

# Improved Bug Localization using Association Mapping and Information Retrieval

Anonymous Author(s)

## ABSTRACT

Bug localization is one of the most challenging tasks undertaken by developers during software maintenance. Most existing studies rely on the lexical similarity between a bug report and the source code for bug localization. Unfortunately, this approach suffers due to vocabulary mismatch issues and often results into unrecognised similarity or spurious matching. In this paper, we propose a bug localization technique namely BLuAMIR that (1) not only uses lexical similarity between a given bug report (query) and the source code documents but also (2) exploits the implicit associations between them based on the bug-fixing history. Experiments using a collection of 3,431 bug reports show that on average our technique can localize buggy files with a Top-10 accuracy of 74.06%, a mean reciprocal rank of 0.52 and a mean average precision of 41% which are highly promising. Comparison with three state-of-the-art approaches and their variants shows that our technique can achieve 32.26% higher MAP and 26.74% higher Top-5 accuracy.

## CCS CONCEPTS

• **Software and its engineering** → **Software verification and validation**; **Software testing and debugging**;

## KEYWORDS

Bug localization, bug report, source code, information retrieval, keyword-source code association

### ACM Reference Format:

Anonymous Author(s). 2019. Improved Bug Localization using Association Mapping and Information Retrieval. In *Proceedings of Proceedings of the 27th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE '19)*. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 INTRODUCTION

Bug localization is a process of locating source code that needs to be changed in order to fix a given bug. Manually locating buggy files is not only time-consuming but also prohibitively costly in terms of development efforts [34]. This is even more challenging for large software systems. Thus, effective automated approaches are highly warranted for localizing software bugs. Traditional Information Retrieval (IR) based bug localization techniques [27, 38] accept a bug report (query) and a subject system as inputs and return a list of buggy entities (e.g., classes, methods). They localize the bugs by simply relying on the *lexical similarity* between a bug report and

the source code. Hence, they are likely to be affected by the quality and content of a submitted query (i.e., bug report). That is, if the query does not contain adequate information, then the retrieved results might not be relevant at all. As existing findings [22, 33] suggest, bug reports could be of low quality and could miss the appropriate keywords. Thus, lexical similarity alone is likely not sufficient enough for solving the bug localization problem.

In order to address the limitations of lexical similarity based approaches, several existing studies [14, 15, 17] derive the underlying semantics of a text document by employing Latent Semantic Analysis (LSA). Marcus et al. and colleagues adopt this technology in the context of concept location [15, 17], program comprehension [14] and traceability link recovery problems [16], and reported higher performance than that of the Vector Space Model (VSM) and probabilistic models. Unfortunately, their approach suffers from a major limitation. Latent Semantic Indexing (LSI) requires the use of a dimensionality reduction parameter that must be tuned for each document collection [10]. The results returned by LSI can also be difficult to interpret, as they are expressed using a *numeric spatial representation*. Other related studies [13, 20] adopt Latent Dirichlet Allocation (LDA) for bug localization. However, they are also subject to their *hyper-parameter tuning* and could even be outperformed by much simpler models (e.g., rVSM [38]).

In this paper, we propose a bug localization approach namely BLuAMIR that (1) not only considers *lexical similarity* between a bug report (the query) and the source code (2) but also captures *implicit association* between them from the bug fixing history. First, we determine the lexical similarity between each source document and the query using Vector Space Model (VSM). Second, we construct association maps between keywords of previously fixed bug reports and their corresponding changed documents using a bipartite graph [2]. Third, we prioritize such source documents that are historically associated with the keywords (from the query at hand) according to these maps. Then, we rank the source code documents based on their *lexical* and *association scores*. Thus, our approach uses implicit association to address the vocabulary mismatch between a bug report (the query) and the source code. That is, unlike traditional IR-based bug localization approaches [27, 38], BLuAMIR can return the buggy documents even if the query does not lexically match with the source code documents. Our approach also does not require the *dimensionality reduction* since we use a finite graph rather than a large sparse term-document matrix [15, 16].

We evaluate our technique from three different aspects using three widely used performance metrics and 3,431 bug reports (i.e., queries) from four open source subject systems. First, we evaluate in terms of bug localization performance metrics, and contrast with two replicated baselines – Latent Semantic Indexing (LSI) [16] and the basic Vector Space Model (VSM) [30]. BLuAMIR localizes bugs with 9%–37% higher accuracy (i.e., Hit@10), 12%–63% higher precision (i.e., MAP), and 11%–64% higher reciprocal ranks (i.e., MRR) than these baselines (Section 4.3). Second, we compare our

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

ESEC/FSE '19, August 26–30, 2019, Tallinn, Estonia

© 2019 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

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

Table 1: A working example of BLuAMIR

VSM (Lexical Similarity Only)			BLuAMIR (Lexical Similarity + Implicit Association)				
Retrieved Documents	S <sub>LS</sub>	GT	Retrieved Documents	S <sub>LS</sub>	S <sub>Assoc</sub>	S <sub>Total</sub>	GT
ClasspathLocation.java	1.00	✗	CompletionEngine.java	0.67	1.00	0.80	✓
JavaCore.java	0.74	✗	AbstractDecoratedTextEditor.java	0.69	0.73	0.71	✗
SearchableEnvironmentRequestor.java	0.73	✗	AntEditor.java	0.69	0.73	0.71	✗
Compiler.java	0.71	✗	JavaCore.java	0.74	0.64	0.70	✗
AccessRuleSet.java	0.70	✗	Engine.java	0.70	0.70	0.70	✗

GT = Ground Truth,  $S_{Total} = (1-\alpha) \times S_{LS} + \alpha \times S_{Assoc}$  and weighting parameter,  $\alpha = 0.4$

technique with three state-of-the-art approaches – BugScout [20], BugLocator [38] and BLUiR [27] (Section 4.3). Our technique can localize bugs with 6%–54% higher accuracy (i.e., Hit@5), 4%–32% higher precision (i.e., MAP) and 8%–27% higher reciprocal ranks (i.e., MRR) than those of the state-of-the-art approaches. Third, in terms of query-wise improvement, BLuAMIR improves the result ranks of 46.37% baseline queries from Eclipse subject system and outperforms the baseline VSM approach (Section 4.3). Thus, this paper makes the following contributions:

- (1) A novel technique that not only considers the *lexical similarity* between a bug report and the source code but also exploits their *implicit associations* through bug-fixing history for bug localization.
- (2) Comprehensive evaluation of the technique using *three* widely used performance metrics and a total of 3,431 bug reports from *four* open source subject systems.
- (3) Comparison with not only *two* baselines [16, 30] but also *three* state-of-the-art approaches – BugScout [20], BugLocator [38] and BLUiR [27] with statistical tests.
- (4) Experimental meta data and our used dataset for replication and third party reuse.

The rest of the paper is organized as follows– Section 2 discusses a motivating example of our proposed approach, and Section 3 presents proposed bug localization method for BLuAMIR, and Section 4 focuses on the conducted experiments and experimental results, and Section 5 identifies the possible threats to validity, and Section 6 discusses the existing studies related to our research, and finally, Section 7 concludes the paper with future plan.

## 2 MOTIVATING EXAMPLE

Let us consider a bug report (ID# 95167) from an Eclipse subsystem namely eclipse.jdt.core. Table 2 shows the *title* and *description* of the bug report. We capture both fields and construct a baseline query by employing standard natural language preprocessing (e.g., stop word removal, token splitting). Then we execute the query with the Vector Space Model (VSM) and our approach–BLuAMIR, and attempt to locate the buggy source documents.

