

Bench4BL: Reproducibility Study on the Performance of IR-Based Bug Localization

Jaekwon Lee
SnT, University of Luxembourg
Luxembourg
jaekwon.lee@uni.lu

Dongsun Kim
SnT, University of Luxembourg
Luxembourg
dongsun.kim@uni.lu

Tegawendé F. Bissyandé
SnT, University of Luxembourg
Luxembourg
tegawende.bissyande@uni.lu

Woosung Jung
Seoul National University of
Education
Seoul, South Korea
wsjung@snue.ac.kr

Yves Le Traon
SnT, University of Luxembourg
Luxembourg
Yves.LeTraon@uni.lu

ABSTRACT

In recent years, the use of Information Retrieval (IR) techniques to automate the localization of buggy files, given a bug report, has shown promising results. The abundance of approaches in the literature, however, contrasts with the reality of IR-based bug localization (IRBL) adoption by developers (or even by the research community to complement other research approaches). Presumably, this situation is due to the lack of comprehensive evaluations for state-of-the-art approaches which offer insights into the actual performance of the techniques.

We report on a comprehensive reproduction study of six state-of-the-art IRBL techniques. This study applies not only subjects used in existing studies (old subjects) but also 46 new subjects (61,431 Java files and 9,459 bug reports) to the IRBL techniques. In addition, the study compares two different version matching (between bug reports and source code files) strategies to highlight our observations related to performance deterioration. We also vary test file inclusion to investigate the effectiveness of IRBL techniques on test files, or its noise impact on performance. Finally, we assess potential performance gain if duplicate bug reports are leveraged.

CCS CONCEPTS

• **Software and its engineering** → **Software testing and debugging; Empirical software validation; Software evolution; Maintaining software;**

KEYWORDS

Reproducibility studies, bug localization, information retrieval

ACM Reference Format:

Jaekwon Lee, Dongsun Kim, Tegawendé F. Bissyandé, Woosung Jung, and Yves Le Traon. 2018. Bench4BL: Reproducibility Study on the Performance of IR-Based Bug Localization. In *Proceedings of 27th ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA'18)*. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3213846.3213856>

1 INTRODUCTION

In software development and maintenance, debugging constitutes one of the most costly activities [5, 42, 50]. To alleviate debugging costs, the research community is striving to produce approaches and tools for automating various tasks. Among these tasks, bug localization [16, 25, 34, 35], i.e., the process of locating program elements which contain bugs leading to abnormal behavior, is a well-studied field. Bug localization leverages information such as issue reports [4] and stack traces describing a bug and its symptom, as well as possible reproduction steps, to identify which source code files, classes, functions, or code chunks are relevant to the bug.

The literature of bug localization proposes several approaches which leverage Information Retrieval (IR) techniques [7, 9, 15, 27, 28, 38] to identify potential bug locations by processing textual bug reports and source code files. Such approaches produce a ranked list of files where highly ranked files are supposed to be most likely to contain the reported bug. As Parnin and Orso pointed out, based on a developer survey, developers are strongly sensitive to the performance of such debugging tools: they do not find the tool useful when it does not help pinpoint the root cause of the bug early in the output ranked list [33].

State-of-the-art IR-based bug localization (IRBL) approaches are reported in the literature with increasingly high accuracy scores. Nevertheless, their benchmarking has not yet reached the level of maturity that is now required in several other sub-fields of software engineering. Our work contributes towards this maturity by (1) exploring a comprehensive evaluation of existing techniques with new large-scale datasets, (2) performing a reproducibility study on the performance of state-of-the-art techniques, and (3) investigating potential directions for improving IRBL performance (e.g., by matching the version of a bug report with source code files, assessing the impact of test file inclusion, or applying IRBL on duplicate bug reports).

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

ISSTA'18, July 16–21, 2018, Amsterdam, Netherlands

© 2018 Copyright held by the owner/author(s). Publication rights licensed to the Association for Computing Machinery.

ACM ISBN 978-1-4503-5699-2/18/07...\$15.00

<https://doi.org/10.1145/3213846.3213856>

Overall, this paper makes the following contributions:

- **A Reproducibility study with new subjects:** We examine potential overfitting scenarios by investigating the difference in performance when applied to the limited set of common (old) subjects used by most studies [37, 44–46, 52, 53] in the literature, and when applied to new subjects.
- **An Empirical Assessment of the impact of IRBL execution configurations:** Our study varies several different execution strategies of IRBL techniques to show the performance of the techniques in different evaluation environments such as version matching strategies and test file inclusion.
- **A Comprehensive benchmark for bug localization:** Finally, we package our datasets and detailed reproduction study results of state-of-the-art approaches into a new benchmark, Bench4BL¹.

Our findings include:

- (1) The performance of IRBL techniques is actually higher than what is reported in the literature. We further show that experiments in the literature are also flawed since they were often performed on out-of-sync versions of the source code.
- (2) While the differences in performance among the state-of-the-art techniques become less apparent with our large dataset, there is a high variability of performance across projects.
- (3) In contrast to the common belief of the community [45], excluding test files is not necessary for the evaluation of IRBL techniques.
- (4) Merging duplicate reports may achieve the performance gain.
- (5) There is still room to improve the performance of IRBL.

The remainder of this paper is organized as follows. After explaining the background of IR-based bug localization and the motivation of this paper in Section 2, we describe the design of our study in Section 3. Section 4 illustrates, interprets the study results, and provides take-away messages for prospective users and researchers of IRBL. Section 5 discusses additional issues. After surveying the related work in Section 6, we conclude with directions for future research in Section 7.

2 BACKGROUND & MOTIVATION

Bug/Fault localization is a software maintenance activity in which a maintainer examines symptom information submitted by users, or obtained by running test cases, to locate suspicious resources (e.g., source code, configuration or data files) in a project. There are two leading research directions on bug (fault) localization, namely spectrum-based fault localization [1, 2, 10, 47] and IR-based bug localization [37, 44–46, 52, 53]. In the former, probable locations of faults are identified by computing ranking metrics, generally based on similarity coefficients and statistical techniques, on succeeding and failing test execution traces. In the latter, on the other hand, approaches only leverage source code information and bug report text to identify suspicious files using IR techniques such as Latent Dirichlet Allocation (e.g., [25]), Vector Space Model (e.g., [53]), Latent Semantic Analysis (e.g., [8]), Clustering (e.g., [24]).

