

Understanding the Triaging and Fixing Processes of Long Lived Bugs

Ripon K. Saha, Sarfraz Khurshid, Dewayne E. Perry

*Center for Advanced Research in Software Engineering (ARiSE)
Department of Electrical and Computer Engineering
The University of Texas at Austin, USA*

Abstract

Context: Bug fixing is an integral part of software development and maintenance. A large number of bugs often indicate poor software quality, since buggy behavior not only causes failures that may be costly but also has a detrimental effect on the user's overall experience with the software product. The impact of *long lived* bugs can be even more critical since experiencing the same bug version after version can be particularly frustrating for user. While there are many studies that investigate factors affecting bug fixing time for entire bug repositories, to the best of our knowledge, none of these studies investigates the extent and reasons of long lived bugs.

Objective: In this paper, we investigate the triaging and fixing processes of long lived bugs so that we can identify the reasons for delay and improve the overall bug fixing process.

Methodology: We mine the bug repositories of popular open source projects, and analyze long lived bugs from five different perspectives: their proportion, severity, assignment, reasons, as well as the nature of fixes.

Results: Our study on seven open-source projects shows that there are a considerable number of long lived bugs in each system and over 90% of them adversely affect the user's experience. The reasons for these long lived bugs are diverse including long assignment time, not understanding their importance in advance etc. However, many bug-fixes were delayed without any specific reasons. Furthermore, 40% of long lived bugs need small fixes.

Conclusion: Our overall results suggest that many long lived bugs can be fixed quickly through careful triaging and prioritization, if developers could predict their severity, change effort, and change impact in advance. We believe our results will help both developers and researchers better to understand factors behind delays, improve the overall bug fixing process, and investigate analytical approaches for prioritizing bugs based on bug severity as well as expected bug fixing effort.

Keywords: Bug tracking, bug triaging, bug survival time

1. Introduction

Software development and maintenance is a complex process. Although developers and testers try their best to make their software error free, in practice software ships with bugs. The number of bugs in software is a significant indicator of software quality since bugs can adversely affect users experience directly. Therefore, developers are generally very active in finding and removing bugs.

To ensure high software quality for each release, developers/managers triage bugs carefully and schedule

the bug fixing tasks based on their severity and priority. Despite such a rigorous process, there are still many bugs that live for a long time. We believe the impact of these *long lived* bugs (for our study, bugs that are not fixed within one year after they are reported) is even more critical since the users may experience the same failures version after version. Therefore, it is important to understand the extent and reasons of these long lived bugs so that we can improve software quality.

A number of previous studies have investigated the overall factors affecting bug fix time. Giger et al. [8] empirically investigated the relationships between bug report attributes and the time to fix. Zhang et al. [31] predicted overall bug fix time in commercial projects. Canfora et al. [6] used survival analysis to determine the relationship between the risk of not fixing a bug within a

Email addresses: ripon@utexas.edu (Ripon K. Saha),
khurshid@ece.utexas.edu (Sarfraz Khurshid),
perry@ece.utexas.edu (Dewayne E. Perry)

given time frame and specific code constructs changed when fixing the bug. Zhang et al. [30] examined factors affecting bug fixing time along three dimensions: bug reports, source code involved in the fix, and code changes that are required to fix the bug.

While these studies are useful in understanding the overall factors related to bug fix time, we know of no study that has specifically investigated long lived bugs to understand why they take such a long time to be fixed and how important they are. We point out that analyzing entire bug datasets using various machine learning or data mining techniques (as done in previous work) is not sufficient in understanding long lived bugs due to the imbalanced dataset.¹ Imbalanced datasets are a major problem in most data mining applications since machine learning algorithms can be biased towards the majority class due to over-prevalence [12]. We expect (and our results also support) that the proportion of long-lived bugs would be lot less than 50% of the total bugs, thus resulting an imbalanced dataset. Therefore, if we automatically analyze all the bug reports using a standard data mining technique, it is highly likely that the main factors behind long lived bugs would get lost. In this paper, we conduct an exploratory study focused solely on long lived bugs to understand their extent and reasons with respect to following research questions:

1. **What proportion of the bugs are long lived?** The answer to this question is important since if there are few long lived bugs, there may be little reason to worry.
2. **How important long lived bugs are in terms of severity?** It is important to understand how crucial these bugs were from the perspective of both developers and users. If they are minor or trivial bugs, their impact would be less on overall software quality.
3. **Where is most of the time spent in the bug fixing process?** The answer to this question is important to identify the time consuming phases so that developers as well as researchers can work on improving the process involving that phase.
4. **What are common reasons for long lived bugs?** To improve the bug fixing process, first we need to understand the underlying reasons for delays. Delineating the common reasons of long lived bugs will help researchers deal with the problem more systematically.

¹A dataset is imbalanced if the classification classes are not approximately equally represented.

5. **What is the nature of long lived bug fixes?** The answer to this question will help us in better understanding the bug fixing process, estimating change efforts, and so on, which will be useful in exploring potential approaches for improving overall bug fixing process.

We study seven open source projects: JDT, CDT, PDE, and Platform from the Eclipse product family, written in Java, ² and the Linux Kernel, WineHQ, and GDB, written in C. Our key observations are summarized below:

1. Despite advances in software development and maintenance processes, there are a significant number of bugs in each project that survive for more than one year.
2. More than 90% of long lived bugs affect users' normal working experiences and thus are important to fix. Moreover, there are duplicate bug reports for these long lived bugs, which indicate the users' demand for fixing them.
3. The average bug assignment time of these bugs was more than one year despite the availability of a number of automatic bug assignment tools that could have been used. The bug fix time after the assignment was another year on average.
4. Reasons for long lived bugs are diverse. While problem complexity, reproducibility, and not understanding the importance of some of the bugs in advance are the common reasons, we observed there are many bug-fixes that got delayed without any specific reason.
5. Unlike previous studies [30], we found that a bug surviving for a year or more does not necessarily mean that it requires a large fix. We found that 40% of long-lived bug fixes involved few changes in only one file.

This paper extends our previous "long lived bugs" paper presented at CSMR-WCRE 2014 [22] in three directions. First, we perform the same set of experiments on three additional popular projects: the Linux Kernel, WineHQ, and GDB, which are not only written in different programming language but also from different domains than our previous subject systems. Second, we provide more details of our previous results. Finally, our new results show that our previous findings hold even for software projects from different domains

²<http://www.eclipse.org>

and written in different languages. Thus this paper reports more generalizable results. We believe these findings will play an important role in developing new approaches for bug triaging as well as improving the overall bug fixing process.

2. Background

In this section, we provide necessary background for our study that includes a brief description of bug tracking systems and a typical bug life cycle.

2.1. Bug Tracking System:

Generally project stakeholders maintain a bug database for tracking all the bugs associated with their projects. There are several online bug tracking systems such as Bugzilla, JIRA, Mantis, etc. These systems enable developers/managers to manage bug database for their projects. Different repositories may have different data structures and follow different life cycles of bugs. The dataset we created and used in our work was extracted from Bugzilla, a popular online bug tracking system. Therefore, the rest of the discussion in this paper regarding the bug tracking system is only limited to Bugzilla.

Any person having legitimate access to a project's bug database can post a change request through Bugzilla. A change request could be either a bug or an enhancement. In Bugzilla, however, both bugs and enhancements are represented similarly and referred as bugs with an exception that for enhancements severity field is set to enhancement. Generally bug reporters provide a bug summary, bug description, the suspected product, and the component name with its severity.

Developers in a particular project can define their own severity level. According to Eclipse Bugzilla documentation, the severity level can be one of the following values, which actually represents the degree of potential harm.³

Blocker: These bugs block the development and/or testing work. There exists no workaround.

Critical: These bugs cause program crashes, loss of data, or severe memory leaks.

Major: These bugs result major loss of function.

Normal: These are regular issues. There are some loss of functionality under specific circumstances.

Minor: These bugs cause minor loss of functionality, or other problems where an easy workaround was present.

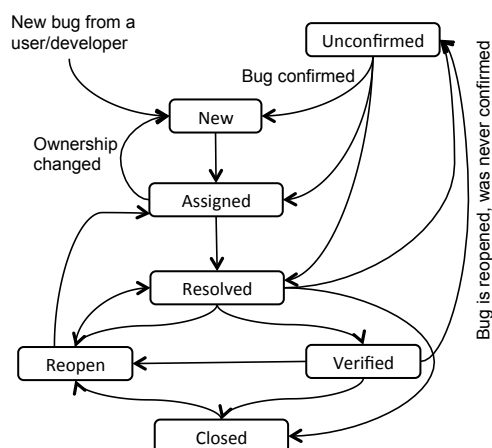


Figure 1: Life Cycle of a Bug in Bugzilla

Trivial: These are generally cosmetic problems such as misspelled words or misaligned text.

The developers in WineHQ also follow the same severity levels. However, the GDB community recognizes three levels of severity: critical, normal, and minor. On the other hand, the Linux community has their own severity level: blocking, high, normal, and low.

In addition to providing severity level, reporters also specify the software version, the platform and operating system where they encountered the bug so that developers can easily reproduce it. Bug reporters also can attach files to the bug report such as screen shots, failing test cases etc. Once a bug is posted, all other related developers can make comments regarding the bug to discuss different issues. Therefore, a bug repository has rich set of information that can be analyzed to gain insight about bugs.

