detection for bug tracking systems. In *2008 IEEE International Conference on Dependable Systems and Networks With FTCS and DCC (DSN)*, pages 52–61. IEEE, 2008.
- [KCZH11] Foutse Khomh, Brian Chan, Ying Zou, and Ahmed E. Hassan. An Entropy Evaluation Approach for Triaging Field Crashes: A Case Study of Mozilla

- Firefox. In *2011 18th Working Conference on Reverse Engineering*, pages 261–270. IEEE, oct 2011.
- [KFF06] Rainer Koschke, Raimar Falke, and Pierre Frenzel. Clone Detection Using Abstract Syntax Suffix Trees. In *2006 13th Working Conference on Reverse Engineering*, pages 253–262. IEEE, oct 2006.
- [KG] C. Kapser and M.W. Godfrey. Aiding comprehension of cloning through categorization. In *Proceedings. 7th International Workshop on Principles of Software Evolution, 2004.*, pages 85–94. IEEE.
- [KG06] Cory Kapser and Michael Godfrey. Cloning Considered Harmful Considered Harmful. In *2006 13th Working Conference on Reverse Engineering*, pages 19–28. IEEE, oct 2006.
- [KH] R. Komondoor and S. Horwitz. Effective, automatic procedure extraction. In *MHS2003. Proceedings of 2003 International Symposium on Micromechanics and Human Science (IEEE Cat. No.03TH8717)*, pages 33–42. IEEE Comput. Soc.
- [KH00] Raghavan Komondoor and Susan Horwitz. Semantics-preserving procedure extraction. In *Proceedings of the 27th ACM SIGPLAN-SIGACT symposium on Principles of programming languages - POPL '00*, pages 155–169, New York, New York, USA, jan 2000. ACM Press.
- [KKI02a] T. Kamiya, S. Kusumoto, and K. Inoue. CCFinder: a multilinguistic token-based code clone detection system for large scale source code. *IEEE Transactions on Software Engineering*, 28(7):654–670, jul 2002.
- [KKI02b] T. Kamiya, S. Kusumoto, and K. Inoue. CCFinder: a multilinguistic token-based code clone detection system for large scale source code. *IEEE Transactions on Software Engineering*, 28(7):654–670, jul 2002.
- [Kon] K. Kontogiannis. Evaluation experiments on the detection of programming patterns using software metrics. In *Proceedings of the Fourth Working Conference on Reverse Engineering*, pages 44–54. IEEE Comput. Soc.
- [Kri63] Saul A. Kripke. Semantical Considerations on Modal Logic. *Acta Philosophica Fennica*, 16(1963):83–94, 1963.

- [Kri01] Jens Krinke. Identifying Similar Code with Program Dependence Graphs. In *Eighth IEEE Working Conference on Reverse Engineering, 2001*, pages 301–309. IEEE Computer Society, oct 2001.
- [Kro99] Thomas Kropf. *Introduction to formal hardware verification*. Springer, 1999.
- [KSN05] Miryung Kim, Vibha Sazawal, and David Notkin. An empirical study of code clone genealogies. *ACM SIGSOFT Software Engineering Notes*, 30(5):187, sep 2005.
- [KTM⁺13] Dongsun Kim, Yida Tao, Student Member, Sunghun Kim, and Andreas Zeller. Where Should We Fix This Bug? A Two-Phase Recommendation Model. *Transaction on Software Engineering*, 39(11):1597–1610, 2013.
- [KWM⁺11] Dongsun Kim, Xinming Wang, Student Member, Sunghun Kim, Andreas Zeller, S C Cheung, Senior Member, and Sooyong Park. Which Crashes Should I Fix First?: Predicting Top Crashes at an Early Stage to Prioritize Debugging Efforts. *TRANSACTIONS ON SOFTWARE ENGINEERING*, 37(3):430–447, 2011.
- [KZN13] Sunghun Kim, Thomas Zimmermann, and Nachiappan Nagappan. Crash Graphs: An Aggregated View of Multiple Crashes to Improve Crash Triage. In *International Conference on Dependable Systems and Networks (DSN)*, pages 486–493, 2013.
- [KZPJ06] Sunghun Kim, Thomas Zimmermann, Kai Pan, and E. Jr. Whitehead. Automatic Identification of Bug-Introducing Changes. In *21st IEEE/ACM International Conference on Automated Software Engineering (ASE’06)*, pages 81–90. IEEE, 2006.
- [KZWG11] Sunghun Kim, Hongyu Zhang, Rongxin Wu, and Liang Gong. Dealing with noise in defect prediction. *Proceeding of the 33rd international conference on Software engineering - ICSE ’11*, page 481, 2011.
- [KZWZ07a] Sunghun Kim, Thomas Zimmermann, E. James Whitehead Jr., and Andreas Zeller. Predicting Faults from Cached History. In *29th International Conference on Software Engineering*, pages 489–498. IEEE, may 2007.
- [KZWZ07b] Sunghun Kim, Thomas Zimmermann, E. James Whitehead Jr., and Andreas Zeller. Predicting Faults from Cached History. In *29th International Conference on Software Engineering (ICSE’07)*, pages 489–498. IEEE, may 2007.