IRBL techniques [37, 44, 46, 52, 53] are commonly known to exhibit low cost and to be scalable in practice since they require

static information and, besides the bug report and the source code, they have no other external dependencies. The basic intuition for adopting IR technique for bug localization is that bug reports and project resources share common words. In other words, we can regard that a certain source file is suspicious if it contains words also written in a given bug report. According to the survey of feature location studies [13], IRBL is one of the tasks to locate feature locations; its type of analysis is ‘textual’, the input is a natural language query, the data source is source code, and the output is a subset of source code files in a project.

2.1 Performance Metrics

When evaluating IRBL techniques, many studies first collect ground truth data, which consist of bug reports already fixed by developers and the corresponding change set (i.e., files changed to fix a bug). Then, they apply the IRBL technique on this dataset and compare output results against the ground truth. The followings are common metrics used in most studies to assess an approach’s performance:

- **Precision:** Also more accurately referred to as *Precision@k*, this metric represents an estimation of how many files are correctly recommended within given top k files. It is computed as follows:

$$P(k) = \frac{\# \text{ of buggy files in top } k}{k} \quad (1)$$

- **Recall:** Also more accurately referred to as *Recall@k*, this metric estimates how many files are correctly recommended within given top k files over the actually fixed files by a developer for a given bug report. It is computed as follows:

$$R(k) = \frac{\# \text{ of buggy files in top } k}{\# \text{ of actually fixed files}} \quad (2)$$

- **Average Precision (AP):** The average precision of a given bug report aggregates precision values of several positively recommended files as:

$$AP = \sum_{i=1}^N \frac{P(i) \cdot pos(i)}{\# \text{ of positive instances}} \quad (3)$$

where N is the number of ranked files by an IRBL technique, i is a rank in the ranked list of recommended files. $pos(i)$ indicates whether the i -th file in the ranked list is a buggy file (i.e., $pos(i) \in \{0, 1\}$). For example, $AP = 0.5$, for a bug report with k actually fixed files, implies that an IRBL technique can make correct recommendations with 50% of probability within top k recommendations.

- **Mean Average Precision (MAP):** The MAP is computed by taking the mean value of AP values across all bug reports:

$$MAP = \frac{1}{M} \sum_{j=1}^M AP(j) \quad (4)$$

where M is the number of all reports. $AP(j)$ is the average precision of report j . If $MAP = p_{map}$, at least one file is likely to be a correct recommendation for every $\frac{1}{p_{map}}$ files in the ranked list.

- **Mean Reciprocal Rank (MRR):** This measure computes the mean value of the position of the first buggy file in the ranked list recommended by an IRBL technique, following this equation:

$$MRR = \frac{1}{M} \sum_{i=1}^M \frac{1}{f\text{-rank}_i} \quad (5)$$

where M is the number of all bug reports and $f\text{-rank}_i$ means the position of the first buggy file in the ranked list for the i -th

¹<https://github.com/exatoo/Bench4BL>

bug report. For example, assuming $MRR = 0.5$, an IRBL technique can locate at least one correct file to fix within top two recommendations (i.e., $\frac{1}{MRR} = \frac{1}{0.5} = 2$).

2.2 Motivation

We motivate this work as a contribution to the community for boosting the research on IR-based bug localization. We have indeed identified two key issues with the current state of research in IRBL, which justify a thorough reproducibility study for understanding which techniques indeed perform well and how the performance of such techniques can be realistically improved:

Performance of IRBL techniques is not solidly established: IRBL research has produced a number of state-of-the-art approaches in recent years. Unfortunately, performance assessments often focused on the same old and limited set of projects (such as old versions of JDT, AspectJ, etc.). Wen et al. [45] have also raised several threats to validity on the recorded performance: there could be some overfitting in the models due to some subject selection bias; and given the limited datasets, the data could be of poor quality to draw conclusions in the sense that they may not include interesting bug cases.

There are few perspectives on how to improve existing techniques: Most approaches in the literature focus on tuning the algorithms to gain fractions of precision performance points. In this work, we attempt to show that there is more room to demonstrate or improve the performance of techniques by consolidating the input data. For example, we investigate the consistency of project version with the bug report when applying IRBL techniques since bug reports indeed reflect the functioning of the program at a given time. We also examine whether including test files when identifying buggy files would help bug localization with IRBL techniques. In addition, duplicate bug reports may help us identify which features in bug reports are effective for IRBL techniques.

3 STUDY DESIGN

Our reproduction study explores different experimental scenarios to address the following research questions:

- **RQ1:** To what extent do IRBL techniques perform on up-to-date subjects?
- **RQ2:** What is the impact of version matching on the performance of IRBL techniques?
- **RQ3:** To what extent are IRBL techniques sensitive to the inclusion of test code files?
- **RQ4:** What potential performance gain can be reached by leveraging duplicate bug reports?

The objective of RQ1 is to (1) reproduce the results of IRBL techniques with the outdated subjects from the literature and (2) produce another set of results of the techniques with larger and recent subjects collected from active projects. With these two objectives, we can eventually establish whether some performance results recorded in the literature were actually over-fitted to a specific set of subjects.

Throughout RQ2 and RQ3, we investigate how IRBL techniques perform under different execution settings. Our study first (RQ2) compares two different version matching strategies: (1) assuming that an IRBL technique identifies files in a *single* version of the target

program for all bug reports and (2) each bug report corresponds to a specific version of the program. Second (RQ3), we give two different search spaces to IRBL techniques: source code files with and without test files.

RQ4 aims at estimating the usefulness of duplicate reports for IRBL. Since duplicate reports contain different descriptions for a single bug, they could help improve performance.

To answer these research questions, we selected IRBL techniques whose implementations were readily available and usable (cf. Section 3.1). We then collected data from a large and diversified set of open source projects (cf. Sections 3.2 and 3.3); the dataset is publicly available at Bench4BL (<https://github.com/exatoa/Bench4BL>). In Sections 3.4 and 3.5, we define execution strategies for IRBL techniques to evaluate them in various configurations.