2.2. Bug Life Cycle

The overall bug fixing process in a system is directly related to the bug life cycle maintained by the bug tracking system. Although different projects may have different schemes for using Bugzilla, a common life cycle for a bug is as follows:⁴

Validation: At the start of each day, each project/component team leader triages NEW bugs to verify if the bug is really a bug and if the provided information is correct. In case of any inconsistencies, the bug triager can correct them. The bug triager also can request further information to validate a bug if it is necessary. If there is no response within a week, the team leader closes the bug marking RESOLVED, INVALID,

³http://wiki.eclipse.org/Eclipse/Bug_Tracking

⁴http://wiki.eclipse.org/Development_Resources/HOWTO/Bugzilla_Use

or `WONTFIX`. However, the reporter can reopen the bug anytime if she has more information.

Prioritization: In this stage, the triager first determines whether a bug is a feature request. If so, the severity of the bug is changed to `enhancement`. Otherwise, she checks the severity level of the bug to make sure that it is consistent with the bug description. Then the priority of the bug is set based on following guidelines:⁵

`P1` : These bugs are a must fix for the indicated target milestone.

`P2` : These bugs are very important for the indicated target milestone. Generally developers try to resolve all the `P2` bugs.

`P3` : It is the default priority. If the bug triager is uncertain about the priority of a bug or it is actually a normal bug, she can set `P3` priority. Then the assigned developer can adjust it if appropriate.

`P4` : These bugs should be fixed if time permits.

`P5` : These are valid bugs, but there are no plans to fix. Also `P5` priority indicates that help is wanted.

Fixing: At this point, a bug remains in the component's "inbox" account until a developer takes the bug, or the team leader assigns it to them. After fixing the bug, the developer mark it as `RESOLVED-FIXED`.

Verification: Once a bug is fixed, it is assigned to another committer on the team to verify. Ideally, all bugs should be verified before the next integration build. Once the verifier tests that the bug is completely resolved, she changes the bug status to `VERIFIED`. Figure 1 represents all possible state transitions of a bug in Bugzilla.

3. Study Setup

This section provides a brief description of the subject systems that we studied, and the metrics and process we used to understand the extent and reason of long lived bugs.

3.1. Subject Systems

We use seven open source projects for our study. Among them, we choose four projects from the Eclipse product family, namely, JDT, CDT, PDE, and Platform, which are written in Java programming language. The other three projects are the Linux Kernel, WineHQ, and GDB, which are written in C programming language. There are mainly three reasons for choosing

⁵http://wiki.eclipse.org/WTP/Conventions_of_bug_priority_and_severity

these projects. First, These projects are highly successful and have been widely used in software engineering research. Second, each project has a long development history. Third, these projects are from different domains.

- JDT and CDT provide a fully functional Integrated Development Environment based on the Eclipse platform for developing Java, and C and C++ applications.^{6,7}
- The Plug-in Development Environment (PDE) provides tools to create, develop, test, debug, build and deploy Eclipse plug-ins, fragments, features, update sites and RCP products.⁸
- The Eclipse Platform defines the set of frameworks and common services that collectively make up infrastructure required to support the use of Eclipse.⁹
- Linux Kernel: The kernel of the Linux operating system.¹⁰
- WineHQ: A compatibility layer, making it possible to run Windows applications on POSIX compliant operating systems.¹¹
- GDB: A debugger for programs written in C, C++, and many other programming languages.¹²

We have created the dataset for C projects. This dataset includes all the bug reports and their histories from their inception to May 2014. For Java projects, we have used Lamkanfi et al's [16] bug dataset to extract the bug information. The Java dataset includes all the bug reports and their histories from their inception to March 2011 for these four projects (extracted from Eclipse Bugzilla database). More detailed descriptions of the dataset is presented in Table 1. In the Table, the last column represents the number of bugs (excluding enhancement and duplicated bugs) that got eventually fixed, which is the actual dataset of this study.

3.2. Terms and Metrics

We make use of bug tracking and version control system's information to calculate metrics that we were interested in. This section defines different terms and metrics that we use in the rest of the paper.

⁶<http://projects.eclipse.org/projects/eclipse.jdt>

⁷<http://www.eclipse.org/cdt/>

⁸<http://www.eclipse.org/pde/>

⁹<https://projects.eclipse.org/projects/eclipse.platform>

¹⁰<https://www.kernel.org/>

¹¹<http://www.winehq.org>

¹²<http://www.sourceware.org/gdb/>

Table 1: Data Set

System	#CR	#Bugs	# Enh.	# Bug Fixed
JDT	46,308	38,520	7,788	18,873
CDT	14,871	12,854	20,17	7,260
PDE	13,677	11,958	1,719	6,854
Platform	90,691	78,120	12,571	33,738
Linux	23,618	23,387	231	5,784
WineHQ	36,691	34,490	2,201	14,338
GDB	17,038	15,562	1,476	7,667

Bug Introduction Time (T_I): This is the timestamp when the buggy code is committed for the first time for a given bug.

Bug Reporting Time (T_R): This is the timestamp when a bug is reported to the Bugzilla system by a user/developer.

Bug Assignment Time (T_A): This is the timestamp when a bug was officially assigned to the right developers through Bugzilla. If a bug is assigned to multiple developers, we use the assignment time of the developer who fixed the bug. If a bug is fixed by multiple developers, we use the assignment time of the developer who committed the last changes.

Bug Severity Realization Time (T_S): This is the timestamp when the actual severity of a given bug was understood by the developers and thus the severity field of that bug was changed for the last time.

Bug Fix Time (T_F): This is the timestamp when a developer officially marked a bug as `FIXED` in Bugzilla through the `resolution` field.

Bug Assignment Period (AP): This is the lapse time between when the bug was opened and when it was assigned to the right developer. Mathematically, $AP = T_A - T_R$

Bug Fixing Period (FP): This is the period of time that developers took to fix a bug. It should be noted that it is not the actual coding time of the bug-fix. Instead, it is the time period between the bug assignment time and the bug fix time. Mathematically, $FP = T_F - T_A$. It should be further noted that we do not deduct the time if a bug is temporarily closed. A bug is temporarily closed when the developers think that the bug is fixed but actually it is not. Therefore, we do not deduct that time, since the bug is still (at least partially) present during the time when the bug reports was closed.

Pre-Severity Realization Period ($Pre-SRP$): This is the period of time developers took to understand the actual severity of the bug. Therefore, pre-severity realization time is the time between bug reporting time and

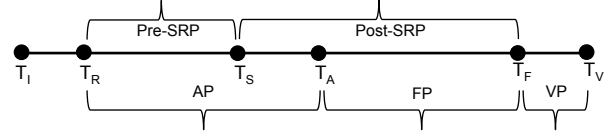


Figure 2: An Example Timeline of a Bug

the time when the severity was changed for the last time. Mathematically, $Pre-SRP = T_S - T_R$.

Post-Severity Realization Period ($Post-SRP$): This is the time developers took to fix the bug after realizing the actual severity time. Mathematically, $Post-SRP = T_F - T_S$.

Bug Verification Period (VP): This is the period of time that a developer took to verify a bug after it is marked as `FIXED` in Bugzilla. Mathematically, $VP = T_V - T_F$.

Bug Survival Period (SP): This is the period that a bug exists in the system. Ideally it should be the time period between the bug introduction time (T_I) and bug fixing time (T_F). However, in our study it is the time period between T_R and T_F . The timestamps of bug introduction (T_I) and bug reporting (T_R) can be certainly different, since a bug can remain dormant for a long time [7]. Although there are some algorithms [14] to identify bug introducing changes, it is difficult to map those changes to the associated bug reports. Therefore, we preferred bug reporting time over bug introduction time. Furthermore, since T_R is always later than T_I (i.e., a bug is always reported after the bug introducing changes are committed), our calculated bug survival period (SP) never overestimates actual SP . However, we do not subtract the time period from SP when a bug was temporarily closed. Figure 2 visually presents all the terms and metrics in a timeline.

3.3. Identification of Faulty Source Code

Previous studies [13] showed that when developers fix bugs they often put the bug id in their commit message. Therefore, to get the version histories and commit messages of these four projects, first we accessed their git repositories. Then using JGit APIs, we extracted all the commit messages from the histories and searched all numbers.^{13,14,15,16} Then we matched each number with the bug IDs. To further ensure that those are indeed

¹³<http://www.eclipse.org/jgit/>

¹⁴[git://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git](https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git)

¹⁵[git://source.winehq.org/git/wine.git/](https://source.winehq.org/git/wine.git/)

¹⁶git://sourceware.org/git/binutils-gdb.git

Table 2: Bug Fix Time

Time	JDT		CDT		PDE		Platform		Linux		WineHQ		GDB	
	#Bugs	[%]	#Bugs	[%]	#Bugs	[%]	#Bugs	[%]	#Bugs	[%]	#Bugs	[%]	#Bugs	[%]
<1 day	5,009	26.54	1,911	26.32	2,154	31.43	8,244	24.44	578	9.99	673	4.69	1,591	20.75
1-7 days	4,788	25.37	1,485	20.45	1,570	22.91	7,420	21.99	881	15.23	1,809	12.62	1,470	19.17
8-30 days	3,704	19.63	1,230	16.94	1,293	18.86	6,442	19.09	1,221	21.11	1,656	11.55	1,358	17.71
1-6 mon.	3,604	19.10	1,426	19.64	1,212	17.68	7,101	21.05	1,573	27.2	2,837	19.79	1,654	21.57
6-12 mon.	855	4.53	567	7.81	353	5.15	2,173	6.44	578	9.99	2,062	14.38	511	6.66
>1 year	913	4.84	641	8.83	272	3.97	2,358	6.99	953	16.48	5,301	36.97	1,083	14.13
Total	18,873		7,260		6,854		33,738		5,784		14,338		7,667	

```
commit 768b107e4b3be0acf6f58e914afe4f337c00932b
Date: Fri May 4 11:29:56 2012 +0200
```

```
drm/i915: disable sdvo hotplug on i945g/gm
```

```
v2: While at it, remove the bogus hotplug_active read,
and do not mask hotplug_active[0] before checking
whether the irq is needed, per discussion with Daniel
on IRC.
```

```
Bugzilla:
https://bugzilla.kernel.org/show_bug.cgi?id=38442
```

Figure 3: A bug fixing commit for #38442 in Linux Kernel

bug IDs, we only accepted those commits that contain additional information. For example, in Java projects, the term `bug(s)` (case insensitive) was present. In the Linux Kernel, we found that developers referred to the Bugzilla URL (as shown in Figure 3), whereas in GDB the bug was referred by the term `PR`. In this way, we reduced the chance of getting false positives, although we might missed some true mappings.