- [LDGG10] Ahmed Lamkanfi, Serge Demeyer, Emanuel Giger, and Bart Goethals. Predicting the severity of a reported bug. In *2010 7th IEEE Working Conference on Mining Software Repositories (MSR 2010)*, pages 1–10. IEEE, may 2010.
- [LDSV11] Ahmed Lamkanfi, Serge Demeyer, Quinten David Soetens, and Tim Verdonck. Comparing Mining Algorithms for Predicting the Severity of a Reported Bug. In *2011 15th European Conference on Software Maintenance and Reengineering*, pages 249–258. IEEE, mar 2011.
- [LHMI07] Simone Livieri, Yoshiki Higo, Makoto Matsushita, and Katsuro Inoue. Very-Large Scale Code Clone Analysis and Visualization of Open Source Programs Using Distributed CCFinder: D-CCFinder. In *29th International Conference on Software Engineering*, pages 106–115. IEEE, may 2007.
- [LLMZ06] Z. Li, S. Lu, S. Myagmar, and Y. Zhou. CP-Miner: finding copy-paste and related bugs in large-scale software code. *IEEE Transactions on Software Engineering*, 32(3):176–192, mar 2006.
- [LLS⁺13] Chris Lewis, Zhongpeng Lin, Caitlin Sadowski, Xiaoyan Zhu, Rong Ou, and E. James Whitehead Jr. Does bug prediction support human developers? findings from a google case study. pages 372–381, may 2013.
- [LNH⁺11] Taek Lee, Jaechang Nam, DongGyun Han, Sunghun Kim, and Hoh Peter In. Micro interaction metrics for defect prediction. In *Proceedings of the 19th ACM SIGSOFT symposium and the 13th European conference on Foundations of software engineering - SIGSOFT/FSE '11*, page 311, New York, New York, USA, 2011. ACM Press.
- [Lo13] D. Lo. A Comparative Study of Supervised Learning Algorithms for Reopened Bug Prediction. In *2013 17th European Conference on Software Maintenance and Reengineering*, pages 331–334. IEEE, mar 2013.
- [LPM⁺] B. Lague, D. Proulx, J. Mayrand, E.M. Merlo, and J. Hudépohl. Assessing the benefits of incorporating function clone detection in a development process. In *Proceedings International Conference on Software Maintenance*, pages 314–321. IEEE Comput. Soc.
- [LvdH11] Nicolas Lopez and André van der Hoek. The code orb. In *Proceeding of the 33rd international conference on Software engineering - ICSE '11*, page 824, New York, New York, USA, may 2011. ACM Press.