3.1 IRBL Techniques

We identified six recent state-of-the-art techniques that target Java projects:

- (2012) - **BugLocator** [53] leverages similar bug reports that have been previously fixed and relies on revised Vector Space Model (rVSM) for the recommendation.
- (2013) - **BLUiR** [37] extracts code entities such as classes, methods, and variable names from bug reports and leverages them to localize files.
- (2014) - **BRTracer** [46] analyzes stack traces shown in bug reports to improve bug localization accuracy.
- (2014) - **Amalgam** [44] utilizes revision history in addition to similar reports and code entities.
- (2015) - **BLIA** [52] combines information such as similar reports, revision history, code entities, and stack trace information all together to improve the performance of IRBL.
- (2016) - **Locus** [45], the most recent technique, leverages code change information.

We select the above six techniques because they are frequently adopted to perform comparative studies against each other. Locus is one of the most recent techniques and its evaluation compares with BRTracer, BLUiR, and Amalgam. The work of BLIA compares its performance with BugLocator, BLUiR, BRTracer, and Amalgam. In the work of Amalgam, the authors evaluated its performance comparing with BugLocator and BLUiR. The performance of BRTracer and BLUiR are compared with BugLocator in their work.

Default parameters: We use the default parameters in the literature by the approach authors for their assessment experiments. As recent IR-based studies [12, 32] pointed out, configuration and parameters of each IRBL technique might have a significant impact on the performance of bug localization. Since our goal is to reproduce the techniques and their results without any modification, we use the same configuration and parameters specified in each paper. If there are multiple configurations or parameters are available, we take the best one that was reported to outperform other settings. We provide more details in Section 6 for these techniques.

Specifically, we leverage the following configurations and parameters. For BugLocator, our setup includes the use of the revised Vector Space Model (rVSM), similar bug reports, the weighting factor α (for similar reports) of 0.3, and a logistic function as the length function. In case of BLUiR, we utilize TFIDF with term frequency

parameter $k_1 = 1.0$ and document scaling factor $b = 0.3$, select the Krovetz stemmer, enable exact identifier name indexing and code structure modeling, and leverage similar bug reports. The execution setup for BRTracer leverages rVSM and similar bug reports. For Amalgam, we use the same setup with BugLocator since the technique is based on BugLocator. This technique takes two additional parameters: version history length $k = 15$ and composer component ratio $b = 0.3$ by default. In case of BLIA, its configuration uses rVSM and three control parameters $\alpha = \beta = 0.2$, $k = 30$ by default as specified in [52]. For Locus, we use VSM and control parameter $\lambda = 0.5$ by default. Refer to the corresponding papers of the six techniques for more details of each parameter.

In addition, we run the above six techniques with the following environment: the machine that has an eight-core with 3.6GHz Intel processors and 16GB memory, and the 64-bit Ubuntu 16.04.

3.2 Subjects and Data Extraction

The first part of Table 1 lists 46 open source projects, denoted as *new subjects* collected for our study and which are (1) written in Java, (2) with publicly available bug reports, and (3) having at least 20 source code files in one of its version. We applied these criteria to the projects of Apache², Spring³, and JBoss⁴; we selected those software communities since they are the largest ones with projects written by Java and well-managed together with issue tracking systems such as Jira⁵.

In addition, we consider projects commonly used in the literature of IRBL [37, 44, 46, 52, 53] as listed in the second part of Table 1 as *old subjects*. Since the previous studies have used slightly different data sets from each other, we collect the union set of reports and source files as long as those data are available. Note that we have collected bug reports and source code files for Eclipse JDT and PDE instead of Eclipse itself since the repository of Eclipse is now separately managed for each sub-project [45]. Although additional benchmarks are available in [11, 14], we focus on existing subjects widely-used in IRBL studies.

In accordance with the settings of previous studies [37, 44–46, 52, 53], we collect only source code files (i.e., *.java files) from each subject. As shown in Table 1, 61,431 files are collected from 46 new subjects and 19,475 files from five old subjects. Note that the number of files is computed for the version with the maximum number of source files for each subject: the actual number of files could be different for each specific version (cf. Section 3.4 for multi-version collection details).

For each subject in Table 1, we collect bug reports from the corresponding issue tracking systems (ITS). Among all reports available in an ITS, we only select bug reports that are explicitly classified as “Bug” by developers and marked as **FIXED** with explicit file changes, yielding 9,459 reports from the 46 new subjects and 558 reports from the five old subjects. From each collected bug report, we extract information only available when it initially submitted, i.e., summary (or title), description, and reporter. We leave out extra information such as comments.

Table 1: Projects used in our study.

Group	Subject	# Source files (Max)	# Major versions	# Bug reports	# Duplicate reports
New subjects					
Apache	CAMEL	14,522	60	1,469	50
	HBASE	2,714	70	836	80
	HIVE	4,651	21	1,241	270
	CODEC	115	9	42	2
	COLLECTIONS	525	7	92	16
	COMPRESS	254	15	113	9
	CONFIGURATION	447	11	133	4
	CRYPTO	82	1	8	0
	CSV	29	3	14	0
	IO	227	13	91	7
	LANG	305	16	217	23
MATH	1,617	15	245	8	
WEAVER	113	1	2	0	
JBoss	ENTESB	252	3	47	1
	JBMETA	858	5	26	1
	ELY	68	3	25	1
	SWARM	727	6	58	1
	WFARQ	126	1	1	0
	WFCORE	3,598	16	361	8
	WFLY	8,990	11	984	27
WFMP	80	1	3	0	
Spring	AMQP	408	33	108	3
	ANDROID	305	2	11	0
	BATCH	1,732	33	432	3
	BATCHADM	243	4	20	0
	DATACMNS	604	33	158	15
	DATAGRAPH	848	4	60	2
	DATAJPA	330	38	147	11
	DATAMONGO	622	40	271	19
	DATAREDIS	551	17	49	1
	DATAREST	414	23	132	12
	LDAP	566	5	53	1
	MOBILE	64	3	11	0
	ROO	1,109	15	714	72
	SEC	1,618	42	541	68
	SECOAUTH	726	7	101	0
	SGF	695	19	107	1
	SHDP	1,102	9	45	0
	SHL	151	3	11	1
	SOCIAL	212	4	15	0
	SOCIALFB	253	5	15	0
SOCIALLI	180	1	4	0	
SOCIALTW	153	5	8	0	
SPR	6,512	12	130	73	
SWF	808	20	134	5	
SWS	925	25	174	12	
Sub-total		61,431	690	9,459	807
Old subjects					
Old subjects	AspectJ	6,485	1	286	35
	JDT	6,842	1	94	24
	PDE	5,273	1	60	20
	SWT	484	1	98	57
	ZXing	391	1	20	-
Sub-total		19,475	5	558	136
Total		80,906	695	10,017	943