For one of our considered projects, WineHQ, however, the above process gave no results. We thus consulted with a developer from the WineHQ community who informed us that in this community the convention is for the bug report to refer back to the commit, rather than the commit referring to the bug report. Indeed, in the WineHQ Bugzilla, there is a dedicated field for a git commit id. However, many of these fields were empty since the field is not required. In this way, we identified the bug fixing commits for those long-lived bugs, where the information was available in the commit messages. Then we used `git diff` to compute following metrics for bug fixes:

Number of Changed Files: It is the number of files that went for changes in the bug fixing commit. If a bug was fixed in multiple commits, it is the total number of distinct files in all commits.

Number of Hunks: A hunk is a chunk of adjacent lines that was changed. For a bug fix spanning over mul-

iple commits, it is total number of hunks in all commits. This is useful to understand how many times developers had to move here and there to fix a bug.

Code Churn: This is the total number of changed lines. Since we use `git diff` itself, the changes in comments were counted as well. For multiple commits, it is the total number of changed lines in all commits. It should be noted that if a line is changed, it is considered as a line deletion first and then addition of another line. Thus the value of code churn for a line change is two.

4. Study Results

In this section, we present the experimental results which answer our research questions.

4.1. RQ1: What proportion of the bugs are long lived?

The first question of any empirical study is how large is the population that we want to study. The answer is important since if the population is small, there may be little reason to worry about them. In this cases, our population of interest is long lived bugs.

The definition of long lived bugs is subjective since the time threshold for deciding whether a bug is long lived or short lived could vary across projects, persons, or studies. In this research question, we analyze the survival time of all the fixed bugs in each subject system and define the long lived bugs more concretely for our study.

Although many of us believe that a bug could be considered as long lived if it survives more than six months, in this study we have considered only those bugs as long lived that survive more than one year. There are two main reasons behind this decision. First, we wanted to be more conservative so that we can investigate really long lived bugs. Second, the release cycle of the subject systems that we considered vary from 2 months

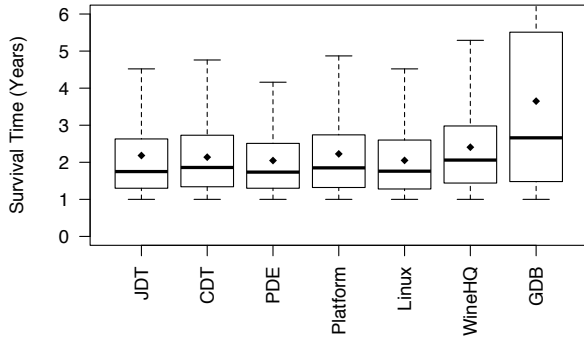


Figure 4: Survival Time of Long Lived Bugs

to one year. More specifically, the Eclipse Foundation has coordinated an *annual simultaneous release* for all projects. GDB release-interval varies from 6 to 10 months. Linux Kernel community make most of the stable releases between 3 months and 1 year. Therefore, if a bug was not fixed in one year, it is expected that the bug propagated through at least two major releases. And it would not be a pleasant experience for a user if s/he experiences the same bug in subsequent major versions of a software.

To investigate long lived bugs, we first investigate how active developers were in fixing bugs in these subject systems. To this end, we group bugs based on their survival period (SP , defined in Section 3.2) and count the number of bugs in each group as shown in Table 2. Results show that around 50%(+/-4%) of the total (fixed) bugs in Java projects were fixed within a week. In C projects, the bug fixing rate was slower than that of Java projects; it took a month to fix around 50% of total bugs (except WineHQ). This indicates that even in open source project, developers are active in fixing bugs. However, as the results show, 10% to 17% of bugs in Java projects and 20% to 50% of bugs in C projects took more than six months to be fixed.

Even after considering such a conservative definition, we found more than 4,184 and 7,337 long lived bugs in Java and C projects respectively. It should be noted that all these bugs eventually got fixed. Therefore, they all are valid bug reports. We believe these are large numbers and thus it is important to investigate them quantitatively and qualitatively.

Figure 4 presents the survival time distribution of long lived bugs. It shows the average (dot in the box), median (line in the box), upper/lower quartile, and 90th/10th percentile of survival time. We limit the values of Y axis to 6 years to better represent the figure.

Table 3: Importance of Long Lived Bugs

System	Blkr.	Critical	Major	Normal	Minor	Trivial
JDT	1	5	49	712	119	27
CDT	2	4	61	523	37	14
PDE	0	6	20	224	13	9
Platform	6	42	221	1856	170	63
Linux	32	130	N/A	743	48	N/A
WineHQ	24	39	185	4200	632	221
GDB	N/A	53	N/A	977	53	N/A

From the figure, we see that at least 25% (upper quartile) of long lived bugs took more than 2.5 years to be fixed. For GDB, it is close to 5.5 years.

Among the total number of bugs that developers fixed, 5%-9% in Java projects and 14%-37% in C projects took more than one year to be fixed.

4.2. RQ2: How important long lived bugs are in terms of severity?

There are two fields in Bugzilla that indicate the importance of a bug: i) severity and ii) priority. However, based on their usage, severity is more important than priority to understand the importance, since severity represents the degree of the impact of the bug on the operation of the system. On the other hand, priority often describes the relative work schedule of fixing a bug set by the developers for a given milestone. For example, if there are 10 critical bugs in a system but developers have time to fix only five bugs, they can set higher priority to any five bugs based on some consideration and set a relatively lower priority to others. Sometimes, developers can set high priority to even a less severe bug, if it is expected to fix easily than a critical bug. Therefore, for this research question, we emphasize on severity over priority. Our initial hypothesis was that *most of the long lived bugs are either minor or trivial, which do not have serious effects.*

4.2.1. Bug Severity

Table 3 presents the importance of long lived bugs based on their severity.

Our results show that almost 90% of the long lived bugs have severity level of `normal` or above both for Java and C projects. Project-wise the proportion varied from 84% to 95%. According to the Eclipse Bugzilla documentation, only `minor` and `trivial` bugs do not interfere with normal work or use, which means that any bugs having severity level `normal` and above adversely

Table 4: Analysis of Severity for Critical and Major Bugs

System	# Bugs	Severity Changed	Proportion	Maximum Changed
JDT	54	23	42.59%	3
CDT	65	21	32.31%	4
PDE	26	11	42.31%	3
Platform	263	115	43.73%	5
Linux	152	14	8.64%	1
WineHQ	248	68	27.42%	3
GDB	53	10	18.89%	1

affect user experiences. Taking that information into account and assuming similar interpretation applies to C projects, we believe that the delay in the long bug fixing process was not due to the fact that they were trivial.

Now let us take a closer look into more severe bugs: `critical` and `major` (blocker bugs generally do not interfere users directly). Our results show that for Java projects, only 1% to 2% of long lived bugs were `critical`, whereas 5% to 10% of long lived bugs were `major` in each system. The absolute number ranged from 4 to 42 for `critical` and 20 to 221 for `major` bugs. For C projects, there are 463 bugs in total that are either `major` or above.

Considering that a `critical` bug causes program crashes and/or data loss and a `major` bug causes major loss of function, these numbers are high, especially since all of them took more than one year to be fixed.

4.2.2. Severity Realization Period (SRP)

As a part of this research question, we are also interested in investigating how long it takes to understand the severity of the bugs. Our initial hypothesis was that *perhaps it took long time to realize the severity of these important (critical or major) bugs. But once the severity was realized it should not take long time to fix them since they are important problems to solve.*

For this analysis, we have considered only those bugs that have severity level of `major` or higher because they are the most important ones. Table 4 represents the number of `critical` and `major` long lived bugs in each system, the number of bugs whose severity level was changed, and the maximum number of time severity level changed for a bug. Results show that for Java projects the severity level of 32%-43% of such bugs was corrected later. For C projects, the change in severity level is only from 8% to 27%. This indicates that the bug reporters could understand the actual severity level of more than 50% of the bugs for Java projects and 70%-90% of the bugs for C projects at the time of bug post-

Table 5: Time Needed for Understanding Bug Severity

System	Pre-SRT (Days)			Post-SRT (Days)		
	Avg.	Med.	Max.	Avg.	Med.	Max.
JDT	374	163	1890	338	320	1712
CDT	80	7	590	760	700	1442
PDE	348	388	1274	300	34	1201
Platform	351	164	2208	498	410	2730
Linux	150	50	1228	644	572	1204
WineHQ	227	55	1756	649	560	2199
GDB	1090	756	2375	377	230	1394

ing. Therefore, it is evident that developers took more than one year to fix a large number of bugs even after they realized that the bugs are very important.