- [LWA07] Lucas Layman, Laurie Williams, and Robert St. Amant. Toward Reducing Fault Fix Time: Understanding Developer Behavior for the Design of Automated Fault Detection Tools. In *First International Symposium on Empirical Software Engineering and Measurement (ESEM 2007)*, pages 176–185. IEEE, sep 2007.
- [Man94] Udi Manber. Finding similar files in a large file system. In *Usenix Winter*, pages 1–10. USENIX Association, jan 1994.
- [MANH14] Shane McIntosh, Bram Adams, Meiyappan Nagappan, and Ahmed E Hassan. Mining Co-Change Information to Understand when Build Changes are Necessary. In *Software Maintenance and Evolution (ICSME)*, 2014.
- [Mar09] Martin Fowler. FeatureBranch, 2009.
- [MBFV13] Joao Eduardo Montandon, Hudson Borges, Daniel Felix, and Marco Tulio Valente. Documenting APIs with examples: Lessons learned with the APIMiner platform. In *2013 20th Working Conference on Reverse Engineering (WCRE)*, pages 401–408. IEEE, oct 2013.
- [MGDL10] N. Moha, Y.-G. Gueheneuc, L. Duchien, and A.-F. Le Meur. DECOR: A Method for the Specification and Detection of Code and Design Smells. *IEEE Transactions on Software Engineering*, 36(1):20–36, jan 2010.
- [MHLN⁺15] A. Maiga, W. Hamou-Lhadj, M. Nayrolles, K. Sabor, and A. Larsson. An Empirical Study on the Handling of Crash Reports in a Large Software Company: An Experience Report. In *International Conference on Software Maintenance and Evolution (ICSME)*, pages 342–351. IEEE, 2015.
- [MLM96] Jean Mayrand, Claude Leblanc, and Ettore M. Merlo. Experiment on the automatic detection of function clones in a software system using metrics. In *Proceedings of International Conference on Software Maintenance ICSM-96*, pages 244–253. IEEE, 1996.
- [MM] A. Marcus and J.I. Maletic. Identification of high-level concept clones in source code. In *Proceedings 16th Annual International Conference on Automated Software Engineering (ASE 2001)*, pages 107–114. IEEE Comput. Soc.

- [MPNC12] Naouel Moha, Francis Palma, Mathieu Nayrolles, and Benjamin Joyen Conseil. Specification and Detection of SOA Antipatterns. In *International Conference on Service Oriented Computing*, pages 1–16. Springer, 2012.
- [MSA04] Roman Manevich, Manu Sridharan, and Stephen Adams. PSE: explaining program failures via postmortem static analysis. In *ACM SIGSOFT Software Engineering Notes*, volume 29, page 63. ACM, nov 2004.
- [NAS09] NASA. Open Mission Control Technologies, 2009.
- [NAW⁺08] Syeda Nessa, Muhammad Abedin, W Eric Wong, Latifur Khan, and Yu Qi. Software Fault Localization Using N -gram Analysis. In *WASA*, pages 548–559, 2008.
- [Nay13] Mathieu; Nayrolles. *Improving SOA Antipattern Detection in Service Based Systems by Mining Execution Traces*. PhD thesis, 2013.
- [Nay14] Mathieu Nayrolles. *Mastering Apache Solr - A Practical Guide to Get to Grips with Apache Solr*. 2014.
- [NB05a] N. Nagappan and T. Ball. Use of relative code churn measures to predict system defect density. In *Proceedings. 27th International Conference on Software Engineering, 2005.*, pages 284–292. IEEE, 2005.
- [NB05b] Nachiappan Nagappan and Thomas Ball. Static analysis tools as early indicators of pre-release defect density. In *Proceedings of the 27th international conference on Software engineering - ICSE '05*, page 580, New York, New York, USA, may 2005. ACM Press.
- [NBMV15] Mathieu Nayrolles, Eric Beaudry, Naouel Moha, and Petko Valtchev. Towards Quality-Driven SOA Systems Refactoring through Planning. In *6th International MCETECH Conference*, 2015.
- [NBZ06] Nachiappan Nagappan, Thomas Ball, and Andreas Zeller. Mining metrics to predict component failures. In *Proceeding of the 28th international conference on Software engineering - ICSE '06*, page 452, New York, New York, USA, may 2006. ACM Press.
- [NHL15] M. Nayrolles and W. Hamou-Lhadj. BUMPER: Bug Metarepository Search Engine for Developers and Researchers. Consortium for Software Engineering Research Fall., 2015.