According to the ground truth constructed from bug-fixing history, one source code document (CompletionEngine.java) was changed to fix the reported bug. As shown in Table 1, we see that the *lexical similarity* based approach, VSM, fails to retrieve any buggy source document within the Top-5 positions. On the contrary, our approach, BLuAMIR, complements the *lexical similarity* based ranking with *implicit association* scores, and then returns the target buggy document at the top most position of the result list. This is not only promising but also the best possible outcome that an automated approach can deliver.

Table 2: An Example Bug Report (#95167, eclipse.jdt.core)

Field	Content
Title	[content assist] Spurious "Access restriction" error during code assist
Description	(1) OSGi, Runtime, SWT, JFace, UI, Text loaded from head, (2) open type on AbstractTextEditor, (3) at start of createPartContro method, type: PartSite <Ctrl+Space>, and (4) it has no effect in the editor, but the status line flashes in red: Access restriction: The type SerializableCompatibility is not accessible due to restriction on required project org.eclipse.swt. The type name doesn't seem to matter. "abcd" has the same effect. I notice that org.eclipse.ui.workbench.texteditor's classpath has an access rule forbidding <code>*/internal/**</code> refs.

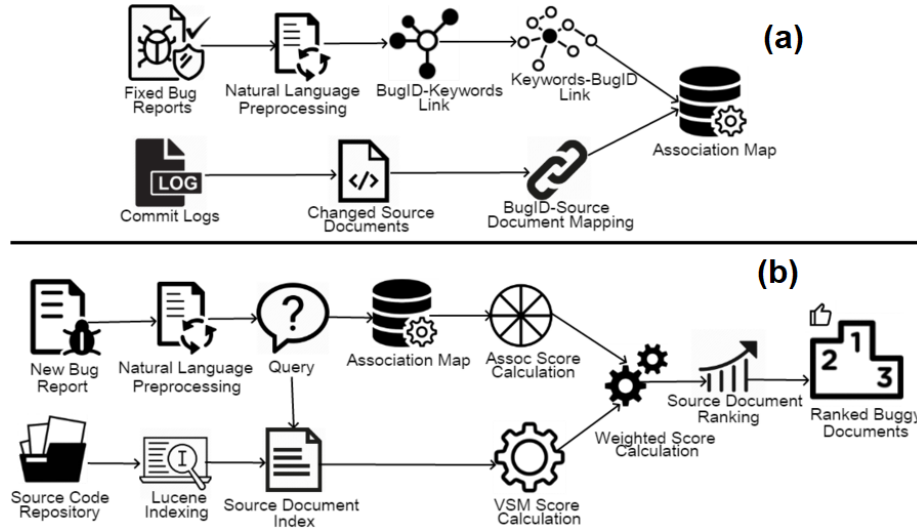
We also investigate why BLuAMIR performs better than VSM in localizing the buggy document(s). Table 1 shows different scores from both approaches. We see that several source documents (e.g., ClasspathLocation.java) that are retrieved by VSM are strongly similar to the query. However, such similarity does not necessarily make them buggy. In fact, VSM returns the ground truth document at the lowest position of the Top-10 results (not shown in Table 1). Thus, *lexical similarity* alone might not be sufficient enough for effective bug localization. However, our approach overcomes such challenge by exploiting the *implicit association* between the query and the buggy documents (Section 3.2), and returns the target ground truth at the top-most position. Although the lexical similarity is low (e.g., 0.67), our approach correctly identifies the buggy document using its strong implicit association score (e.g., 1.00) with the given query.

## 3 BLuAMIR: PROPOSED APPROACH FOR BUG LOCALIZATION

Figure 1 shows the schematic diagram of our proposed approach. Furthermore, Algorithms 1, 2 show the pseudo-code of our approach. First, we (a) construct an association map between the bug reports and their corresponding changed source documents with the help of bug-fixing history. Then we (b) retrieve buggy source code documents for a given query (bug report) by leveraging not only their lexical similarity but also their implicit associations derived from the association map above. We discuss different parts of our proposed approach in the following sections.

### 3.1 Construction of Keyword–Document Association Map

We construct an association map between keywords from previously fixed bug reports and their corresponding changed source documents. The map construction involves three steps. We not only



**Figure 1: Schematic diagram of BLuAMIR—(a) Construction of association map between keywords and source documents, and (b) Bug Localization using VSM and implicit association**

show different steps of our map construction (Fig. 1-(a)) but also provide the corresponding pseudo-code (Algorithm 1). We discuss each of these steps as follows:

**(1) Extraction of Keywords from Bug Reports:** We collect *title* and *description* of each bug report from a subject system. Then we perform standard natural language preprocessing, and remove punctuation marks, stop words and small words (i.e.,  $length \leq 2$ ) from them. Stop words convey very little semantics in a sentence. We use appropriate regular expressions to discard all the punctuation marks and a standard list<sup>1</sup> for stop word removal. Finally, we select a list of remaining keywords from each bug report. We then create a *map* (e.g.,  $M_{bk}$ ) between the ID of each bug report and their corresponding keywords (Lines 5–6, Algorithm 1). We also construct an *inverted index* (e.g.,  $M_{kb}$ ) between keywords and their corresponding bug report IDs (Lines 7–8, Algorithm 1) by using the above constructed map.

**(2) Extraction of Changed Source Documents from Bug-Fix Commits:** We analyse the version control history of each subject system, and identify the bug-fixing commits using appropriate regular expressions [3, 35]. In particular, we go through all the commits and identify such commits that contain keywords related to bug fix or resolution in their title messages. Then, we collect the *changeset* (i.e., list of changed documents) from each of these bug-fix commits, and construct a Bug ID–document map (e.g.,  $M_{bs}$ ) for our study (Lines 9–10, Algorithm 1). We use several utility commands such as `git`, `clone` and `log` on Git-based repository of each system for collecting the above information.

**(3) Construction of Mapping between Keywords and Changed Source Documents:** The above two steps deliver (1) an inverted index ( $M_{kb}$ ) that maps each individual keyword to numerous bug report IDs, and (2) a map ( $M_{bs}$ ) that links each bug report ID to its corresponding changed source documents based on the bug-fixing history. Since both of these maps are connected through bug report IDs, keywords from each bug report also enjoy an implicit,

<sup>1</sup><https://www.ranks.nl/stopwords>

#### Algorithm 1 Construction of Association Map

```

1: procedure MAPCONSTRUCTION( $BRC, BFC$ )
2:    $\triangleright BRC$ : a collection of past bug reports
3:    $\triangleright BFC$ : bug-fix commit history
4:    $M_{KS} \leftarrow \{\}$   $\triangleright$  an empty association map
5:    $\triangleright$  Create map between Bug ID and keywords
6:    $M_{bk} \leftarrow \text{createBugIDtoKeywordMap}(BRC)$ 
7:    $\triangleright$  Create an inverted index between keywords and ID
8:    $M_{kb} \leftarrow \text{createKeywordtoBugIDMap}(M_{bk})$ 
9:    $\triangleright$  Map between Bug ID and changed source documents
10:   $M_{bs} \leftarrow \text{createBugIDtoSourceMap}(BFC)$ 
11:   $\triangleright$  Mapping keywords to the changed source documents
12:   $KW \leftarrow \text{collectKeywords}(M_{kb})$ 
13:  for Keyword  $kw \in KW$  do
14:     $ID_{kw} \leftarrow \text{extractBugIDs}(kw, M_{kb})$ 
15:    for BugID  $id_{kw} \in ID_{kw}$  do
16:       $\triangleright$  Get the linked source documents
17:       $SD_{link} \leftarrow \text{getDocuments}(M_{bs}, id_{kw})$ 
18:       $\triangleright$  Map all source documents to this keyword
19:       $M[kw].link \leftarrow \{M[kw].link \cup SD_{link}\}$ 
20:    end for
21:  end for
22:   $\triangleright$  collect all keyword–source document mappings
23:   $M_{KS} \leftarrow M[KW]$ 
24:  return  $M_{KS}$ 
25: end procedure

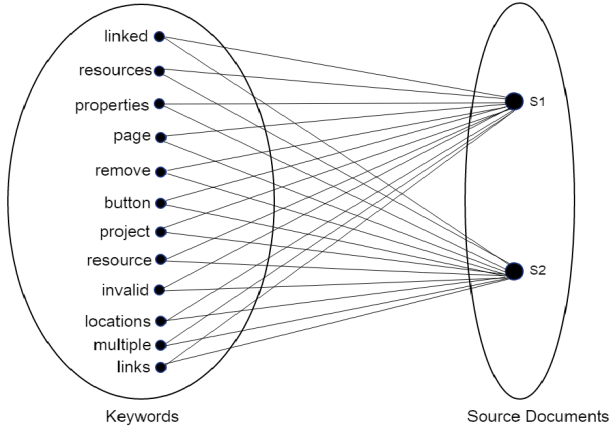
```