Duplicate bug reports: We searched the ITS for duplicate bug reports of all bug reports. Relationships between master bug reports (i.e., those marked as fixed) and duplicate bug reports are $1 \leftarrow n$ (i.e., individual duplicate reports point to one master report independently). In some cases, the relationships are transitive: master report (M) \leftarrow first duplicate (D1) \leftarrow second duplicate (D2). We regard them as a group of duplicate reports for the master report and dealt with as pairs of $\langle M, D^* \rangle$. Note that we cannot collect duplicate bug reports for ZXing as shown in Table 1 because its ITS for old reports is no longer available.

In this work, to avoid introducing bias related to data quality, we do not leverage any of the advanced approaches [17, 36, 40, 41] to detect duplicate reports. Instead, we use the link explicitly specified in each bug report. This method may miss potential duplicate bug

²Apache, <http://www.apache.org/>

³Spring, <https://spring.io/>

⁴JBoss, <http://www.jboss.org/>

⁵Jira, <https://www.atlassian.com/jira>

reports but there might be no false positives. Overall, we collected 807 duplicate report pairs⁶ for the new subjects and 136 pairs for the old subjects. Note that duplicate reports often are not attached with fixed files since those reports tend to be classified as duplicates before linking source files to change. Thus, we assume, when using for our experiments, those duplicate reports have the same set of files that the master reports have changed to fix the bugs.

3.3 Bug Linking

To build experimental ground truth we must link every fixed bug report with the associated files where the fix is performed. The correct link information is crucial for evaluating bug localization techniques since the techniques rely on the link information in their performance evaluation. The researchers of most IRBL studies also collected the link information between bug reports and source code commits to evaluate the techniques.



Figure 1: Examples of two bug linking methods.

Although there exist other more advanced approaches to bug linking (e.g., ReLink [49], MLink [31], and RCLinker [21]), we do not consider them to avoid their performance limitations [6]. Wu et al. [49] reported relatively high precision results for half of the subjects, but their technique showed lower precision (0.682–0.858) for another half. Although MLink [31] achieved a higher precision than ReLink, its evaluation was performed a limited number/size of subjects. This limitation is also discussed in [21] with an evaluation of additional subjects; its precision values are ranged from 0.3 to 0.785 (average: 0.564). RCLinker [21] achieves 0.509 precision on average. Due to the low precision, we cannot adopt the bug linking techniques in our study. The six IRBL techniques used in our study did not leverage the above approaches to avoid many false positives.

In this study, we only use explicit links in bug reports and commit log messages since our goal is to avoid false positives (a link between a report and file but they are not actually relevant) and to allow false negatives (missing links). Developers often put link information in bug reports and commit logs to keep track of bugs and corresponding source code changes [20]. They use explicit IDs to identify bug reports and code commits. Figure 1a shows an example of a commit ID recorded in a bug report. Most issue tracking systems provide a feature to explicitly designate commit ID (with a hyperlink). In case of commit logs, developers manually write bug report ID in its message in a specific format as shown in Figure 1b.

We extracted link information in bug reports, which is explicitly provided by an issue tracking system. The subjects listed in Table 1 are all managed by Jira. Thus, we collected bug reports written in

⁶The number of pairs and duplicate reports are different since some duplicate reports can be grouped by a master report as explained in Section 3.2.

Jira’s format and we identified commit IDs after parsing them. If a bug report contains a designated commit ID, we link the source code files (in this study, “*.java” files) to the bug report. When multiple commit IDs are recorded in a single bug report, we use the latest commit ID to link source code files.

When identifying bug report IDs in a commit log, we use a regular expression. Bug report IDs of the subjects are often recorded in a format of “PROJECT-###”, which can be detected by a regular expression. If there are multiple IDs, we use the first occurrence as the bug report to link source code files in the commit.

3.4 Version Matching

Bug reports are submitted by users in reference to specific project versions. In most cases, the target version of the bug is the most recently released version at the time of submission. When attempting to localize potential buggy files relevant to a given bug report, it is necessary to clearly define a search space of the project files. If version information is not correctly specified, the result of a bug localization technique becomes useless.

In this study, we apply two different strategies of version matching: *single version matching* vs. *multiple version matching*. In the single version matching strategy, a bug localization technique assumes that the search space of potential source code files is the latest version of a software project as shown in Figure 2a. This strategy makes the search space simple. Thus, most bug localization studies have used this strategy in their evaluation. The work of the six techniques adopted in this study also used this strategy in their own evaluation. With multiple version matching, a bug localization technique constructs a set of files of the version specified in a bug report as the search space. Thus, different bug reports correspond to different sets of source code files as shown in Figure 2b.

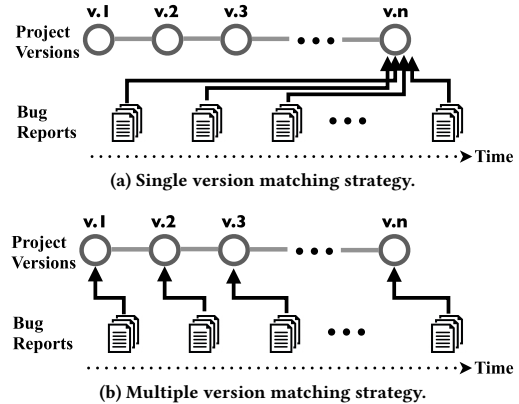


Figure 2: Two version matching strategies used in this study.