- [NHL16a] M. Nayrolles and W Hamou-Lhadj. A Bug-Fix Metarepository for Developers and Researcher. In *Mining Software Repositories - Data Track (MSR'16)*. IEEE, 2016.
- [NHL16b] M. Nayrolles and W. Hamou-Lhadj. BUMPER: A Tool to Cope with Natural Language Search of Millions Bugs and Fixes. In *International Conference on Software Analysis, Evolution, and Reengineering - Tool Track (SANER'16)*. IEEE, 2016.
- [NHL15] Mathieu Nayrolles, Abdelwahab Hamou-Lhadj, Tahar Sofiene, and Alf Larsson. JCHARMING : A Bug Reproduction Approach Using Crash Traces and Directed Model Checking. In *Proceedings of the 22nd International Conference on Software Analysis, Evolution, and Reengineering, IEEE, 2015*, pages 101–110, 2015.
- [NHL16] M. Nayrolles, W. Hamou-Lhadj, S. Tahar, and A. Larsson. A Bug Reproduction Approach Based on Directed Model Checking and Crash Traces. *Journal of Software: Evolution and Process. Wiley (Accepted with revisions)*., 2016.
- [NMH15] Mathieu Nayrolles, Abdou Maiga, Abdelwahab Hamou-lhadj, and Alf Larsson. A Taxonomy of Bugs : An Empirical Study. pages 1–10.
- [NMV13] Mathieu Nayrolles, Naouel Moha, and Petko Valtchev. Improving SOA Antipatterns Detection in Service Based Systems by Mining Execution Traces. In *Working Conference on Reverse Engineering*, number i, pages 321–330. IEEE, 2013.
- [NNN⁺12] Anh Tuan Nguyen, Tung Thanh Nguyen, Tien N. Nguyen, David Lo, and Chengnian Sun. Duplicate bug report detection with a combination of information retrieval and topic modeling. In *Proceedings of the 27th IEEE/ACM International Conference on Automated Software Engineering - ASE 2012*, page 70, New York, New York, USA, 2012. ACM Press.
- [Nor13] Donald A. Norman. *The Design of Everyday Things: Revised and Expanded Edition*. Basic Books, 2013.
- [NPC05] Satish Narayanasamy, Gilles Pokam, and Brad Calder. BugNet: Continuously Recording Program Execution for Deterministic Replay Debugging. In *Proceedings of the 32nd annual International Symposium on Computer Architecture*, volume 33, pages 284–295. ACM, may 2005.

- [Obj05] Object Refinery Limited. JFreeChart, 2005.
- [OP99] L. Opyrchal and A. Prakash. Efficient Object Serialization in Java. Lukasz Opyrchal and Atul Prakash. In *Proceedings. 19th IEEE International Conference on Distributed Computing Systems*, pages 96–101, 1999.
- [Ora11] Oracle. Throwable (Java Platform SE6), 2011.
- [OWB05] T.J. Ostrand, E.J. Weyuker, and R.M. Bell. Predicting the location and number of faults in large software systems. *IEEE Transactions on Software Engineering*, 31(4):340–355, apr 2005.
- [Pal13] Francis Palma. *Detection of SOA Antipatterns*. PhD thesis, Ecole Polytechnique de Montreal, 2013.
- [Pan07] Lucas D. Panjer. Predicting Eclipse Bug Lifetimes. In *Fourth International Workshop on Mining Software Repositories (MSR’07:ICSE Workshops 2007)*, pages 29–29. IEEE, may 2007.
- [PHM⁺04] Jian Pei, Jiawei Han, Senior Member, Behzad Mortazavi-asl, Jianyong Wang, Helen Pinto, and Qiming Chen. Mining Sequential Patterns by Pattern-Growth : The PrefixSpan Approach. *IEEE Transactions on Knowledge and Data Engineering*, 16(10):1424—1440, 2004.
- [PHMa00] Jian Pei, Jiawei Han, and Behzad Mortazavi-asl. Mining Access Patterns Efficiently from Web Logs *. In *Knowledge Discovery and Data Mining. Current Issues and New Applications*, volume 0, pages 396–407, 2000.
- [PMDL99] J.-F. Patenaude, E. Merlo, M. Dagenais, and B. Lague. Extending software quality assessment techniques to Java systems. In *Proceedings Seventh International Workshop on Program Comprehension*, pages 49–56. IEEE Comput. Soc, 1999.
- [PO11] Chris Parnin and Alessandro Orso. Are automated debugging techniques actually helping programmers? *Proceedings of the 2011 International Symposium on Software Testing and Analysis - ISSTA ’11*, page 199, 2011.
- [Pra01] Michael J. Pratt. Introduction to ISO 10303the STEP Standard for Product Data Exchange. *Journal of Computing and Information Science in Engineering*, 1(1):102, mar 2001.