transitive relationship with the corresponding changed source code documents. We leverage such transitive relationship and construct a bipartite graph (e.g., Fig. 2) by explicitly connecting the keywords from each bug report to their corresponding changed source documents (Lines 11–22, Algorithm 1). Here, one or more keywords could be linked to single buggy source code document. Conversely, single source document could also be linked to one or more keywords from multiple bug reports.



**Table 3: Example Bug Report (#322401, eclipse.ui.platform)**

Field	Content
Title	[LinkedResources] Linked Resources properties page should have a Remove button
Description	Project properties > Resource > Linked Resources: Especially for invalid locations, a Remove button would be handy to remove one or multiple links.

**Figure 2: An example association map using bipartite graph**

**An Example Association Map with Bipartite Graph:** In Mathematics, *bipartite graph* is defined as a special graph that (1) has two disjoint sets of nodes and a set of edges and (2) each of its edges connects two nodes from the two different sets but not from the same set [2]. In our study, these two sets correspond to keywords (from bug reports) and source code documents. Thus, no connection between any two keywords or any two source documents is allowed. We construct an example bipartite graph for a bug report (ID 322401) from Eclipse UI Platform. Table 3 shows the *title* and *description* of the bug report. According to the bug-fix history, these two documents -  $S_1$ : IDEWorkbenchMessages.java and  $S_2$ : LinkedResourceEditor.java - constitute the ground truth of the bug report above. Fig. 2 shows the bipartite graph for the bug report along with its ground truth. We see that twelve keywords and two source documents represent the two sets of nodes and each node is connected to all the nodes from another set. In our approach, we iteratively update such a bipartite graph with new nodes and new connections generated from each of the bug reports of a subject system (Lines 15–20, Algorithm 1).

### 3.2 Bug Localization using VSM and Implicit Association

Fig. 1-(b) shows the schematic diagram and Algorithm 2 presents the pseudo-code of our bug localization component. We have constructed an association map (e.g.,  $M_{KS}$ ) that connects the keywords from a bug report to its corresponding changed source documents (Section 3.1). We leverage this association map, and return a list of buggy source documents that are not only lexically similar but also strongly associated with a given query (i.e., bug report at hand).

**Algorithm 2 Proposed Bug Localization Approach**

```

1: procedure BUGLOCALIZATION( $Q, SCrepo, M_{KS}$ )
2:    $\triangleright q$ : a given query (bug report)
3:    $\triangleright SCrepo$ : a source code repository
4:    $\triangleright M_{KS}$ : keyword-source document map
5:    $RL \leftarrow \{\}$   $\triangleright$  list of buggy source documents
6:    $\triangleright$  Create an index of source documents
7:    $Index \leftarrow \text{createLuceneIndex}(SCrepo)$ 
8:    $\triangleright$  Collect lexically similar documents
9:    $VSM \leftarrow \text{getLexSimDocuments}(q, Index)$ 
10:   $\triangleright$  Collecting associated documents
11:   $Assoc \leftarrow \text{getAssociatedDocuments}(q, M_{KS})$ 
12:   $\triangleright$  Combine both candidate document lists
13:   $C \leftarrow \{VSM.doc \cup Assoc.doc\}$ 
14:  for SourceDocument  $d \in C$  do
15:     $\triangleright$  Combine lexical and association scores
16:     $C[d].score \leftarrow VSM[d].score$ 
17:     $C[d].score \leftarrow C[d].score + Assoc[d].score$ 
18:  end for
19:   $\triangleright$  Sort the candidates and collect Top-K documents
20:   $RL \leftarrow \text{getTopKDocuments}(\text{sortByScore}(C))$ 
21:  return  $RL$ 
22: end procedure

```

We thus calculate two different scores for each candidate source document and then suggest the Top-K buggy documents as follows:

**Lexical Similarity Score:** Source code documents often share a major *overlap* in the vocabulary with a submitted bug report. Many of the existing studies [26, 27, 35, 38] consider such vocabulary overlap (i.e., lexical similarity) as a mean to localize the buggy source documents. These studies generally employ Vector Space Model (VSM) [30] for calculating the vocabulary overlap. VSM is a classical approach for constructing vector representation of any text document (e.g., bug report, source code document) [30]. First, it encodes a document collection (a.k.a., corpus) using a term-by-document matrix where each row represents a term and each column represents a document. Second, each matrix cell is defined as the frequency of a term within a specific document (i.e., term frequency). Thus, *lexical similarity* between a given query (bug report) and a candidate source document is computed as the *cosine* or *inner product* between their corresponding vectors from the matrix. Since term frequency (TF)-based vector representation might be biased toward large documents, several studies [27, 38] also represent their vectors using TF-IDF. It stands for term frequency times inverse document frequency. TF-IDF assigns higher weights to such terms that are frequent within a document but not frequent across the whole document collection [28].

In classic VSM, term frequency  $tf(t, d)$  and inverse document frequency  $idf(t)$  are defined as follows:

$$tf(t, d) = \frac{f_{td}}{\#terms}, \quad idf(t) = 1 + \log\left(\frac{\#docs}{n_t}\right)$$

Here  $f_{td}$  refers to the frequency of each unique term  $t$  in a document  $d$ ,  $n_t$  denotes the document frequency of term  $t$ ,  $\#docs$  is the total number of documents in the corpus and  $\#terms$  refers to the total number of unique terms in the corpus. Thus, the lexical similarity between a given query  $q$  (i.e., given bug report) and a candidate

source code document  $d$  is calculated as follows:

$$\text{lexicalSimScore}(q, d) = \text{cosine}(q, d) = \frac{\sum_{t \in \{q \cap d\}} tf(t, q) \times tf(t, d) \times idf(t)^2}{\sqrt{\sum_{t \in q} (tf(t, q) \times idf(t))^2} \times \sqrt{\sum_{t \in d} (tf(t, d) \times idf(t))^2}}$$

Here  $\text{lexicalSimScore}$  takes a value between 0 and 1 where 0 means total lexical dissimilarity and 1 means strong lexical similarity between the query  $q$  and the candidate source code document  $d$ . We use *Apache Lucene* for first creating the corpus index and then for calculating the above lexical similarity (Lines 6–9, Algorithm 2).

**Implicit Association Score:** We analyse implicit associations between a given query (bug report) and each candidate source document using our constructed association map (from Section 3.1). In this map, while each keyword could be linked to multiple candidate documents, each document could also be linked to multiple keywords across multiple bug reports. We perform standard natural language preprocessing on a given query, and extract a list of query keywords. We then identify such source documents in the map that are linked to each of these keywords. Since these links were established based on the bug-fixing history, they represent an *implicit relevance* between the keywords and the candidate source documents. It should be noted that such relevance does not *warrant for any lexical similarity*. We analyse such links for all the keywords ( $\forall t \in q$ ) of a search query  $q$ , and determine how frequently each candidate source document  $d$  was associated with these keywords in the past as follows:

$$\text{associationScore}(q, d) = \sum_{t \in q} \sum_{id \in ID_t} \#link(t, d, id)$$

Here  $\#link(t, d, id)$  returns the frequency of association between the keyword  $t \in q$  and the source document  $d$  for the bug report with ID  $id \in ID_t$ . That is,  $\text{associationScore}$  assigns a score to each candidate document by capturing their historical co-occurrences with the query keywords across the bug-fixing history of a system.