Now we analyze the bugs, whose severity has corrected later. Table 5 presents the average, median, and maximum Pre-SRT and Post-SRT (defined in Section 3.2) values. Our results show that it took almost a year on average to realize the correct severity level of the bug in three of the four Java projects. The only exception is CDT, where the average *Pre-SRP* was 80 days. The maximum *Pre-SRP* of each system shows that for some bug it took several years to realize the severity. On the other hand, for these bugs, it took another year on average to be fixed. For CDT, which was the best in terms of average *Pre-SRP*, *Post-SRP* was more than two years. From the maximum *Post-SRP*, we see that some bugs took even three to eight years to be fixed after developers realized the actual severity level. Therefore, our results indicate that for most long lived bugs in Java projects, *Post-SRP* was high regardless of their *Pre-SRP*.

We see a varying Pre-SRT in C projects. For Linux and WineHQ, severity levels were fixed within six months and a year respectively. For GDB, it took more than three years to understand the actual severity level. However, it should be noted that the severity level of most of the important bugs in C project was known at the time bug posting. However, the Post-SRT was more than a year regardless of whether the actual severity level was understood in advance or later. From the maximum Post-SRT, we see that some `major` or higher severe bugs were fixed after six years.

4.2.3. Duplicate Bugs

Severity is certainly the most reliable information to understand the importance of a bug since it is determined by the bug reporters and supported by developers. However, a large number of duplicate bugs also may express their importance since they often indicate

Table 6: Duplicate Bugs

System	# Bugs	# Duplicate Bugs	# Duplicate Bugs (NOD)						Max NOD
			1	2	3	4	5	>5	
JDT	913	210	101	42	27	10	13	17	26
CDT	641	52	36	8	4	3	0	1	6
PDE	272	50	32	8	2	0	2	6	15
Platform	2,358	495	271	102	51	23	15	33	20
Linux	953	63	46	11	2	2	0	2	10
WineHQ	5,301	603	386	91	41	22	21	42	46
GDB	1,083	102	87	13	1	0	1	0	5

that the scope of the master bug is large and/or the affected users/other developers are getting frustrated [4]. Therefore, in addition to the severity level, we also investigated the number of duplicated bugs.

Table 6 presents an overview of duplicated bugs of long lived bugs. Results show that for 9% to 23% of long lived bugs in Java projects, users/developers submitted multiple bug reports. For C projects, the proportion of duplicate bugs varied between 6% and 11%. From the maximum number of duplicate bugs, we see that some bugs have more than 20 duplicated bug reports in Java projects. In WineHQ, there is a bug (id #6971), for which 46 duplicate bug reports were submitted. The middle columns present more fine grained results of duplicated bugs.

More than 90% of long lived bugs affect users' normal working experiences and thus are important to fix. However, it took a long time to fix these bugs even after realizing their severity. Moreover, there are multiple bug reports for these long lived bugs, which indicate the users' demand for fixing them.

4.3. RQ3: Where was most of the time spent in the bug fixing process?

A bug fixing process majorly can be divided into three phases in terms of activity: i) assignment phase ii) fixing phase, and iii) verification phase. In this research question, we analyze the time taken by team leads/developers in each phase. Our initial hypothesis was that *perhaps it took a long time to assign long lived bugs to the appropriate developers. But once the bugs are assigned, it should not take too long to fix them.*

Table 7 presents the average, median, and maximum time of both assignment period (*AP*) and fixing period (*FP*) in terms of days for all long lived bugs. Our results show that it took more than 1.5 years on average to assign the bugs to the appropriate developers in Java projects. The median *AP* also shows that the data is

Table 7: Bug Assignment Time Vs. Bug Fixing Time

System	Assign. Period (AP)		Fixing Period (FP)			
	Avg. Med.	Max.	Avg. Med.	Max		
JDT	463	374	2745	407	376	2854
CDT	603	552	2035	330	97	1815
PDE	484	482	2728	437	393	1622
Platform	459	373	3326	407	409	2854
Linux	347	275	1617	413	363	2472
WineHQ	489	327	3144	606	478	2563
GDB	915	605	4732	269	78	3224

Table 8: Reassignments of Long-lived Bugs

System	Reassigned	Proportion	Max. Reasn.
JDT	771	84.45%	14
CDT	478	74.57%	8
PDE	174	63.97%	10
Platform	2066	87.61%	12
Linux	437	45.86%	8
WineHQ	399	6.33%	5
GDB	367	33.89%	6

fairly normally distributed. The maximum *AP* shows that it can take more than six years to assign sum bugs to the correct developers. From Table 8, we can also observe that most of the long lived bugs in Java projects are reassigned at least once. The proportion of bugs reassigned ranged from 64% to 88%. More than 10% of long lived bugs were reassigned 5 times or more.

For C projects, the information regarding the first bug assignment was not present for all bugs. We found the bug assignment time only for those bugs that was reassigned (439, 401, and 366 bugs in Linux, WineHQ, and GDB respectively). Based on the results from those bugs, we see that the average bug assignment time in C projects varies from one year to three years.

Guo et al.[10] have conducted a study to investigate the reasons for bug reassignment. They observed that reassignments are not always harmful. In fact many reassignments happened to find the appropriate develop-

ers. However, they also observed that the required time for bug fixes increased with the increase of number of reassignments. Therefore, they concluded that excessive reassignments are harmful. They delineated five reasons for bug reassignments: finding the root cause, determining ownership, poor bug report quality, hard to determine proper fix, and workload balancing. Therefore, taking the aforementioned findings into account, our results indicate that the assignment of these long lived bugs was complex and time consuming, supporting our initial hypothesis.

However, unlike our expectations, the average *FP* for all systems was quite high: around a year. By seeing the median *FP* for CDT, we understand the data is skewed. But for the other three subjects, it is not the case. Also the maximum *FP* shows that, like the bug assignment, it took more than five years for some bugs to be fixed after they assigned to the right developers.

On the other hand, for the verification period, we found that most of the bugs were never verified, at least according to Bugzilla data. However, if they do get verified, the verification time is pretty small: less than a month for most of the subject systems.

Bug assignment and bug fixing are still time intensive processes, despite the availability of automatic bug assignment tools that could have been used.

4.4. RQ:4 What are common reasons for long lived bugs?

To answer this question, we first manually analyzed all the `critical` and `major` bug reports from JDT. We have intentionally chosen the highly severe bugs, since they should be taken seriously by the developers and thus, we will be able to identify the actual reasons of delay. We also analyzed 50 recent (`critical` or `major`) long lived bugs from PDE and Platform. Since JDT and CDT are from similar domain, we did not take any bugs from CDT. Finally, we manually analyze 20 bugs of Linux kernel (10 oldest + 10 recent long lived bugs with high severity) to check if we find any new category. In this way, we identified a set of 125 (= 55+25+25+20) bug reports for manual analysis.

Tagging Methodology: As we discussed in Section 2, each bug report contains a summary, description and a list of developers comments, which often provide rich information about the problems associated the bug. In order to identify the underlying reasons, first, we read the bug summary and description to understand the nature of bugs. Second, we carefully analyze developers' comments to understand the reasons for any delays

since developers often discuss different problems associated with a given bug through comments. For most of the cases, the actual reasons were easily identifiable.

To categorize the reasons for delays, we followed an open-ended taxonomy. We incrementally analyzed all the bug reports. For any given bug report, first we identified the high level reason and checked if the reason already fits into any of the existing categories. Otherwise, we create a new category. We have quoted several key comment(s) for most of the categories to better understand the tagging procedure. In the few cases where the reasons were ambiguous, we relied on contextual information. The following summarizes a taxonomy of common reasons for long lived bugs that we found in the subject systems.

1) Hard to understand: Understanding/locating buggy statements/files in a software project is hard. Sometimes, identifying even the buggy component can be hard. For example, there is a bug (#128563) in JDT, where developers had hard time in understanding if it is a VM or JDT bug. The following comments explains the situation:

"I found something quite interesting. If you move the classes from the two output folders into the same directory and you run from there, it works fine. We generate exactly the same bytecodes in both cases. The VM should behave the same. Might be a VM bug."

After two years, another developer commented—"I believe this is our bug, we should not reference a non accessible type in our bytecode. The fact it works at times feel like unspecified behavior from the VM."

2) Uncertain how to fix: Sometimes developers may know how to solve a bug, but need to wait for making the solution consistent/robust with other parts of the software. The following comments in bug # 3849 represents such a scenario:

"I would like to defer this until we know how we will implement the new Code Manipulation Infrastructure. This is only possible if we get a better undo story. Currently we can only push undo commands on the refactorings undo stack if a file is save. Otherwise the next save would flush the current undo stack which would remove the undo object for extract method."

3) Hard to fix: This kind of bug is hard to fix. There were lots of group discussions for a long time regarding different alternative solutions and finally the group agreed on some specific solution.