Since previous studies did not implement multiple version matching, we perform the following procedure to determine which version corresponds to each report: (1) we collect all tagged version releases for each subject, then (2) we check, for each bug report, that it specifies the release version(s) that are concerned by this bug. When several versions are listed, we consider the oldest one. If there is no version specified in a bug report, we ignore the report. Note that the numbers in “# Bug reports” in Table 1 indicate the number of bug reports with explicit version information. For old subjects, version matching is not applied since many major

releases associated with reports used in previous studies are no longer available.

3.5 Test File Inclusion

This study uses two different strategies for search space construction: (1) with test code files and (2) without test files. Most IRBL techniques select the first strategy; a bug localization technique scans all source code files (e.g., `*.java`) in a target project to create a search space. Among the six techniques adopted in this study, only Locus [45] does not include test code files. The technique excludes test code files since the authors assumed that test files are not the target of bug localization. Specifically, the authors stated that the inclusion of test files may cause bias or noise since some test files contain specific bug identifiers and they are created after the bug is fixed. Note that this situation happens only when using single version matching.

To prepare the search space without test code files, we use the following procedure: (1) filter out files if its path contains `.../test/...` or `.../tests/...`, and (2) exclude source code files if its file name ends with `... Test.java`. Although this procedure may cause false positives (i.e., non-test files can be excluded) or false negatives (i.e., test files can remain in the search space), most projects in the subjects listed in Table 1 follow the path and file naming rules and the above procedure can be effective.

4 ANALYSIS RESULTS

4.1 Baseline Performance

We apply each considered IRBL technique to every subject listed in Table 1, and check the output ranked list of source code files identified as potential bug locations against the collected ground truth dataset previously collected from the subjects. As described in Section 3.1, we use the default parameters specified in the original work of the six IRBL techniques.

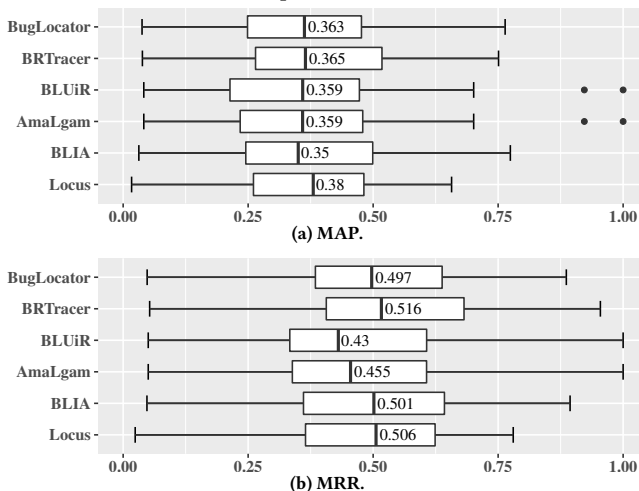


Figure 3: Distribution of MAP/MRR values of subjects for each technique: All (old/new) subjects listed in Table 1 with single version matching and test files included.

Note that the performance values shown in Figure 3 are measured with the following setting: (1) all subjects (old + new) listed in Table 1, (2) single version matching strategy, and (3) including test

code files. This setting is similar to the configurations used in the evaluation of the six techniques described in Section 3.1 (cf., Locus does not include test files by default).

We measured the performance (MAP and MRR) by applying each IRBL technique to each subject. Figure 3 shows the overall distributions of MAP and MRR for each technique. MAP values of all techniques are ranged from 0 to 1. While BLUIR and Amalgam have higher outliers, BLIA has the best top value (0.774 of the IO subject). The average values are 0.37, 0.39, 0.37, 0.38, 0.36, and 0.36, respectively (in sequence of BugLocator, BRTracer, BLUIR, Amalgam, BLIA, and Locus). MRR (Figure 3b) values are also ranged from 0 to 1. The average values are 0.50, 0.52, 0.47, 0.48, 0.49, and 0.47, respectively (with the same sequence of MAP’s average values).

In spite of small differences observed above, none of the six techniques substantially outperforms the others, all presenting a tight range of MAP and MRR values; the median values of MAP and MRR are ranged as 0.350 – 0.380 and 0.430 – 0.516, respectively.

From the perspective of practice, users of IRBL techniques can expect a similar performance by using any of the six techniques. Although each technique yields better performance for some projects (MAP: 0.6 – 0.75 and MRR: 0.75 – 1.0), it is not feasible to figure out whether an IRBL technique works better for a specific project without a prior knowledge. This encourages to characterize project characteristics and investigate the sensitivity of IRBL techniques on the characteristics.

In addition, the overall MAP and MRR values suggest that the field of IR-based bug localization still has much room for improvement from the research perspective. Despite recent efforts in bug localization, 35 – 50% precision and recall values (even for file-level) might not acceptable to practitioners. Further investigation should include (1) the characteristics of those projects for which the IRBL techniques ineffectively perform, (2) the acceptable level of performance to users, and (3) effective granularity of localization.

On the need for a further analysis

There are no significant differences, in terms of performance, among the six investigated IRBL techniques. A thorough analysis of the limitations can reveal various opportunities for research directions for improving the state-of-the-art.

4.2 RQ1: Subject Groups

To compare the performance of the techniques between old and new subject groups, we compute MAP and MRR values for each subject. Our results are summarized in two different ways: (1) Table 2 summarizes the results of each technique assuming all the bug reports are aggregated into a single (virtual) project and (2) Table 3 shows individual performance values each pair of project and technique. In the single virtual evaluation setting, we compute MAP and MRR values assuming all bug reports are in a single project. In the latter setting, the values are computed by each project. Therefore, the average MAP/MRR values in Table 3 are different from the values in Table 2. In Table 2, we use Mann-Whitney U test [26] to find out whether the differences are significant. This statistical test is also used for comparing the average of the following experiments. In Table 3, the highest MAP and MRR values for each subject across the

Table 9: Execution time (minutes) for the tools.

Technique	Min	Max	Median	Mean	Total
BugLocator	1	70	10	15	702
BRTracer	1	71	10	15	710
BLUiR	1	209	9	22	1,034
AmaLgam	1	217	9	24	1,095
BLIA	1	208	12	21	944
Locus	1	9,636	13	373	17,163

in the subject. Note that AmaLgam’s execution time is computed by adding BLUiR’s time and AmaLgam’s pure execution time since AmaLgam uses the results of BLUiR.