- [Pre05] Roger S. Pressman. *Software Engineering: A Practitioner's Approach*. Palgrave Macmillan, 2005.
- [PS93] Perry, Dewayne E. and Carol S. Stieg. Software faults in evolving a large, real-time system: a case study. In *Software EngineeringESEC*, pages 48–67, 1993.
- [RAN07] Per Runeson, Magnus Alexandersson, and Oskar Nyholm. Detection of Duplicate Defect Reports Using Natural Language Processing. In *29th International Conference on Software Engineering*, pages 499–510. IEEE, may 2007.
- [RC08] C.K. Roy and J.R. Cordy. NICAD: Accurate Detection of Near-Miss Intentional Clones Using Flexible Pretty-Printing and Code Normalization. In *2008 16th IEEE International Conference on Program Comprehension*, pages 172–181. IEEE, jun 2008.
- [RGK90] Jane Radatz, Anne Geraci, and Freny Katki. IEEE standard glossary of software engineering terminology. *IEEE Std*, 610121990(121990):3, 1990.
- [RM09] Neha Rungta and Eric G. Mercer. Guided model checking for programs with polymorphism. In *Proceedings of the 2009 ACM SIGPLAN workshop on Partial evaluation and program manipulation - PEPM '09*, page 21, New York, New York, USA, 2009. ACM Press.
- [ROY09] CHANCHAL K. ROY. *DETECTION AND ANALYSIS OF NEAR-MISS SOFTWARE CLONES*. PhD thesis, Queen's University, 2009.
- [RZF⁺13] Jeremias Rößler, Andreas Zeller, Gordon Fraser, Cristian Zamfir, and George Candea. Reconstructing Core Dumps. In *Proceedings of the 6th International Conference on Software Testing, Verification and Validation, IEEE*, 2013.
- [SCFP00] John Steven, Pravir Chandra, Bob Fleck, and Andy Podgurski. jRapture: A Capture/Replay Tool for Observation-Based Testing. In *Proceedings of the International Symposium on Software Testing and Analysis.*, number August, pages 158–167, 2000.
- [SIK⁺10] Emad Shihab, Akinori Ihara, Yasutaka Kamei, Walid M. Ibrahim, Masao Ohira, Bram Adams, Ahmed E. Hassan, and Ken-ichi Matsumoto. Predicting Re-opened Bugs: A Case Study on the Eclipse Project. In *2010 17th Working Conference on Reverse Engineering*, pages 249–258. IEEE, oct 2010.

- [SK03] R. Subramanyam and M.S. Krishnan. Empirical analysis of CK metrics for object-oriented design complexity: implications for software defects. *IEEE Transactions on Software Engineering*, 29(4):297–310, apr 2003.
- [SKP14] Ripon K. Saha, Sarfraz Khurshid, and Dewayne E. Perry. An empirical study of long lived bugs. In *2014 Software Evolution Week - IEEE Conference on Software Maintenance, Reengineering, and Reverse Engineering (CSMR-WCRE)*, pages 144–153. IEEE, feb 2014.
- [SLKJ11] Chengnian Sun, David Lo, Siau-Cheng Khoo, and Jing Jiang. Towards more accurate retrieval of duplicate bug reports. *2011 26th IEEE/ACM International Conference on Automated Software Engineering (ASE 2011)*, pages 253–262, nov 2011.
- [SLW⁺10] Chengnian Sun, David Lo, Xiaoyin Wang, Jing Jiang, and Siau-Cheng Khoo. A discriminative model approach for accurate duplicate bug report retrieval. In *Proceedings of the 32nd ACM/IEEE International Conference on Software Engineering - ICSE '10*, volume 1, page 45, New York, New York, USA, may 2010. ACM Press.
- [SNH13] Weiyi Shang, Meiyappan Nagappan, and Ahmed E. Hassan. Studying the relationship between logging characteristics and the code quality of platform software. *Empirical Software Engineering*, pages 1–27, 2013.
- [TBG] M. Toomim, A. Begel, and S.L. Graham. Managing Duplicated Code with Linked Editing. In *2004 IEEE Symposium on Visual Languages - Human Centric Computing*, pages 173–180. IEEE.
- [TG06] Robert Tairas and Jeff Gray. Phoenix-based clone detection using suffix trees. In *Proceedings of the 44th annual southeast regional conference on - ACM-SE 44*, page 679, New York, New York, USA, mar 2006. ACM Press.
- [The99] The Apache Software Foundation. Log4j 2 Guide - Apache Log4j 2, 1999.
- [TLS12] Yuan Tian, David Lo, and Chengnian Sun. Information Retrieval Based Nearest Neighbor Classification for Fine-Grained Bug Severity Prediction. In *2012 19th Working Conference on Reverse Engineering*, pages 215–224. IEEE, oct 2012.
- [TNAKN11a] Ahmed Tamrawi, Tung Thanh Nguyen, Jafar Al-Kofahi, and Tien N. Nguyen. Fuzzy set-based automatic bug triaging. In *Proceeding of the 33rd*

international conference on Software engineering - ICSE '11, page 884, New York, New York, USA, 2011. ACM Press.