4) Risky to fix: Sometimes, bugs are caught just before the release. Then if developers think that it would be risky to change the relevant code, they generally defer it for the next release although the bug is important. Then it takes a long time to fix the bug. The follow-

ing developers' comments on bug #80,000 in JDT represents such a scenario:

"Will investigate during RC2 whether there's a low risk fix for this."

Two weeks later, the same developer commented: *"Sorry, too risky to touch at this point."*

5) Incomplete fix: This is considered as one of the common problems for taking a long time to fix bugs. Developers often miss corner cases while bug fixing and need to re-fix again until the problem is not fully solved. Here is a developer's comment regarding a bug fix for #38746 in JDT.

"The fix for this problem is not sufficiently robust. Please see Bug 75454 for more information as to how things can go wrong. Not only does the situation described there happen once, but it happens 60 times on start-up (10 minutes of start-up time)."

There are also lots of other reasons for incomplete bug fixes. For a comprehensive set of reasons for incomplete bug fixes, please refer to [20].

6) Importance was not realized until duplicate bugs were reported: We found many bugs where there were some activities around the bug for some time, which we observed by reading developers' comments. After that there was no activity for a long time. Then somebody pointed out some duplicate bugs and everybody started talking again; the bug was fixed quickly. The following comment on bug # 16114 in PDE represents such an example:

"this one is experienced by several users (see the duplicates for more info). Looks like something causes certain fragment files on the disk to be in use and when we try to delete the project (even with 'force' option), we fail. This leave us with a partially deleted project that causes more trouble after that."

7) Reproducibility: There are some bugs that take a long time to reproduce, but once it is reproduced, it is fixed quickly. For example, it took 1 year and 4 months to reproduce the bug #268833 in Platform but took only one day to fix. This problem often happens from low quality bug report, execution difference due to platforms, and so on.

There are some interesting bugs, where users know how to reproduce the bug but it happens for some special cases and thus needs some time to reproduce. For this kind of bug, if users submit the bug without concrete data, it takes a long time to reproduce the required data that developers need to analyze the bug. Therefore, the bug fix gets delayed although the responsible developer is ready to fix it. For example, to debug an "out of memory" problem in JDT (# 54831), developers needed a heap dump, which was not submitted when the bug

was posted. When the assigned developer asked for it, the bug reporter (who is actually another developer of JDT) was busy with his own work and could not submit the heap dump on time. As a result, it took a long time to fix the bug.

8) Schedule issue: Sometimes developers also feel that a bug is important to fix. However, they have more important bugs at hand that should be fixed earlier. Therefore, although the other bugs are important, they are generally deferred. For example, there is a blocker bug (#10800) in JDT that prevented users to put space in VM arguments. Blocker bugs are considered as the most severe bugs. However, such a severe bug was deferred due to scheduling issues. Certainly, other developers were not very happy about that. The following developers' comments illustrate the scenario more clearly.

"Can more explanation be given as to why this issue has been marked as LATER? Does this mean it will not be fixed any time soon? If so, I find it very unfortunate as this is a very serious bug and requires nasty work arounds. If not, then my apologies..."

In reply, the responsible team leader said: *"In this case, 'LATER' means probably not for the final 2.0 release (tentatively scheduled for sometime in May). Quite simply, this problem was not deemed as critical as a lot of other problems that need to be solved for 2.0. The debug committers have A LOT to do before 2.0. But the beauty of an open source project is that if someone feels strongly about a particular feature or bug, they can make a contribution. If you would like to contribute a fix, I would be happy to review it."*

Finally, it took more than two years to fix the bug.

9) Not aware of fix or reopened due to misunderstanding: We are not completely certain if these bugs are really long lived. In this category, some bugs perhaps fixed earlier but developers have not changed the status in Bugzilla. Therefore, the reporters or other developers were not aware of the fix. Later some other developers just closed the bugs mentioning that probably the bugs have been already fixed. Another case is that sometimes reporters misunderstood something and reopened a given bug again. But then some other developer clarified the mistakes the reporter was making and finally again marked it as `FIXED` and `RESOLVED`.

10) Infrequent use case: This kind of bug is important to fix considering their destruction ability. However, they are not too frequent use cases. Therefore, developers just defer it for next milestone. For example, due to the bug # 130874 in JDT, a user can lose his/her Java code template references. However, developer deferred it by making following comment.

“We should definitely fix this during 3.5. Too late for 3.4 and really not a very common case.”

11) Others: There are also other reasons for delay in bug fixing such as expert developers are on vacation, dependency on other bugs to be fixed, and various document fixing.

12) As-usual delay: We have not found any specific reasons for these bugs by analyzing developers’ comments and thus we considered them as as-usual delay. If there are some specific reasons (mentioned above) to make delay, it is highly likely that developers will discuss it like the other bugs. However, it is also possible that the fixes were deferred due to scheduling issues. The following comment for bug #149316 in JDT can be served as an example of as-usual delay.

“Thanks for the good examples, sorry for the wait. fixed > 20080422.”

One may think that since these bug-fixes were delayed without any specific reasons, they are probably not that important or relevant. However, it should be noted that for this investigation we analyzed only critical or major bugs. So they are already marked as important by the reporters or developers. Furthermore, since we ourselves are users of these software systems, we understand that many such bugs can be frustrating. The following represent two summaries of such bugs:

Bug # 26556 (JDT): “PDE Junit does not read plugin info from the plugins directory.”

Bug # 43153 (Linux): “Random SATA drives on PMPs on sata_sil24 cards not being detected at boot since 3.2/3.4.”

We encountered an interesting finding while analyzing the bug reports manually. We started our manual investigation with JDT and listed all the common reasons from there. We have not found any new common reason when analyzing the bug reports for PDE, Platform, or the Linux Kernel. Therefore, we believe that this is a comprehensive list of reasons for long lived bugs. It should be noted that these reasons are not mutually exclusive.

Reasons for long lived bugs are diverse. While problem complexity, problems in reproducing errors, and not understanding the importance of some of the bugs in advance are the common reasons, we observed there are many bug-fixes that were delayed without any specific reason.

Table 9: Reasons of Sampled Long Lived Bugs

Reason #	Bugs	Bug IDs
1)	13	113870 (PDE), 128563 (JDT), 241241 (Platform), 245008 (Platform), 247766 (PDE), 268833 (Platform), 278598 (PDE), 3022 (Linux), 5534 (Linux), 5637 (Linux), 9905 (Linux), 13484 (Linux), 38442 (Linux)
2)	3	3849 (JDT), 36204 (JDT), 133072 (PDE)
3)	16	3849 (JDT), 24951 (PDE), 36204 (JDT), 38746 (JDT), 40243 (JDT), 46407 (JDT), 67425 (JDT), 82850 (JDT), 99137 (JDT), 233643 (PDE), 233773 (Platform), 266651 (JDT), 273450 (Platform), 295200 (JDT), 38442 (Linux), 44161 (Linux),
4)	2	80000 (JDT), 102780 (JDT)
5)	6	1766 (JDT), 33035 (JDT), 36204 (JDT), 136135 (PDE), 3410 (Linux), 46171 (Linux),
6)	14	36204 (JDT), 46216 (JDT), 50735 (JDT), 109636 (JDT), 117698 (JDT), 156168 (JDT), 175226 (JDT), 224880 (Platform), 243894 (Platform), 257202 (Platform), 266651 (JDT), 267649 (Platform), 273450 (Platform), 278598 (PDE)
7)	11	1766 (JDT), 39222 (JDT), 54831 (JDT), 82850 (JDT), 83473 (JDT), 195183 (JDT), 262032 (Platform), 294650 (Platform), 298795 (Platform), 2979 (Linux), 31602 (Linux),
8)	7	3920 (JDT), 19251 (PDE), 46216 (JDT), 67425 (JDT), 224880 (Platform), 235572 (Platform), 277638 (Platform)
9)	12	6437 (JDT), 19248 (PDE), 28637 (JDT), 44035 (JDT), 61744 (PDE), 132333 (PDE), 158589 (PDE), 271373 (Platform), 14563 (Linux), 43981 (Linux), 45031 (Linux), 46161 (Linux)
10)	2	34033 (PDE), 130874 (JDT)
11)	8	12955 (JDT), 16686 (JDT), 20919 (PDE), 24951 (PDE), 29799 (PDE), 34399 (PDE), 231936 (PDE), 290324 (PDE)
12)	34	21100 (PDE), 26556 (JDT), 38288 (PDE), 39803 (JDT), 51862 (PDE), 89347 (JDT), 95288 (JDT), 97541 (JDT), 111419 (JDT), 128303 (PDE), 129689 (PDE), 149316 (JDT), 154823 (JDT), 175133 (JDT), 181954 (JDT), 209537 (JDT), 226595 (Platform), 234623 (Platform), 235554 (Platform), 236104 (Platform), 237025 (PDE), 238943 (JDT), 258952 (Platform), 262032 (Platform), 267173 (Platform), 275910 (Platform), 277638 (Platform), 279781 (Platform), 285101 (Platform), 5637 (Linux), 11509 (Linux), 38312 (Linux), 41682 (Linux), 43153 (Linux)

Table 10: Analysis of Bug Fixes

System	# Bugs	#Number of Files (NOF)						Med NOF	Max NOF
		1	2	3	4	5	>5		
JDT	223	80	41	37	10	15	40	2	47
CDT	185	50	32	19	18	8	58	3	237
PDE	105	34	13	13	12	4	29	3	211
Platform	740	349	114	72	47	36	122	2	91
Total	1253	513	200	141	87	63	249	-	-
(%)	-	40.94	15.96	11.25	6.94	5.03	19.87	-	-
Linux	171	117	28	12	4	3	7	1	19
WineHQ	609	231	203	70	37	29	39	2	51
GDB	33	0	10	2	1	0	20	7	222
Total	813	348	241	84	42	32	66	-	-
(%)	-	42.8	29.64	10.33	5.17	3.94	8.12	-	-

4.5. RQ5: What is the nature of bug fixes?

In this research question, we investigate the nature of bug fixes in terms of source code changes. To do this, first we identify the bug fixing changes in the source code for the long lived bugs. Using the *methodology* described in Section 3.3, we were able to identify fixed files for 280, 264, 154, and 1,109 bugs in JDT, CDT, PDE, and Platform respectively, and 171, 609, and 33 bugs in Linux, WineHQ, and GDB respectively. Then, we compute the number of changed files, number of hunks, and code churns for each bug-fix, as described in Section 3.2. These metrics are often used to get a rough idea about change effort, although understanding the actual change effort is difficult and depends additionally on the implemented algorithm and code complexity itself. Analogous to previous study results [30], our initial hypothesis was that *the required source code changes to fix most of the long lived bugs would be large*.