**Final Score Calculation:** The above two sections deliver two different scores (*lexical similarity*, *implicit association*) for each of the candidate source code documents. Since these scores could be of different ranges, we normalize both of them between 0 and 1. We then combine both scores using a weighted summation, and calculate the final relevance score for each candidate document  $d$  (Lines 14–18, Algorithm 2) as follows:

$$\text{FinalScore}(q, d) = (1 - \alpha) \times \text{Norm}(\text{lexicalSimScore}) + \alpha \times \text{Norm}(\text{associationScore})$$

Here, the weighting factor  $\alpha$  varies from 0.2 to 0.4, and the detailed justification is provided in the discussion section (Section 4.3).

Once the final score is calculated for each of the candidate source documents from the corpus, we return the only Top-K ( $1 \leq K \leq 10$ ) results as the buggy source documents for the given query  $q$  (i.e., Lines 19–21, Algorithm 2).

## 4 EXPERIMENT

We evaluate our technique in three different aspects using three widely used performance metrics and 3,431 bug reports (i.e., queries) from four open source subject systems. We compare not only with

**Table 4: Experimental Dataset**

Project	Version	Study Period	#Bugs	#Documents
Eclipse	v3.1	Oct 2004 - Mar 2011	3071	11,831
SWT	v3.1	Oct 2004 - Apr 2010	98	484
AspectJ	-	July 2002 - Oct 2006	244	3519
ZXing	-	Mar 2010 - Sep 2010	20	391

two baseline techniques [16, 30] but also three state-of-the-art studies [20, 27, 38] from the literature. In particular, we answer three research questions with our experiment as follows:

- **RQ<sub>1</sub>:** How does the proposed approach perform in bug localization compared to the baseline approaches?
- **RQ<sub>2</sub>:** Can the proposed approach outperform the state-of-the-art studies in bug localization?
- **RQ<sub>3</sub>:** (a) Can implicit association make any significant difference in IR-based bug localization? and (b) can BLuAMIR overcome the challenges with large documents?

### 4.1 Experimental Dataset

Table 4 shows our experimental dataset. We use a total of 3,431 bug reports from four open source systems—Eclipse, AspectJ, SWT and ZXing—for our experiment. These systems are collected from two existing, frequently used public benchmarks [27, 38]. *Eclipse* is a well-known large-scale system which is frequently used in empirical Software Engineering research. *SWT* is a component of Eclipse IDE. *AspectJ* is a part of iBUGs dataset provided by University of Saarland. *ZXing* is an android based system maintained by Google. We collect *title* and *description* from each of these 3,431 bug reports. We perform standard natural language preprocessing on them, and remove stop words, punctuation marks, and small words from them. Then each of these preprocessed bug reports is used as the *baseline query* for our experiment.

**Ground Truth Selection:** We analyse the bug-fixing commits from each subject system, and select the ground truth for our experiment. In particular, we go through all commit messages of a system and identify such commits that deals with bug fixing or resolution (i.e., bug fixing commits) using appropriate regular expressions [3]. We then extract the changed source documents from each of these commits, and map them to the fixed bug ID. Such changed documents are then used as the *ground truth* for the corresponding bug reports (i.e., queries). The same approach has been widely used by the literature [26, 35, 38] for the ground truth selection.

### 4.2 Performance Metrics

We choose three performance metrics for evaluation and validation. These metrics have been widely used by the relevant literature [27, 38]. Thus, they are highly appropriate for our experiments as well. They are defined as follows:

**Recall at Top K / Hit@K:** It represents the percentage of bug reports for each of which at least one ground truth buggy document is successfully retrieved within the Top-K results. The higher the Hit@K value is, the better the bug localization performance is.

**Mean Reciprocal Rank (MRR):** The reciprocal rank of a query is the multiplicative inverse of the rank of the first buggy document within the result list. Hence, mean reciprocal rank is calculated as

**Table 5: Comparison between Baseline VSM and BLuAMIR**

#Bugs	Approach	Hit@1	Hit@5	Hit@10	MRR	MAP
3071	VSM	23.09%	47.54%	57.57%	0.33	0.24
	BLuAMIR	27.45%	<b>53.79%</b>	<b>63.30%</b>	<b>0.38</b>	0.27

the mean of reciprocal ranks of a set of queries  $Q$  as follows:

$$MRR(Q) = \frac{1}{|Q|} \sum_{q \in Q} \frac{1}{firstRank(q)}$$

where  $firstRank(q)$  is the rank of the first buggy source document within the returned result list for a query  $q \in Q$ . Reciprocal rank takes a value between 0 and 1. Here, 0 value suggests that no buggy document is retrieved whereas 1 value suggests that the first buggy document is retrieved at the top most position of the result list.

**Mean Average Precision (MAP):** Mean Average Precision is a commonly used metric for evaluating the ranking approaches. Unlike MRR, it considers the ranks of all buggy documents retrieved in the result list into the consideration. Average Precision of a query  $q$  can be computed as follows:

$$AP(q) = \sum_{k=1}^K \frac{P(k) \times buggy(k)}{|R|}, \quad P(k) = \frac{\#buggyDocs}{k}$$

where  $k$  is a rank in the ranked results,  $K$  is the number of retrieved documents, and  $R$  is the set of true positive instances.  $buggy(k)$  function indicates whether the  $k_{th}$  document is buggy or not.  $P(k)$  returns the calculated precision at a given rank position. Since AP is a metric for single query  $q$ , MAP can be calculated for a set of queries  $Q$  as follows:

$$MAP(Q) = \frac{1}{|Q|} \sum_{q \in Q} AP(q)$$

MAP takes a value between 0 and 1. The higher the metric value is, the better the bug localization performance is.

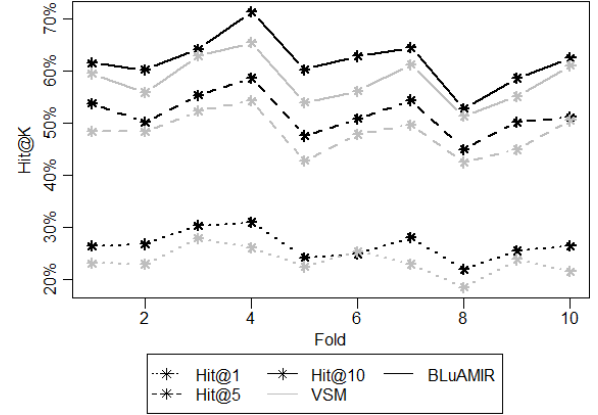
### 4.3 Evaluation & Validation

We use a total of 3,431 bug reports from four subject systems (Section 4.1) and evaluate our approach using three performance metrics (Section 4.2). We employ 10-fold cross validation, and the whole bug report collection is divided into 10 folds. Out of these 10 folds, nine folds are used for training (i.e., construct of association map) and the remaining fold is used for testing (i.e., bug localization). We repeat this process 10 times, determine our performance each time, and then report the average performance of our approach. In the following sections, we discuss our experimental results and answer our research questions (RQ<sub>1</sub>, RQ<sub>2</sub>) as follows:

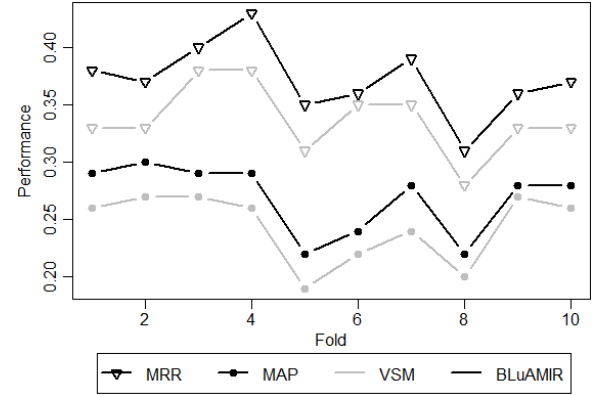
#### Answering RQ<sub>1</sub>—Comparison with Baseline Approaches:

We compare the performance of BLuAMIR with two baseline approaches – Vector Space Model (VSM) [30] and Latent Semantic Indexing (LSI) [16]. We replicate both approaches in our development environment, evaluate with our dataset, and then determine their performance. We discuss the comparison with these baselines in details as follows:

**Baseline VSM vs BLuAMIR:** While VSM approach simply relies on the lexical similarity between a given query (bug report) and the source code document, our approach additionally considers their implicit associations based on the bug-fixing history. We thus



**Figure 3: Comparison between baseline VSM and proposed approach, BLuAMIR, in Hit@K**



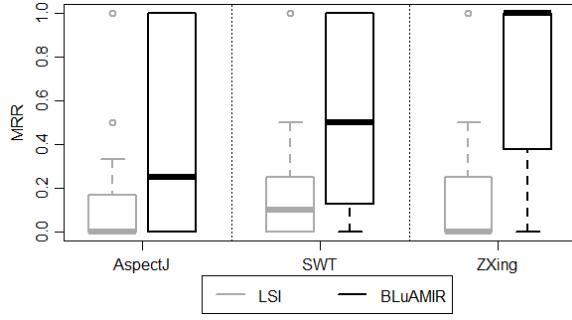
**Figure 4: Comparison between baseline VSM and proposed approach, BLuAMIR, in MRR and MAP**

compare BLuAMIR with VSM, and investigate whether the addition of association score could improve the overall bug localization performance or not. Table 5 and Figures 3, 4 show the comparison between VSM and our approach for *Eclipse* system. From Table 5, we see that the baseline VSM achieves only 23% Hit@1, 48% Hit@5 and 58% Hit@10. On the contrary, our approach achieves 28% Hit@1, 54% Hit@5 and 63% Hit@10 which are 19%, 13% and 10% higher respectively. BLuAMIR also achieves 15% higher MRR and 13% higher MAP than the baseline measures. Since our experiment involves 10-fold cross validation, we also compare our performance with that of baseline VSM for each of these folds. As shown in Fig. 3, we see that our Hit@K is better than the baseline for each fold of dataset. Fig. 4 shows how BLuAMIR outperforms the baseline VSM in MRR and MAP respectively for each of these folds. We also perform Wilcoxon Signed-Rank tests, and found that BLuAMIR achieves statistically significant improvement over the baseline in terms of both MRR (i.e.,  $p\text{-value}=0.005<0.05$ ) and MAP (i.e.,  $p\text{-value}=0.005<0.05$ ).

**Baseline LSI vs BLuAMIR:** Traditional VSM often suffers from vocabulary mismatch problem. Latent Semantic Indexing (LSI) attempts to overcome such problem by extracting the underlying meaning of a document rather than simply relying on the individual words from the document. As existing study [4] suggests,

**Table 6: Comparison between Baseline LSI and BLuAMIR**

#System	Approach	Hit@1	Hit@5	Hit@10	MRR	MAP
AspectJ	LSI	9.42%	24.59%	29.10%	0.15	0.07
	BLuAMIR	33.20%	54.92%	66.39%	0.43	0.23
SWT	LSI	13.26%	30.61%	53.06%	0.22	0.17
	BLuAMIR	<b>45.93%</b>	<b>75.00%</b>	<b>82.29%</b>	<b>0.58</b>	<b>0.50</b>
ZXing	LSI	15.00%	35.00%	45.00%	0.23	0.21
	BLuAMIR	55.00%	80.00%	85.00%	0.67	0.62

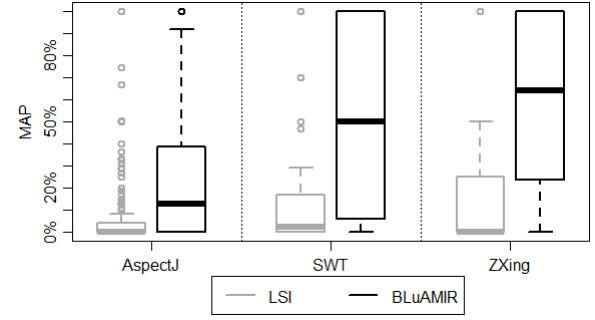
**Figure 5: Comparison between baseline LSI and proposed approach, BLuAMIR, in MRR**

individual words often do not provide reliable evidence about the conceptual topic or meaning of a document. LSI has been used for document retrieval in several Software Engineering contexts [15, 16] which makes it an attractive baseline for our evaluation. We replicate LSI in our development environment following these steps. First, we create a term-by-document matrix by capturing both bug reports (queries) and source code documents from a subject system. Second, we apply Singular Value Decomposition (SVD) on this matrix, and construct a subspace (i.e., LSI subspace) [29]. Third, we compute cosine similarity between queries (bug reports) and candidate source documents using corresponding vectors from this subspace. Fourth, Top-K source code documents are then returned for each given query (bug report) based on their similarity. We compare BLuAMIR with baseline LSI using three subject systems—AspectJ, SWT and ZXing.

Table 6 contrasts our approach against LSI in terms of Hit@K, MRR and MAP. We see that our approach outperforms LSI in all the cases. For example, baseline LSI performs the best with SWT system, and achieves 53% Hit@10 with a MRR of 0.22 and a MAP of 0.17. On the contrary, our approach achieves 82% Hit@10, a MRR of 0.58 and a MAP of 0.50 which are 55%, 163% and 194% higher respectively. Box plots on MRR (Fig. 5) and MAP (Fig. 6) also demonstrate that our approach outperforms the baseline LSI with large margins for each of the subject systems.

**Summary of RQ<sub>1</sub>:** Our approach outperforms two baseline approaches [16, 30] with statistically significant, large margins. BLuAMIR can even deliver **10%–55%** higher Hit@10, **15%–163%** higher MRR and **13%–194%** higher MAP than the baseline.

**Answering RQ<sub>2</sub>– Comparison with the State-of-the-Art:** While the above evaluation clearly shows that our approach is promising, we still wanted to see where our method stands compared to the state-of-the-art methods. Thus, we compare BLuAMIR

**Figure 6: Comparison between baseline LSI and proposed approach, BLuAMIR, in MAP****Table 7: Comparison with State-of-the-art**

System	Approach	Hit@1	Hit@5	Hit@10	MRR	MAP
Eclipse	BugScout	14.00	24.00	31.00	–	–
	BugLocator	24.36	46.15	55.90	0.35	0.26
	BLUiR	30.96	53.20	62.86	0.42	0.32
	BLuAMIR	27.45	<b>53.79</b>	<b>63.30</b>	0.38	0.27
AspectJ	BugScout	11.00	26.00	35.00	–	–
	BugLocator	22.73	40.91	55.59	0.33	0.17
	BLUiR	32.17	51.05	60.49	0.41	0.24
	BLuAMIR	<b>33.20</b>	<b>54.92</b>	<b>66.39</b>	<b>0.43</b>	0.23
SWT	BugLocator	31.63	65.31	77.55	0.47	0.40
	BLUiR	55.10	76.53	87.76	0.65	0.56
	BLuAMIR	45.93	75.00	82.29	0.58	0.50
ZXing	BugLocator	40.00	55.00	70.00	0.48	0.41
	BLUiR	40.00	65.00	70.00	0.49	0.38
	BLuAMIR	<b>55.00</b>	<b>80.00</b>	<b>85.00</b>	<b>0.67</b>	<b>0.62</b>

with three state-of-the-art techniques – BugScout [20], BugLocator [38] and BLUiR [27] – using our benchmark subject systems. BugScout employs Latent Dirichlet Allocation (LDA) for localizing software bugs. BugLocator combines a revised Vector Space Model (rVSM) and past bug reports for improving the IR-based bug localization. Finally, BLUiR exploits the structures from both bug reports (queries) and source code documents, and then improves the bug localization with structured information retrieval. Table 7 summarizes our comparison details. Nguyen et al. [20] used two subject systems (Eclipse and AspectJ) in order to evaluate BugScout. Both Saha et al. [27] and Zhou et al. [38] employ all four subject systems for their evaluation. Since we use the same dataset from the earlier studies [27, 38], we compare with their published results.