BugLocator and BRTracer need less time (15 min on average) to execute than the rest. BLUiR, AmaLgam, and BLIA took slightly more time (22, 24, and 21 min on average, respectively) than the first two techniques. Running Locus requires significantly more time (373 min on average). Based on our inspection, the original implementation of Locus has a performance issue due to text concatenation in the middle of its execution. We fixed the issue and confirmed that it can significantly improve the performance. However, we do not report here since the goal of this study is to reproduce the results of the original implementations.

Overall, for all six techniques, the execution time is proportional to the number of bug reports in a subject. BLUiR’s analysis time is higher than BugLocator and BRTracer due to its use of Indri [39], an open source IR toolkit, which accounts for most of the execution time. Considering revision history incurs more overhead for AmaLgam, BLIA, and Locus. The low runtime performance of Locus is mainly due to change hunk scanning as the technique leverages hunk-level change information even for file-level bug localization.

5.2 Threats to Validity

External validity: Our study examines only Java subjects as listed in Table 1. However, the same process in the study can be applied to other subjects implemented by another programming language. Another threat to the validity of our study is that our subjects are all based on an open source development model. The practice in the software industry may involve projects with specific characteristics that may be even more suitable (or in contrast unsuitable) for IRBL techniques. The provided replication package should help practitioners validate IRBL techniques in their context.

Internal validity: We use the Mann-Whitney U test [26] to examine statistical significance. This method may, however, present limitations. Nevertheless, the methods are commonly used in the literature to figure out the significance.

6 RELATED WORK

Bug localization techniques. Topic modeling and semantic analysis are common techniques used in IRBL. PROMESIR [34] utilizes Latent Semantic Analysis (LSI) [9] to identify buggy files. Lukins et al. [25] adopted Latent Dirichlet Allocation (LDA) [7] to their approach that models source code topics and showed its effectiveness with a small number of case studies. BugScout [30], on the other hand, builds topic models for both source code and bug reports and compares their distribution to locate files to fix a bug.

Stack traces are regarded as a promising information source in bug localization. Wong et al. proposed a BRTracer [46] which further considers stack traces in similarity scores. Lobster [29] also

uses stack traces to compare with code elements in source code files. CrashLocator [48] focuses more on stack traces together with function call graphs.

Other IRBL techniques consider machine learning. Ye et al. [51] proposed a learning-to-rank approach to bug localization based features representing the degree of suspiciousness. Kim et al. [18] dealt with bug report quality to improve bug localization with a two-phase model focusing on high-quality bug reports.

IRBL-related studies. Closely related to our work, Le et al. [22, 23] have proposed a study where they attempt to predict whether the ranked list produced by a bug localization tool is likely to be relevant to the given bug. They extract various textual and metadata features from 3 old projects and test on two IRBL techniques. They indeed find that it is possible, to some extent, to predict the effectiveness of the considered techniques. Our work is a generalized and large-scale investigation into the question of IRBL performance.

Wang et al. [43] have conducted an analytical study and a user study on IRBL techniques to assess their usefulness. Focusing on a single technique, BugLocator, and four common projects from previous studies, they report that the information needed for IR-based techniques to be effective is often not available in bug reports.

7 CONCLUSION

We presented a comprehensive reproducibility analysis of state-of-the-art IR-based bug localization techniques as a contribution to the community towards (1) demonstrating the actual performance of current approaches, (2) providing an updated benchmark for furthering this research field, and (3) showing the performance variations of different evaluation strategies.

Our reproducibility study have yielded several findings for the practice and research around bug localization. Overall, while IRBL approaches exhibit similar performance scores, subjects are not all equivalently adapted for each technique. All techniques are not overfitted to outdated subjects and perform better for up-to-date subjects as well. In contrast to the common belief, including test cases does not degrade the performance and, rather, even improves localization results. In addition, we reveal a potential research direction that leveraging duplicate reports together would enhance the performance of IRBL techniques.

Our future work will include (1) investigating relationships between project/report/file characteristics and the performance of different IRBL techniques (cf. D&C approach [19]), (2) building a decision model that predicts which IRBL technique performs better than others for a given project of file, and (3) improving preprocessing steps of IRBL techniques to reduce noise.

We provide a **replication package** with datasets and scripts as Bench4BL, at <https://github.com/exatoa/Bench4BL>.

ACKNOWLEDGEMENTS

This work was supported by the Fonds National de la Recherche (FNR), Luxembourg, under projects RECOMMEND C15/IS/10449467, FIXPATTERN C15/IS/9964569. This work was also supported by the National Research Foundation of Korea (NRF) grant funded by the Korea government (Ministry of Science, ICT & Future Planning) (No. 2015R1C1A1A01054994).