4.5.1. Changes at File Level

To understand the nature of bug fixes, we first analyze the source code changes in terms of number of changed files. As we stated in our initial hypothesis, we expected a large number of changed files for most long lived bugs. Surprisingly, from the results in Table 10, we see that for both Java and C projects, more than 40% of fixes involved only one source code file. This proportion varied from 27% to 47% among the Java projects, whereas it varied from 38% to 68% in C projects (except GDB). For only 30% of long lived bugs, the required changes spanned over more than three files in Java projects, whereas it was only 17% for C projects. From our results it is also noticeable that the maximum

number of changed files in each system is quite high. However, when we manually investigated such changes, we found they are often moving files from one directory to another, or adding a test suite to test a specific or multiple bugs.

4.5.2. Number of Hunks

Now we analyzed the changes in terms of hunk size to get more fine grained results. A large number of hunks indicates that developers needed to modify a lot of different places to fix the bugs. Figure 5 presents the number of bugs for each hunk size in Java projects. Our results show that 43% to 53% of bugs in Java projects was fixed by changing five hunks or less. In C projects, the proportion was even higher. 76% and 67% of bug fixes in the Linux Kernel and WineHQ respectively involved 5 hunks or less. More than 70% of long lived bugs for JDT and Platform was fixed within 10 hunks, whereas the numbers are 16 and 17 for PDE and CDT respectively. The median number of hunks for all long lived bugs is only 5 or less for all projects except GDB. Considering that, this is an overestimated measurement of source code changes since we have analyzed textual *diff* results (so changes in comments have been also counted), we believe the number of hunk is low. It should be noted that we presented all the bugs that involved more than 100 hunks in the graph at the end. Therefore, there is a spike at the end of some graphs.

4.5.3. Code Churn

We investigate further low level changes (at line level). Figure 6 presents the distribution of code changes in terms of code churns (defined in Section 3.2). From the figure, we see that a consider-

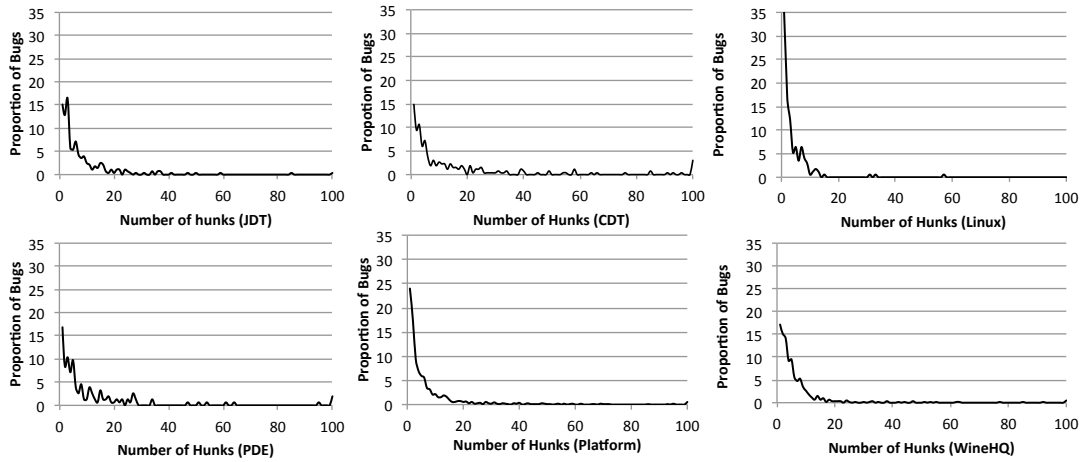


Figure 5: Number of Hunks Vs. Proportion of Bugs

able proportion of bugs required major changes. More specifically, for 23% of bugs, the value of code churn was more than 100. However, there are even a larger proportion of bugs that required changes less than 20 lines. For example, for 8-14% of bug fixes in Java projects, the value of code churn was from 1 to 5, for 7-12% the value of code churn was from 6 to 10, and for 9-13% the value of code churn was from 11 to 20. For the Linux Kernel and WineHQ, the proportion was even higher. 40% of long lived bugs in the Linux Kernel and 20% of bugs in WineHQ required changes in only 10 or less lines of code. Recalling Section 3.2, it should be noted that the code churn value of one line change is 2, whereas an addition or deletion of line is 1. Therefore, a value of code churn value of 10 may be changes in only five lines. We understand that some smaller bug fixes can be complex. But at the same time, we also stress that many long lived bugs could be fixed quickly through careful prioritization.

It should be noted that the overall bug fixing changes in GDB seems to be larger than other projects. However, as Table 10 shows, we were able to identify the bug fixing changes for only 33 bugs in GDB. Therefore, it is very difficult to draw any conclusion from GDB.

4.5.4. A Qualitative Analysis

Now we present three bug fixes to discuss how a simple fix can take a long time to be fixed. To select these example-fixes we considered the following criteria: i) they are from different subject systems, ii) they are important, i.e., critical or major (or high for linux kernel), and iii) their descriptions are concise enough to present in the paper.

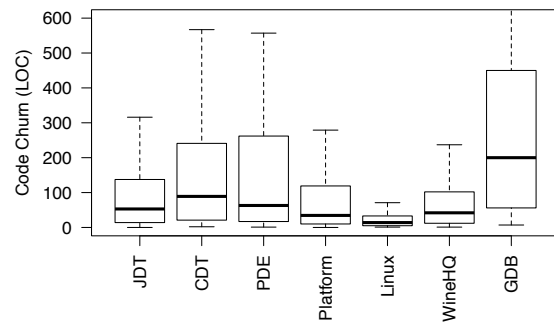


Figure 6: Code Churn of Long Lived Bugs

Bug # 38260 (SWT): This is a bug in the SWT component of Eclipse Platform. The bug was first reported as having a normal severity level. However, within 15 days, it was reconsidered to be critical. The following is the bug description provided by the reporter: “When I use the CCombo with the dialog, the dropdown list shown in the back of the Dialog.

I did as following.

1. Create one sample application with a Group.
2. The parent of group is one shell.
3. I create a dialog.
4. I put my application to this dialog. and change the Group’s parent as the dialog’s parent.

Now the above described pbm occurred.”

From the bug comments, we found that the bug was


```

--- a/bundles/org.eclipse.swt/Eclipse SWT Custom Widgets/common/org.eclipse.swt/custom/CCombo.java
+++ b/bundles/org.eclipse.swt/Eclipse SWT Custom Widgets/common/org.eclipse.swt/custom/CCombo.java
@@ -76,7 +76,7 @@ public CCombo (Composite parent, int style) {
    if ((style & SWT.READ_ONLY) != 0) textStyle |= SWT.READ_ONLY;
    if ((style & SWT.FLAT) != 0) textStyle |= SWT.FLAT;
    text = new Text (this, textStyle);
-   popup = new Shell (getShell (), SWT.NO_TRIM);
+   popup = new Shell (getDisplay(), SWT.NO_TRIM | SWT.ON_TOP);
    int listStyle = SWT.SINGLE | SWT.V_SCROLL;
    if ((style & SWT.FLAT) != 0) listStyle |= SWT.FLAT;
    if ((style & SWT.RIGHT_TO_LEFT) != 0) listStyle |= SWT.RIGHT_TO_LEFT;

```

Figure 7: Bug Fixing Changes for # 38260 in SWT Component of Platform

reproduced within three days and the actual problem was identified within a month. However, it took more than ten months to fix the bug. Figure 7 presents the changed code for fixing this bug and it was only a one line change.

Bug # 195183 (JDT): This is a major bug in the Debug component in JDT. The bug summary and description (condensed) are as follows: “*JavaClassPath.performApply() uses original instead of working copy causes NPE*”

“*Steps To Reproduce: I use JavaClassPath in my custom launch config and has some code like:”*[some code snippet] “*and as soon as performApply() is called isDefaultClasspath() fails since it is passed in a null as a launchconfig even though I passed in a newly created one. This worked fine in eclipse 3.2 and it seem the culprit is that wc.getOriginal() is used instead of just wc. Resulting in a NPE.*”

From the bug description, we can see that the bug report was very specific. The reporter clearly pointed out that there are some problems with the *getOriginal()* API. Interestingly, from the bug-fix (Figure 8), we found that only one line was changed and the change was the removal of the *getOriginal()* API. But by that time, more than two years had passed.