From Table 7, we see that BLuAMIR outperforms both BugScout and BugLocator consistently across all the systems (Eclipse, AspectJ, SWT and ZXing) and all performance metrics. BugLocator is a common baseline for a number of existing studies [27, 34, 35]. BugLocator performs the best with SWT system, and achieves 78% Hit@10 with a MRR of 0.47 and a MAP of 0.40. On the contrary, BLuAMIR achieves 82% Hit@10 with a MRR of 0.58 and a MAP of 0.50 which are 6%, 23% and 25% higher respectively. We also see that our approach performs better than BLUiR with three out of four systems—Eclipse, AspectJ and ZXing—in terms of several performance metrics (i.e., emboldened measures). While BLUiR performs better than ours with SWT system, it only contains 98 bug reports. On the contrary, AspectJ contains 2.5 times more bug reports, and



**Table 8: Comparison of Result Rank Improvement with State-of-the-art on Eclipse System**

Approach	#Bugs	Improvement				Worsening				# Preserved
		#Improved	Mean	Min.	Max.	#Worsened	Mean	Min.	Max.	
BugLocator	2911	1315 (45.17%)	68.46	1	949	1084 (37.23%)	172.87	1	8030	512 (17.59%)
BLuAMIR	2864	1328 (46.37%)	50.90	1	909	763 (26.64%)	52.10	1	811	773 (26.99%)

**Table 9: Impact of Weighting Parameter on BLuAMIR (with SWT System)**

$\alpha$	Hit@1	Hit@5	Hit@10	MRR	MAP
0.1	43.75	73.96	83.33	0.57	0.50
0.2	48.96	75.00	83.33	0.60	0.51
0.3	47.92	72.92	85.42	0.59	0.51
0.4	45.93	75.00	82.29	0.58	0.50

our approach outperforms BLuIR with this system. For example, while BLuIR achieves 60% Hit@10, our approach improves upon this metric by 10% which is promising. We also investigate why our approach fails to outperform BLuIR with SWT system. In particular, we calibrate the weighting parameter  $\alpha$  and demonstrate that BLuAMIR could perform comparably with BLuIR with the right  $\alpha$  value. Table 9 summarizes our investigation results. We see that our approach delivers improved Hi@1, Hit@5, MAP and MRR at  $\alpha = 0.2$ . However, it achieves the highest Hit@10 (85%) at  $\alpha = 0.3$ , which is comparable to that (88%) of BLuIR [27]. Thus, choosing the right  $\alpha$  needs significant trade-off which we leave as an area for future study. However, such investigation clearly states that our approach has a high potential for bug localization.

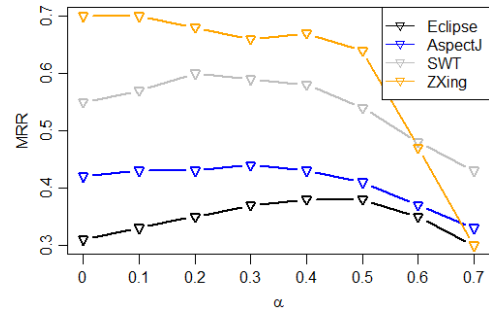
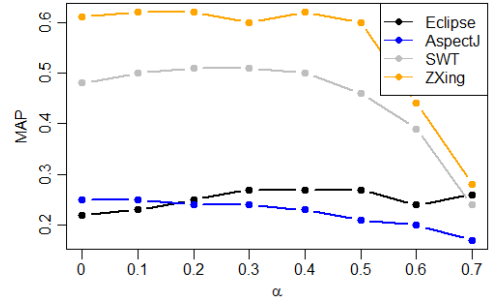
We also investigate how each of the three approaches – BugLocator, BLuIR and BLuAMIR – improve upon baseline VSM in terms of result rank improvement for the queries of Eclipse system. In particular, we (1) collect the rank of first correct buggy document within the results returned by baseline VSM (i.e., baseline rank), and (2) collect similar rank for each of these approaches (i.e., changed rank), and (3) then determine the result improvement and worsening by comparing the changed ranks with the baseline ranks. Table 8 contrasts our approach with the state-of-the-art in terms of result rank improvement and worsening. We see that BugLocator improves upon VSM for 45.17% of the queries and worsens 37.23% of the queries. On the contrary, BLuAMIR improves 46.37% and worsens 26.64% of the queries which are 2.66% higher and 39.75% lower respectively. All these findings clearly suggest the high potential of our approach in bug localization.

**Summary of RQ<sub>2</sub>:** Our approach outperforms two of the state-of-the-art studies [20, 38] and performs comparably with one study [27]. BLuAMIR not only achieves 23% higher MRR but also delivers 25% higher MAP than the state-of-the-art.

**Answering RQ3-(a) Impact of Association Score on Bug Localization:** While RQ<sub>1</sub> and RQ<sub>2</sub> suggest that our proposed approach has high potential for bug localization, at this point, we wanted to investigate in-depth why our proposed approach performs better than the lexical similarity-based approaches. We make use of not only the lexical similarity but also the implicit association between a bug report (query) and the source code that is derived from the bug fixing history. Thus, in this RQ, we wanted to see whether such implicit association has any significant impact on

**Table 10: Impact of Weighting Parameter on BLuAMIR (with Eclipse System)**

$\alpha$	Hit@1	Hit@5	Hit@10	MRR	MAP
0.0	21.20	43.05	54.58	0.31	0.22
0.1	22.83	46.47	58.02	0.33	0.23
0.2	24.23	49.53	60.44	0.35	0.25
0.3	25.72	51.87	62.45	0.37	0.27
0.4	27.45	53.79	63.30	0.38	0.27

**Figure 7: The impact of  $\alpha$  on BLuAMIR's MRR****Figure 8: The impact of  $\alpha$  on BLuAMIR's MAP**

the bug localization performance or not. In particular, we tune our approach with different  $\alpha$  values, and record the responses. Since  $\alpha$  is the *relative weight of association score*, such tuning actually evaluates the impact of association score. Table 10 shows how the gradual increment in  $\alpha$  value improves the overall performance. We see that each of the metrics improves consistently up to  $\alpha = 0.4$  with Eclipse system. We also repeat the same investigation with different  $\alpha$  values with all four systems. From Fig. 7, we see that BLuAMIR initially performs low with  $\alpha = 0$ , i.e., no association score. Then our MRR improves continuously with the increment of  $\alpha$ . We see a somewhat stationary state during  $0.2 \leq \alpha \leq 0.4$  followed by a sharp decrease in MRR with all systems except Eclipse. This performance reaches the lowest at  $\alpha = 0.7$ . We also note a pretty similar scenario in Fig. 8 when investigated with MAP. Such findings above suggest that association score has a noticeable impact upon the bug location performance. That is, lexical similarity +



association score is consistently better than lexical similarity alone. Our findings in RQ<sub>1</sub> also support this observation with empirical evidence. Such findings also justify our choice for  $\alpha$  between 0.2 and 0.4 in our approach.