REFERENCES

- [1] R. Abreu, P. Zoetewij, and A.J.C. van Gemund. 2007. On the Accuracy of Spectrum-based Fault Localization. In *Testing: Academic and Industrial Conference Practice and Research Techniques - MUTATION, 2007. TAICPART-MUTATION 2007*. 89–98.
- [2] Rui Abreu, Peter Zoetewij, and Arjan J. C. van Gemund. 2006. An Evaluation of Similarity Coefficients for Software Fault Localization. In *12th Pacific Rim International Symposium on Dependable Computing, 2006. PRDC '06*. IEEE, 39–46.
- [3] N. Bettenburg, R. Premraj, T. Zimmermann, and Sunghun Kim. 2008. Duplicate bug reports considered harmful ... really?. In *IEEE International Conference on Software Maintenance, ICSM 2008*. IEEE, 337–345.
- [4] T. F. Bissyandé, D. Lo, L. Jiang, L. Réveillère, J. Klein, and Y. L. Traon. 2013. Got issues? Who cares about it? A large scale investigation of issue trackers from GitHub. In *2013 IEEE 24th International Symposium on Software Reliability Engineering (ISSRE)*. 188–197.
- [5] Tegawendé F Bissyandé, Laurent Réveillère, Julia L Lawall, and Gilles Muller. 2012. Diagnosys: automatic generation of a debugging interface to the linux kernel. In *Automated Software Engineering (ASE), 2012 Proceedings of the 27th IEEE/ACM International Conference on*. IEEE, 60–69.
- [6] Tegawendé F Bissyandé, Ferdian Thung, Shaowei Wang, David Lo, Lingxiao Jiang, and Laurent Reveillere. 2013. Empirical evaluation of bug linking. In *17th European Conference on Software Maintenance and Reengineering (CSMR)*. IEEE, 89–98.
- [7] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. 2003. Latent dirichlet allocation. *Journal of Machine Learning Research* 3 (March 2003), 993–1022.
- [8] Brendan Cleary, Chris Exton, Jim Buckley, and Michael English. 2009. An empirical analysis of information retrieval based concept location techniques in software comprehension. *Empirical Software Engineering* 14, 1 (01 Feb 2009), 93–130.
- [9] Scott Deerwester, Susan T. Dumais, George W. Furnas, Thomas K. Landauer, and Richard Harshman. 1990. Indexing by latent semantic analysis. *Journal of the American Society for Information Science* 41, 6 (Sept. 1990), 391–407.
- [10] Nicholas DiGiuseppe and James A. Jones. 2014. Fault density, fault types, and spectra-based fault localization. *Empirical Software Engineering* 20, 4 (March 2014), 928–967.
- [11] B. Dit, A. Holtzhauer, D. Poshyvanik, and H. Kagdi. 2013. A dataset from change history to support evaluation of software maintenance tasks. In *2013 10th Working Conference on Mining Software Repositories (MSR)*. 131–134.
- [12] B. Dit, E. Moritz, M. Linares-Vásquez, and D. Poshyvanik. 2013. Supporting and Accelerating Reproducible Research in Software Maintenance Using TraceLab Component Library. In *2013 IEEE International Conference on Software Maintenance*. 330–339.
- [13] Bogdan Dit, Meghan Revelle, Malcom Gethers, and Denys Poshyvanik. 2013. Feature location in source code: a taxonomy and survey. *Journal of Software: Evolution and Process* 25, 1 (Jan. 2013), 53–95.
- [14] Bogdan Dit, Meghan Revelle, and Denys Poshyvanik. 2013. Integrating information retrieval, execution and link analysis algorithms to improve feature location in software. *Empirical Software Engineering* 18, 2 (April 2013), 277–309.
- [15] William B. Frakes and Ricardo Baeza-Yates. 1992. *Information Retrieval: Data Structures and Algorithms* (1 ed.). Prentice Hall.
- [16] G. Gay, S. Haiduc, A. Marcus, and T. Menzies. 2009. On the use of relevance feedback in IR-based concept location. In *2009 IEEE International Conference on Software Maintenance*. 351–360.
- [17] N. Jalbert and W. Weimer. 2008. Automated duplicate detection for bug tracking systems. In *IEEE International Conference on Dependable Systems and Networks With FTCS and DCC, DSN 2008*. IEEE, 52–61.
- [18] Dongsun Kim, Yida Tao, Sunghun Kim, and A. Zeller. 2013. Where Should We Fix This Bug? A Two-Phase Recommendation Model. *IEEE Transactions on Software Engineering* 39, 11 (Nov. 2013), 1597–1610.
- [19] Anil Koyuncu, Tegawendé F. Bissyandé, Kui Liu, Dongsun Kim, Jacques Klein, Yves Le Traon, and Martin Monperrus. 2018. *D&C: A Divide-and-Conquer, IR-based, Multi-Classifer Approach to Bug Localization*. Technical Report TR-2018-SerVAL-01. University of Luxembourg.
- [20] Anil Koyuncu, Tegawendé F. Bissyandé, Dongsun Kim, Jacques Klein, Martin Monperrus, and Yves Le Traon. 2017. Impact of Tool Support in Patch Construction. In *Proceedings of the 26th ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA 2017)*. ACM, 237–248.
- [21] T. D. B. Le, M. Linares-Vásquez, D. Lo, and D. Poshyvanik. 2015. RCLinker: Automated Linking of Issue Reports and Commits Leveraging Rich Contextual Information. In *2015 IEEE 23rd International Conference on Program Comprehension*. 36–47.
- [22] Tien-Duy B. Le, Ferdian Thung, and David Lo. 2014. Predicting effectiveness of ir-based bug localization techniques. In *Software Reliability Engineering (ISSRE), 2014 IEEE 25th International Symposium on*. IEEE, 335–345.
- [23] Tien-Duy B. Le, Ferdian Thung, and David Lo. 2016. Will this localization tool be effective for this bug? Mitigating the impact of unreliability of information retrieval based bug localization tools. *Empirical Software Engineering* (2016), 1–43.
- [24] Xiaoyong Liu and W. Bruce Croft. 2004. Cluster-based Retrieval Using Language Models. In *Proceedings of the 27th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '04)*. ACM, New York, NY, USA, 186–193.
- [25] Stacy K. Lukins, Nicholas A. Kraft, and Letha H. Etzkorn. 2010. Bug localization using latent Dirichlet allocation. *Information and Software Technology* 52, 9 (Sept. 2010), 972–990.
- [26] H. B. Mann. 1947. On a Test of Whether one of Two Random Variables is Stochastically Larger than the Other. *The Annals of Mathematical Statistics* 18, 1 (March 1947), 50–60.
- [27] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. 2008. *Introduction to Information Retrieval* (1 edition ed.). Cambridge University Press, New York.
- [28] Christopher D. Manning and Hinrich Schütze. 1999. *Foundations of Statistical Natural Language Processing* (1 edition ed.). The MIT Press, Cambridge, Mass.
- [29] L. Moreno, J. J. Treadway, A. Marcus, and Wuwei Shen. 2014. On the Use of Stack Traces to Improve Text Retrieval-Based Bug Localization. In *2014 IEEE International Conference on Software Maintenance and Evolution (ICSME)*. 151–160.
- [30] Anh Tuan Nguyen, Tung Thanh Nguyen, J. Al-Kofahi, Hung Viet Nguyen, and T. N. Nguyen. 2011. A topic-based approach for narrowing the search space of buggy files from a bug report. In *2011 26th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 263–272.
- [31] Anh Tuan Nguyen, Tung Thanh Nguyen, Hoan Anh Nguyen, and Tien N. Nguyen. 2012. Multi-layered approach for recovering links between bug reports and fixes. In *Proceedings of the ACM SIGSOFT 20th International Symposium on the Foundations of Software Engineering*. ACM, New York, NY, USA, 63:1–63:11.
- [32] A. Panichella, B. Dit, R. Oliveto, M. D. Penta, D. Poshyvanik, and A. D. Lucia. 2016. Parameterizing and Assembling IR-Based Solutions for SE Tasks Using Genetic Algorithms. In *2016 IEEE 23rd International Conference on Software Analysis, Evolution, and Reengineering (SANER)*, Vol. 1. 314–325.
- [33] Chris Parnin and Alessandro Orso. 2011. Are automated debugging techniques actually helping programmers?. In *Proceedings of the International Symposium on Software Testing and Analysis*. 199–209.
- [34] D. Poshyvanik, Y. G. Gueheneuc, A. Marcus, G. Antoniol, and V. Rajlich. 2007. Feature Location Using Probabilistic Ranking of Methods Based on Execution Scenarios and Information Retrieval. *IEEE Transactions on Software Engineering* 33, 6 (June 2007), 420–432.
- [35] Shivani Rao and Avinash Kak. 2011. Retrieval from Software Libraries for Bug Localization: A Comparative Study of Generic and Composite Text Models. In *Proceedings of the 8th Working Conference on Mining Software Repositories*. ACM, New York, NY, USA, 43–52.
- [36] Per Runeson, Magnus Alexandersson, and Oskar Nyholm. 2007. Detection of Duplicate Defect Reports Using Natural Language Processing. In *Proceedings of the 29th International Conference on Software Engineering, ICSE 2007 (ICSE '07)*. IEEE Computer Society, Washington, DC, USA, 499–510.
- [37] R. K. Saha, M. Lease, S. Khurshid, and D. E. Perry. 2013. Improving bug localization using structured information retrieval. In *2013 28th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. 345–355.
- [38] Gerard Salton and Michael J. McGill. 1986. *Introduction to Modern Information Retrieval*. McGraw-Hill, Inc., New York, NY, USA.
- [39] T Strohman, D Metzler, H Turtle, and WB Croft. 2004. Indri: A language model-based search engine for complex queries. In *Proceedings of the International Conference on Intelligence Analysis*. 2–6.
- [40] Chengnian Sun, David Lo, Xiaoyin Wang, Jing Jiang, and Siau-Cheng Khoo. 2010. A discriminative model approach for accurate duplicate bug report retrieval. In *Proceedings of the 32nd ACM/IEEE International Conference on Software Engineering - Volume 1, ICSE 2010*. ACM, Cape Town, South Africa, 45–54. ACM ID: 1806811.
- [41] A. Sureka and P. Jalote. 2010. Detecting Duplicate Bug Report Using Character N-Gram-Based Features. In *17th Asia Pacific Software Engineering Conference, APSEC 2010*. IEEE, 366–374.
- [42] G. Tasse. 2002. *The Economic Impacts of Inadequate Infrastructure for Software Testing: Final Report*. Diane Publishing Company. <https://books.google.lu/books?id=juSgPAAACAAJ>
- [43] Qianqian Wang, Chris Parnin, and Alessandro Orso. 2015. Evaluating the usefulness of IR-based fault localization techniques. In *Proceedings of the 2015 International Symposium on Software Testing and Analysis*. ACM, 1–11.
- [44] Shaowei Wang and David Lo. 2014. Version History, Similar Report, and Structure: Putting Them Together for Improved Bug Localization. In *Proceedings of the 22Nd International Conference on Program Comprehension*. ACM, New York, NY, USA, 53–63.
- [45] Ming Wen, Rongxin Wu, and Shing-Chi Cheung. 2016. Locus: locating bugs from software changes. In *Proceedings of the 31st IEEE/ACM International Conference on Automated Software Engineering*. ACM, Singapore, Singapore, 262–273.
- [46] C. P. Wong, Y. Xiong, H. Zhang, D. Hao, L. Zhang, and H. Mei. 2014. Boosting Bug-Report-Oriented Fault Localization with Segmentation and Stack-Trace Analysis. In *2014 IEEE International Conference on Software Maintenance and Evolution*.