Bug # 3410 (Linux Kernel): This is a bug with high severity in the Linux Kernel. The bug summary and description provided by the reporter are as follows:

Summary: “passive mode is not left, once entered”

Description: “Steps to reproduce:

echo “xx:xx:low_value:xx:xx:” >/proc/acpi/thermal_zone/trip_points*

echo “xx:xx:high_value:xx:xx:” >/proc/acpi/thermal_zone/trip_points*

passive mode still active, even the temperature is far below the trip point

With this bug, there was an attempt to fix it on the same day the bug was reported. However, the fix was incomplete, which was identified on the next day. Although there was some discussion regarding the bug

around that time, it remained unfixed for almost one and half year. Figure 9 shows the second fix of the bug, which made the first fix complete. The patch involved changes in only two lines of code.

We believe that one year or more is too long time to fix the bugs like these examples, especially considering that they were considered as very important.

Unlike previous studies, we found that a bug surviving for a year or more does not necessarily mean that it requires a large fix. We found that 40% of long-lived bug fixes involved few changes in only one file.

5. Developers’ Survey

Since we have not found any specific reasons for a significant proportion of long lived bugs (classified as “as-usual delay”), and many long lived bugs involved small fixes, we believe that many such delays could have been avoided if developers could predict the severity and change effort in advance. To investigate what developers think about our assumption, we conducted a survey. We sent the survey to all the developers who contributed to the subject projects from year 2010. In total, we sent the survey to 104 developers. Among them, 38 developers do not use their email address anymore.

In the survey, we briefly explained our results and possible actions (e.g. predicting severity, change effort, prioritization, etc.), and asked developers if they agree with us or not. We also asked what actions can be taken to minimize the number of long lived bugs, if they think otherwise. We got responses from five developers. Three developers agreed that tool support to predict severity and change effort may be helpful. Furthermore, some of them think that understanding the impact of change is even more critical and often a major reason for delay in making bug fixing commits. One developer asked us if it is really possible to develop such tools instead of answering our survey. Another developer said

```

--- a/org.eclipse.jdt.debug.ui/ui/org/eclipse/jdt/debug/ui/launchConfigurations/JavaClasspathTab.java
+++ b/org.eclipse.jdt.debug.ui/ui/org/eclipse/jdt/debug/ui/launchConfigurations/JavaClasspathTab.java
@@ -272,7 +272,7 @@ public class JavaClasspathTab extends AbstractJavaClasspathTab {
    public void performApply(ILaunchConfigurationWorkingCopy configuration) {
        if (isDirty()) {
            IRuntimeClasspathEntry[] classpath = getCurrentClasspath();
-           boolean def = isDefaultClasspath(classpath, configuration.getOriginal());
+           boolean def = isDefaultClasspath(classpath, configuration);
            if (def) {
                configuration.setAttribute(IJavaLaunchConfigurationConstants.ATTR_DEFAULT_CLASSPATH,
                    (String)null);
                configuration.setAttribute(IJavaLaunchConfigurationConstants.ATTR_CLASSPATH,
                    (String)null);
            }
        }
    }
}

```

Figure 8: Bug Fixing Changes for # 195183 in Debug Component of JDT

```

--- a/drivers/acpi/processor_thermal.c
+++ b/drivers/acpi/processor_thermal.c
@@ -102,8 +102,8 @@ static int cpu_has_cpufreq(unsigned int cpu)

    struct cpufreq_policy policy;
    if (!acpi_thermal_cpufreq_is_init || cpufreq_get_policy(&policy, cpu))
-       return -ENODEV;
-       return 0;
+       return 0;
+       return 1;
}

```

Figure 9: Bug Fixing Changes for # 3410 in Linux Kernel

that duplicate bug report detection or improving bug report quality can be more useful.

6. Threats to Validity

This section discusses the validity and generalizability of our findings. In particular, we discuss Construct Validity, Internal Validity, and External Validity.

Construct Validity: We used two artifacts: bug reports from the bug tracking system and source code changes from the version history, which are generally well understood. We have used also well known metrics in our data analysis such as various time periods, the number of changed files, the number of hunks, code churns, which are straightforward to compute. Both the used dataset and version histories are also publicly available, which enable the replication of this study. Therefore, we argue for a strong construct validity.

Internal Validity: In our study, we relied on the information from the bug tracking system and version histories. However, the information in these systems may not be completely accurate. For example, a change request can be actually an enhancement but it could be misclassified as a bug [11]; the severity level associated with some bugs may not reflect the actual severity levels. Furthermore, a developer may commit a bug fixing change a long time after she actually fixed the bug. Similarly, a tester may change the bug status from FIXED to VERIFIED a long time after she actually verified the bug. Although it is very difficult to completely eliminate these threats, we performed extensive manual investiga-

tions and qualitative analyses, and provided many concrete examples throughout the paper to minimize these threats.

To delineate the common reasons of long lived bugs, we manually analyzed bug reports. There might have been some unintentional misinterpretations during the manual verification due to the lack of domain knowledge or the lack of useful contextual knowledge. However, we held extensive discussions to minimize this threat.

We used traditional heuristics to find mappings between bug fixing changes and associated bug reports. Although, in this way, we missed many bug fixing changes, the precision of our result is very high, which is important for our study. ReLink [29] is a more advanced algorithm that improves the recall quite a bit for finding the mappings. However, it also sacrifices some precision. In future, we would like to use ReLink to see how it affects our results.

The phenomena studied had yearly major releases. Systems with more frequent release cycles may well exhibit different phenomena, although there will still be long lived bugs. The number of release cycles and the lapse times for long lived bugs in this context are likely to be different.

External Validity: We have used seven subject systems in our experiment and all of them are open source projects. Although, they are very popular projects, our findings may not be generalizable to other open source projects or industrial projects. However, the Java dataset that we used has more than 165,000 bug reports, and the

C dataset we created contains more than 77,000 bug reports, which is large. Furthermore, the consistent findings from both Java and C projects make our results more generalizable for large open source projects. However, additional confidence could be achieved by adding more subject systems (both open source and industrial).

7. Related Work

The study of software bugs/faults has been an active research area for nearly two decades. Perry and Stieg [21] were among the first to analyze software faults in a large evolving software system. Since then, researchers analyzed various software artifacts relevant to bugs (e.g. bug report, bug fixing changes) to understand and to improve different steps (e.g. bug reporting, triaging, localizing, fixing) of the bug fixing processes.

Thung et al. [26] investigated when a bug should be reported. Bettenburg et al. [3] studied the qualities of a good bug report. In another study, Bettenburg et al. [4] investigated the extents and reasons of duplicated bug reports. They presented empirical evidences that duplicate bug reports are not necessarily bad. They often provide additional information, which is important for automatic bug triaging, bug assignment, and localization. Guo et al. [9] characterized and predicted which bugs get fixed. They noted that in addition to the importance of bugs, there are several other factors that affect whether a bug would get fixed, e.g., the reputation of bug reporters, influence of seniority, personal relations and trust, etc. In another study [10], the same authors investigated the reasons for bug reassignment (described in more detail in Section 4.3). Lamkanfi et al. [15] and Tian et al. [27] predicted the severity and priority of bugs respectively. Anvik et al. [2] and Shokripour et al. [25] proposed approaches for automatic bug assignment. Saha et al. [23] and Zhou et al. [32] proposed different approaches for automatic bug localization. To complement these studies, in this paper, we focused on long lived bugs to understand their characteristics and reasons.

The work closest to ours is the study of bug-fix time prediction, since these studies also identify the factors that are correlated to bug fixing time. Weiss et al. [28] considered the text (summary and description) in the bug report as the prime factor and used that to predict bug fix time. Panjer [19] observed that commenting activity, bug severity, product, component, and version are the most influential factors in predicting bug fix time. Giger et al. [8] found that the assigned developer, the bug reporter, and the month when the bug was reported have the strongest influence on the bug

fixing time. Zhang et al. [31] also found the same results for commercial projects. Anbalagan et al. [1] found a strong relationship between bug fixing time and the number of people participating in the bug report. Marks et al. [17] observed different results for different projects. They found that bug fixing time is important for Mozilla project, whereas, bug severity is the key for Eclipse.

While the aforementioned studies vary in terms of analysis and techniques used, some of the common approaches used in these studies are that researchers used various machine learning or data mining techniques to analyze the whole bug dataset in identifying the overall factors affecting bug fixing time. Bhattacharya and Neamtiu [5] pointed out that most attributes used by prior work do not correlate with bug-fix time when analyzed in isolation, and thus they emphasized on finding new attributes that correlate with bug-fix time in isolation. We stress that it is also important to analyze various kinds of bug-fixes in isolation to gain better insight about specific group of bugs. For example, Shihab et al. [24] studied and predicted reopened bugs, Park et al. [20] investigated supplementary bug-fixes, and Ngyuen et al. [18] analyzed recurring bug-fixes. In this study, we analyzed long lived bugs to advance empirical knowledge further regarding long-term delays in the bug fixing process.