**(b) Overcoming the Challenge with Large Documents:** We further wanted to investigate whether our approach could work well with large source documents. This is in particular important to investigate since traditional VSM approaches are mostly biased towards such large source documents during bug localization. Although these documents have a major overlap with a given bug report (query), they are often noisy and less relevant in practice. We discuss how BLuAMIR can overcome this challenge during bug localization with an example as follows:

**A Case Study with Large Documents:** Let us consider a bug report (ID #95561) from `eclipse.ui.platform` subsystem of Eclipse. The bug report discusses about a synchronization issue with Eclipse workbench. We preprocess the *title* and *description* of the bug report, and use them as the query for bug localization. According to the bug-fixing history, `org.eclipse.ui.internal.WorkbenchPage.java` constitutes the ground truth for this query, which is a large document with 4K+ lines of code and  $\approx 14$ K words. Baseline VSM [30] and BugLocator [38] return this ground truth at the 30<sup>th</sup> and 26<sup>th</sup> positions respectively. On the contrary, our approach returns the same ground truth at the second position of the result list which is promising. We then perform in-depth analysis on the ranking mechanism of each approach, and make interesting observations. First, the lexical score for this document is 0.45 due to noise, which is much lower than that of the other smaller candidate documents. Such a low score possibly pushes the document down the ranked list. On the contrary, BLuAMIR complements such low lexical scores with a strong association score of 0.95. The document and the query are strongly associated according to the bug-fixing history. Thus, final score for the document is calculated by BLuAMIR as follows:  $(0.95 \times 0.40 + 0.45 \times 0.60) = 0.65$ . Such a final score places the candidate source document at the second position in our case. Thus, *implicit association* dimension can actually help overcome the challenges (noises) with large source documents.

**Summary of RQ<sub>3</sub>:** Our empirical and analytical findings suggest that **implicit association** can make a significant difference to IR-based bug localization by (a) **improving** upon baseline VSM and (b) **overcoming** the challenges with large candidate documents.

## 5 THREATS TO VALIDITY

This section discusses the validity and generalizability of our findings. In particular, we discuss the threats to internal validity, construct validity, and external validity as follows:

**Internal Validity:** We used three artifacts of a software repository: bug reports, source code documents and bug fixing history, which are generally well understood. Our evaluation uses four open source subject systems - Eclipse, AspectJ, SWT, and ZXing. These systems are collected from two existing, frequently used public benchmarks [27, 38]. Bug reports provide crucial information for developers to fix the bugs. A bug report could be of low quality and could miss the appropriate keywords. Hence, such a “bad” bug report could cause a delay in the bug fixing process. Our proposed

approach makes use of both lexical similarity as well as implicit association between a bug report and their corresponding source code documents. Therefore, if a bug report does not provide adequate information, or provides inappropriate keywords, the performance of BLuAMIR could be adversely affected.

**External Validity:** Threats to external validity is related to the generalizability of our findings. To reduce this threat, we have analyzed 3,431 bug reports from four popular subject systems. Still in the future, we plan to reduce these threats further by analyzing more bug reports from more subject systems written in multiple programming languages.

**Construct Validity:** In our experiment, we use three evaluation metrics – Hit@K, MAP and MRR, and one statistical test, i.e., Wilcoxon signed-rank test. These metrics have been widely used in the literature to evaluate bug localization performances [27, 38]. Thus, threats to construct validity are likely to be mitigated.

**Parameter Tuning:** We tune our weighting parameter using controlled iterative experiments, and report our findings (Section 4.3). However, such iterative experiments might not guarantee an optimal parameter value. In future, we plan to use more sophisticated parameter tuning methods such as Genetic Algorithms.

## 6 RELATED WORK

Existing automatic bug localization or automatic debugging approaches can be broadly categorized into two types - dynamic and static techniques. Generally, dynamic approaches can localize a bug much more precisely than static approaches. These techniques usually contrast the program spectra information (such as execution statistics) between passed and failed executions to compute the fault suspiciousness of individual program elements (i.e., statements, branches, and predicates), and rank these program elements by their fault suspiciousness. Developers may then locate faults by examining a list of program elements sorted by their suspiciousness. Some of the well known dynamic approaches are spectrum-based fault localization [1, 9, 11, 25], model-based fault localization [6, 18], dynamic slicing [37], and delta debugging [36]. These approaches require a test case suite and need to execute the program for collecting the passed and failed execution traces. Moreover, these approaches are computationally expensive. Our proposed approach neither requires a test case nor requires the execution of the program.

Static approaches, on the other hand, do not require any program test cases or execution traces. In most cases, they need only program source code documents and bug reports. They are also computationally efficient. The static approaches usually can be categorized into two groups: program analysis based approaches and IR-based approaches. FindBugs is a program analysis based approach that locates a bug based on some predefined bug patterns [8]. Therefore, FindBug does not even need a bug report. However, it often detects too many false positives and misses many real bugs [32]. On the contrary, our proposed approach is not program analysis based, and thus, does not depend on any predefined bug patterns.

Traditional Information Retrieval based techniques localize the bugs by simply relying on the *lexical similarity* between a bug report (query) and the source code. There are three traditionally-dominant IR paradigms: TFIDF [28], the probabilistic approach known as BM25 [24], and more recently, language modeling [21]. However,

using an empirical study, Fang et al. [5] show that all three approaches perform comparably when well-tuned. On the other hand, Rao and Kak [23] investigated many standard information retrieval techniques for bug localization and found that simpler techniques, e.g., TFIDF and SUM, perform the best. However, If a bug report does not contain adequate information (i.e., appropriate keywords), lexical similarity alone might not be sufficient enough to localize the bugs successfully. In order to address the issue with lexical similarity, Latent Semantic Indexing (LSI) introduces latent concepts by exploring deeper methods matching “concepts”. While a probabilistic variant of LSI has been devised [7], its probability model was found to be deficient. On the other hand, several other studies [14, 15, 17] derive the underlining semantic of a text document by employing Latent Semantic Analysis (LSA). Unfortunately, their approach suffers from a major limitation. Latent Semantic Indexing (LSI) requires the use of a dimensionality reduction parameter that must be tuned for each document collection [10]. The results returned by LSI can also be difficult to interpret, as they are expressed using a numeric spatial representation. Our approach does not require dimensionality reduction since we use a finite graph rather than a large sparse term-document matrix [15, 16]. Lukins et al. [13] and [20] adopt Latent Dirichlet Allocation (LDA) for bug localization. LDA is not able to predict the appropriate topic because it follows a generative topic model in a probabilistic way [12]. Besides, they are also subject to their hyper-parameters and could even be outperformed by simpler models (e.g., rVSM [38]). LDA is also computationally expensive. In this paper, we compare with a baseline LSI [16] and a LDA based technique namely BugScout [20] and the detailed comparison can be found in Section 4.3 (i.e., RQ1 and RQ2 respectively). Our approach outperforms both these baseline [16, 20] with large margins across all subject systems.