- [TNAKN11b] Ahmed Tamrawi, Tung Thanh Nguyen, Jafar M. Al-Kofahi, and Tien N. Nguyen. Fuzzy set and cache-based approach for bug triaging. In *Proceedings of the 19th ACM SIGSOFT symposium and the 13th European conference on Foundations of software engineering - SIGSOFT/FSE '11*, page 365, New York, New York, USA, sep 2011. ACM Press.
- [TSL12] Yuan Tian, Chengnian Sun, and David Lo. Improved Duplicate Bug Report Identification. In *2012 16th European Conference on Software Maintenance and Reengineering*, pages 385–390. IEEE, mar 2012.
- [VHB⁺03] Willem Visser, Klaus Havelund, Guillaume Brat, SeungJoon Park, and Flavio Lerda. Model Checking Programs. In *Automated Software Engineering*, volume 10, pages 203–232. Springer, 2003.
- [VPK04] Willem Visser, Corina S. Psreanu, and Sarfraz Khurshid. Test input generation with java PathFinder. *Proceedings of the 2004 ACM SIGSOFT International Symposium on Software Testing and Analysis*, page 97, 2004.
- [WAC12] James A. Whittaker, Jason Arbon, and Jeff Carollo. *How Google Tests Software*. Addison-Wesley, 2012.
- [Wel13] Brian Wellington. *dnsjava*, 2013.
- [WM05] R. Wettel and R. Marinescu. Archeology of code duplication: recovering duplication chains from small duplication fragments. In *Seventh International Symposium on Symbolic and Numeric Algorithms for Scientific Computing (SYNASC'05)*, page 8 pp. IEEE, 2005.
- [WPZZ07] Cathrin Weiss, Rahul Premraj, Thomas Zimmermann, and Andreas Zeller. How Long Will It Take to Fix This Bug? In *Fourth International Workshop on Mining Software Repositories (MSR'07:ICSE Workshops 2007)*, pages 1–1. IEEE, may 2007.
- [WSWF] V. Wahler, D. Seipel, J. Wolff, and G. Fischer. Clone detection in source code by frequent itemset techniques. In *Source Code Analysis and Manipulation, Fourth IEEE International Workshop on*, pages 128–135. IEEE Comput. Soc.

- [WZKC11] Rongxin Wu, Hongyu Zhang, Sunghun Kim, and SC Cheung. Relink: recovering links between bugs and changes. In *Proceedings of the 19th ACM SIGSOFT symposium and the 13th European conference on Foundations of software engineering.*, pages 15–25, 2011.
- [Xst11] Xstream. Xstream, 2011.
- [ZGV13] Hongyu Zhang, Liang Gong, and Steve Versteeg. Predicting bug-fixing time: an empirical study of commercial software projects. pages 1042–1051, may 2013.
- [ZN08] Thomas Zimmermann and Nachiappan Nagappan. Predicting defects using network analysis on dependency graphs. In *Proceedings of the 13th international conference on Software engineering - ICSE '08*, page 531, New York, New York, USA, may 2008. ACM Press.
- [ZNGM12] Thomas Zimmermann, Nachiappan Nagappan, Philip J. Guo, and Brendan Murphy. Characterizing and predicting which bugs get reopened. In *Proceedings of the 34th International Conference on Software Engineering, IEEE*, pages 1074–1083. IEEE, jun 2012.
- [ZPZ07] Thomas Zimmermann, Rahul Premraj, and Andreas Zeller. Predicting Defects for Eclipse. In *Third International Workshop on Predictor Models in Software Engineering (PROMISE'07: ICSE Workshops 2007)*, pages 9–9. IEEE, may 2007.
- [ZZL12] Jian Zhou, Hongyu Zhang, and David Lo. Where should the bugs be fixed? More accurate information retrieval-based bug localization based on bug reports. In *Proceedings of the 34th International Conference on Software Engineering, IEEE*, pages 14–24. IEEE, jun 2012.