There is another group of studies that investigated the actual source code changes for bug fixes to study bug fixing time. Canfora et al. [6] found relationships between different program constructs and bug survival time. For example, exception handling leads to low bug survival time. Zhang et al. [30] found that bug fixing time increases with the increase of code churns. In our study, we have also analyzed the source code changes for long lived bug-fix and showed that many long lived bugs involved only few changes in one file.

8. Conclusion

Bug fixing is a fundamental and critical activity in the software development and maintenance phases since buggy behavior may cause not only costly failures but also can affect user's overall experiences with the software product. In this paper, we showed that although the software development and maintenance processes have advanced a lot, there are still a significant number of bugs in each project that survive for more than a year. More than 90% of these long lived bugs may have affected users' normal working experience. The average bug assignment time was more than one year and the bug fix time after the assignment was another year on

average. When we analyzed the bug descriptions and the developers' comments around these bugs, we found that the reasons for long lived bugs are diverse. While problem complexity, problems in reproducing, and not understanding the importance of some of the long lived bugs in advance are the common reasons, we observed there are many bugs that were delayed without any specific reasons. Finally, by investigating the actual source code changes for these long lived bugs, we noted that a bug surviving for a year or more does not necessarily mean that it requires a large fix. In fact, we found 40% of long-lived bug fixes that involved few changes in only one file. Most importantly, all of the findings are consistent across the projects that we considered regardless of domains or programming languages.

In summary, our results indicate that the overall bug fixing time of many, if not all, long lived bugs can be reduced through careful prioritization, and by predicting their severity, change effort, and change impact. Our findings also indicate that although there are a number of tools for supporting bug triaging and fixing (e.g. automatic bug assignment, bug fix time prediction), we appear to realize very few benefits from them. There may be two possible reasons: i) developers are not aware that these tools exist, or ii) the tools do not meet developers needs or expectations. In the future, we plan to conduct a developers survey to understand the reasons for this phenomenon. We believe all of these findings together will play an important role in developing new and more effective approaches for bug triaging as well as improving the overall bug fixing process.

Acknowledgement: We thank Ahmed Lamkanfi at the University of Antwerp for extending their bug Java dataset significantly for our study, and Julia Lawall at Inria/LIP6/UPMC/Sorbonne University, France, for her help creating C dataset. This research was supported in part by NSF Grants CCF-0820251 and CCF-0845628.

- [1] P. Anbalagan and M. Vouk. On predicting the time taken to correct bug reports in open source projects. In *Proceeding of the International Conference on Software Maintenance*, pages 523–526, 2009.
- [2] J. Anvik, L. Hiew, and G. C. Murphy. Who should fix this bug? In *Proceeding of the International Conference on Software Engineering*, pages 361–370. ACM, 2006.
- [3] N. Bettenburg, S. Just, A. Schröter, C. Weiss, R. Premraj, and T. Zimmermann. What makes a good bug report? In *Proceeding of the International Symposium on the Foundations of Software Engineering*, pages 308–318. ACM, 2008.
- [4] N. Bettenburg, R. Premraj, T. Zimmermann, and S. Kim. Duplicate bug reports considered harmful...really? In *Proceeding of the International Conference on Software Maintenance*, pages 337–345, 2008.
- [5] P. Bhattacharya and I. Neamtii. Bug-fix time prediction models: can we do better? In *Proceeding of the Working Conference on Mining Software Repositories*, pages 207–210. ACM, 2011.
- [6] G. Canfora, M. Ceccarelli, L. Cerulo, and M. Di Penta. How long does a bug survive? an empirical study. In *Proceeding of the International Conference on Software Maintenance*, pages 191–200. IEEE Computer Society, 2011.
- [7] T.-H. Chen, M. Nagappan, E. Shihab, and A. E. Hassan. An empirical study of dormant bugs. In *Proceedings of the 11th Working Conference on Mining Software Repositories*, pages 82–91. ACM, 2014.
- [8] E. Giger, M. Pinzger, and H. Gall. Predicting the fix time of bugs. In *Proceeding of the International Workshop on Recommendation Systems for Software Engineering*, pages 52–56. ACM, 2010.
- [9] P. J. Guo, T. Zimmermann, N. Nagappan, and B. Murphy. Characterizing and predicting which bugs get fixed: an empirical study of microsoft windows. In *Software Engineering, 2010 ACM/IEEE 32nd International Conference on*, volume 1, pages 495–504. IEEE, 2010.
- [10] P. J. Guo, T. Zimmermann, N. Nagappan, and B. Murphy. Not my bug! and other reasons for software bug report reassignments. In *Proceedings of the ACM 2011 conference on Computer supported cooperative work*, pages 395–404. ACM, 2011.
- [11] K. Herzig, S. Just, and A. Zeller. It's not a bug, it's a feature: How misclassification impacts bug prediction. In *Proceedings of the 2013 International Conference on Software Engineering*, pages 392–401. IEEE Press, 2013.
- [12] N. Japkowicz. Learning from imbalanced data sets: A comparison of various strategies. In *Proceeding of the AAAI Workshop on Learning from Imbalanced Data Sets*, pages 10–15. AAAI, 2000.
- [13] S. Kim, E. J. Whitehead, Jr., and Y. Zhang. Classifying software changes: Clean or buggy? *IEEE Transaction on Software Engineering*, 34(2):181–196, Mar. 2008.
- [14] S. Kim, T. Zimmermann, K. Pan, and E. J. Whitehead. Automatic identification of bug-introducing changes. In *Proceeding of the Automated Software Engineering*, pages 81–90. IEEE, 2006.
- [15] A. Lamkanfi, S. Demeyer, E. Giger, and B. Goethals. Predicting the severity of a reported bug. In *Proceeding of the Working Conference on Mining Software Repositories*, pages 1–10, 2010.
- [16] A. Lamkanfi, J. Pérez, and S. Demeyer. The eclipse and mozilla defect tracking dataset: a genuine dataset for mining bug information. In *Proceeding of the Working Conference on Mining Software Repositories*, pages 203–206. IEEE Press, 2013.
- [17] L. Marks, Y. Zou, and A. E. Hassan. Studying the fix-time for bugs in large open source projects. In *Proceeding of the Promise*, pages 11:1–11:8. ACM, 2011.
- [18] T. T. Nguyen, H. A. Nguyen, N. H. Pham, J. Al-Kofahi, and T. N. Nguyen. Recurring bug fixes in object-oriented programs. In *Proceeding of the International Conference on Software Engineering*, pages 315–324, 2010.
- [19] L. D. Panjer. Predicting eclipse bug lifetimes. In *Proceeding of the Working Conference on Mining Software Repositories*, pages 29–32. IEEE Computer Society, 2007.
- [20] J. Park, M. Kim, B. Ray, and D.-H. Bae. An empirical study of supplementary bug fixes. In *Proceeding of the Working Conference on Mining Software Repositories*, pages 40–49, 2012.
- [21] D. E. Perry and C. S. Stieg. Software faults in evolving a large, real-time system: A case study. In *Proceeding of the ESEC*, pages 48–67, 1993.
- [22] R. K. Saha, S. Khurshid, and D. E. Perry. An empirical study of long lived bugs. In *Proceeding of the IEEE CSMR-18/WCRE-21 Software Evolution Week*, pages 144–153. IEEE, 2014.
- [23] R. K. Saha, M. Lease, S. Khurshid, and D. E. Perry. Improving bug localization using structured information retrieval. In *Proceeding of the Automated Software Engineering*, pages 345–

- 355, 2013.
- [24] E. Shihab, A. Ihara, Y. Kamei, W. Ibrahim, M. Ohira, B. Adams, A. Hassan, and K.-i. Matsumoto. Studying re-opened bugs in open source software. *Empirical Software Engineering*, 18(5):1005–1042, 2013.
- [25] R. Shokripour, J. Anvik, Z. M. Kasirun, and S. Zamani. Why so complicated? simple term filtering and weighting for location-based bug report assignment recommendation. In *Proceeding of the Working Conference on Mining Software Repositories*, pages 2–11, 2013.
- [26] F. Thung, D. Lo, L. Jiang, Lucia, F. Rahman, and P. Devanbu. When would this bug get reported? In *Proceeding of the International Conference on Software Maintenance*, pages 420–429, 2012.
- [27] Y. Tian, D. Lo, and C. Sun. Drone: Predicting priority of reported bugs by multi-factor analysis. In *Software Maintenance (ICSM), 2013 29th IEEE International Conference on*, pages 200–209. IEEE, 2013.
- [28] C. Weiss, R. Premraj, T. Zimmermann, and A. Zeller. How long will it take to fix this bug? In *Proceeding of the Working Conference on Mining Software Repositories*, pages 1–8, 2007.
- [29] R. Wu, H. Zhang, S. Kim, and S.-C. 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. ACM, 2011.
- [30] F. Zhang, F. Khomh, Y. Zou, and A. E. Hassan. An empirical study on factors impacting bug fixing time. In *Proceeding of the International Conference on Software Maintenance*, pages 225–234. IEEE Computer Society, 2012.
- [31] H. Zhang, L. Gong, and S. Versteeg. Predicting bug-fixing time: an empirical study of commercial software projects. In *Proceeding of the International Conference on Software Engineering*, pages 1042–1051. IEEE Press, 2013.
- [32] J. Zhou, H. Zhang, and D. Lo. Where should the bugs be fixed? - more accurate information retrieval-based bug localization based on bug reports. In *Proceeding of the International Conference on Software Engineering*, pages 14–24, 2012.