- 181–190.
- [47] W. E. Wong, R. Gao, Y. Li, R. Abreu, and F. Wotawa. 2016. A Survey on Software Fault Localization. *IEEE Transactions on Software Engineering* 42, 8 (Aug. 2016), 707–740.
- [48] Rongxin Wu, Hongyu Zhang, Shing-Chi Cheung, and Sunghun Kim. 2014. CrashLocator: Locating Crashing Faults Based on Crash Stacks. In *Proceedings of the 2014 International Symposium on Software Testing and Analysis*. ACM, New York, NY, USA, 204–214.
- [49] Rongxin Wu, Hongyu Zhang, Sunghun Kim, and Shing-Chi Cheung. 2011. 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*. ACM, New York, NY, USA, 15–25.
- [50] Min Xie and Bo Yang. 2003. A study of the effect of imperfect debugging on software development cost. *IEEE Transactions on Software Engineering* 29, 5 (2003), 471–473.
- [51] X. Ye, R. Bunescu, and C. Liu. 2016. Mapping Bug Reports to Relevant Files: A Ranking Model, a Fine-Grained Benchmark, and Feature Evaluation. *IEEE Transactions on Software Engineering* 42, 4 (April 2016), 379–402.
- [52] K. C. Youm, J. Ahn, J. Kim, and E. Lee. 2015. Bug Localization Based on Code Change Histories and Bug Reports. In *2015 Asia-Pacific Software Engineering Conference (APSEC)*. 190–197.
- [53] Jian Zhou, Hongyu Zhang, and David Lo. 2012. Where should the bugs be fixed? - more accurate information retrieval-based bug localization based on bug reports. In *Proceedings of the 2012 International Conference on Software Engineering*. IEEE Press, Piscataway, NJ, USA, 14–24.