Recently automatic bug localization techniques have gone beyond traditional Information Retrieval based techniques by collecting additional information from software repositories [19, 27, 31, 34, 38]. Sisman and Kak [31] propose a history-aware IR-based bug localization solution to achieve a better result. Zhou et al. [38] propose BugLocator, which leverages similarities among bug reports and uses a refined vector space model to perform bug localization. Saha et al. [27] build BLUIR that considers the structure of bug reports and source code files and employs structured retrieval to achieve a better result. Moreno et al. [19] use a text retrieval based technique and stack trace analysis to perform bug localization. To locate buggy files, they combine the textual similarity between a bug report and a program entity (e.g., class) and the structural similarity between the stack trace and the entity. Different from the existing IR-based bug localization approaches, Wang and Lo [34] proposed Amalgam, a new method for locating relevant buggy files that puts together version history, similar reports, and structure, to achieve better performance. Later, they proposed Amalgam+ [35], which is a method for locating relevant buggy files that puts together five sources of information i.e., version history, similar reports, structure, stack traces, and reporter information. However, in our proposed technique BLuAMIR, we consider only three sources of information i.e., bug reports, source code documents and bug fixing history. Hence, we performed extensive comparison with two closely related studies – BugLocator [38] and BLUIR [27]. The details of our comparison can be found in Section 4.3 (i.e., RQ2). Unlike the

above traditional IR-based approaches [19, 27, 31, 34, 35, 38], our proposed approach could deliver the buggy documents even if the bug report (the query) does not lexically match with source code documents. Our proposed approach outperforms BugLocator [38] consistently across all the subject systems. BLuAMIR also performs better than BLUIR [27] with three out of four subject systems and perform comparably with the remaining system.

## 7 CONCLUSION AND FUTURE WORK

During the evolution of a software system, a large number of bug reports are submitted. For a large software project, developers typically need to examine a large number of source code files in order to locate the buggy files responsible for a bug, which is tedious and expensive work. In this paper, we propose BLuAMIR, a new bug localization technique not only based on lexical similarity but also on an implicit association mapping between bug report keywords and source code documents. We perform large-scale experiments with four subject systems, namely Eclipse, SWT, AspectJ and ZXing to localize 3,431 bug reports. We also compare our technique with two baseline techniques i.e., [30], [16], and three state-of-the-art IR-based bug localization techniques i.e., BugScout [20], BugLocator [38] and BLUIR [27]. Our experiment with those dataset shows that on average our technique can locate buggy files with a Top-10 accuracy of 74.06% and a mean reciprocal rank of 0.52 and a mean precision average of 41%, which is highly promising. Our technique can localize bugs with 6%–54% higher accuracy (i.e., Hit@5), 4%–32% higher precision (i.e., MAP) and 8%–27% higher reciprocal ranks (i.e., MRR) than these state-of-the-art approaches. All these empirical findings suggest the high potential of our proposed approach for bug localization.

In the future, we will explore whether other bug related information such as bug report structure, source code structure, stack traces, reporter information, and similar bug information could be integrated into our approach or not. Such external resources could also assist in handling the vocabulary mismatch issues and the challenges with poor queries.

## REFERENCES

- [1] R. Abreu, P. Zoeteveij, R. Golsteijn, and A. J. C. van Gemund. 2009. A Practical Evaluation of Spectrum-based Fault Localization. *J. Syst. Softw.* 82, 11 (Nov. 2009), 1780–1792.
- [2] A. S. Asratian, T. M. J. Denley, and R. Häggkvist. 1998. *Bipartite Graphs and Their Applications*. Cambridge University Press, New York, NY, USA.
- [3] A. Bachmann and A. Bernstein. 2009. Software Process Data Quality and Characteristics: A Historical View on Open and Closed Source Projects. In *Proc. IW/PSE*. 119–128.
- [4] S. Deerwester, S. T. Dumais, G. W. Furnas, T. K. Landauer, and R. Harshman. 1990. Indexing by Latent Semantic Analysis. *Journal of the American Society for Information Science* 41, 6 (1990), 391–407.
- [5] H. Fang, T. Tao, and C. Zhai. 2004. A Formal Study of Information Retrieval Heuristics. In *Proc. SIGIR*. 49–56.
- [6] A. Feldman and A. van Gemund. 2006. A Two-step Hierarchical Algorithm for Model-based Diagnosis. In *Proc. AAAI*. 827–833.
- [7] Thomas Hofmann. 2017. Probabilistic Latent Semantic Indexing. *SIGIR Forum* 51, 2 (2017), 211–218.
- [8] D. Hovemeyer and W. Pugh. 2004. Finding Bugs is Easy. *SIGPLAN Not.* 39, 12 (2004), 92–106.
- [9] J. A. Jones and M. J. Harrold. 2005. Empirical Evaluation of the Tarantula Automatic Fault-localization Technique. In *Proc. ASE*. 273–282.
- [10] A. Kontostathis. 2007. Essential Dimensions of Latent Semantic Indexing (LSI). In *Proc. HICSS*. 73–73.
- [11] Lucia, D. Lo, L. Jiang, and A. Budi. 2010. Comprehensive Evaluation of Association measures for Fault Localization. In *Proc. ICSM*. 1–10.
- [12] S.K. Lukins, N.A. Kraft, and L.H. Etzkorn. 2008. Source Code Retrieval for Bug Localization Using Latent Dirichlet Allocation. In *Proc. WCRE*. 155–164.
- [13] S. K. Lukins, N. A. Kraft, and L. H. Etzkorn. 2010. Bug Localization Using Latent Dirichlet Allocation. *Inf. Softw. Technol.* 52, 9 (2010), 972–990.
- [14] J. I. Maletic and A. Marcus. 2001. Supporting Program Comprehension Using Semantic and Structural Information. In *Proc. ICSE*. 103–112.

- [15] A. Marcus and J. I. Maletic. 2001. Identification of High-level Concept Clones in Source Code. In *Proc. ASE*. 107–114.
- [16] A. Marcus and J. I. Maletic. 2003. Recovering Documentation-to-Source-Code Traceability Links Using Latent Semantic Indexing. In *Proc. ICSE*. 125–135.
- [17] A. Marcus, A. Sergeyev, V. Rajlich, and J. I. Maletic. 2004. An Information Retrieval Approach to Concept Location in Source Code. In *Proc. WCRE*. 214–223.
- [18] W. Mayer and M. Stumptner. 2007. Model-Based Debugging – State of the Art And Future Challenges. *Electron. Notes Theor. Comput. Sci.* 174, 4 (2007), 61–82.
- [19] 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 *Proc. ICSME*. 151–160.
- [20] A. T. Nguyen, T. T. Nguyen, J. Al-Kofahi, H. V. Nguyen, and T. N. Nguyen. 2011. A Topic-based Approach for Narrowing the Search Space of Buggy Files from a Bug Report. In *Proc. ASE*. 263–272.
- [21] J. M. Ponte and W. B. Croft. 98. A Language Modeling Approach to Information Retrieval. In *Proc. SIGIR*. 275–281.
- [22] M. M. Rahman and C. K. Roy. 2018. Improving IR-Based Bug Localization with Context-Aware Query Reformulation. In *Proc. ESEC/FSE*. 12.
- [23] S. Rao and A. Kak. 2011. Retrieval from Software Libraries for Bug Localization: A Comparative Study of Generic and Composite Text Models. In *Proc. MSR*. 43–52.
- [24] S. E. Robertson, S. Walker, and M. Beaulieu. 2000. Experimentation As a Way of Life: Okapi at TREC. *Inf. Process. Manage.* 36, 1 (2000), 95–108.
- [25] D. Saha, M. G. Nanda, P. Dhoolia, V. K. Nandivada, V. Sinha, and S. Chandra. 2011. Fault Localization for Data-centric Programs. In *Proc. ESCE/FSE*. 157–167.
- [26] R.K. Saha, J. Lawall, S. Khurshid, and D.E. Perry. 2014. On the Effectiveness of Information Retrieval Based Bug Localization for C Programs. In *Proc. ICSME*. 161–170.
- [27] R.K. Saha, M. Lease, S. Khurshid, and D.E. Perry. 2013. Improving Bug Localization using Structured Information Retrieval. In *Proc. ASE*. 345–355.
- [28] G. Salton and C. Buckley. 1988. Term-weighting Approaches in Automatic Text Retrieval. *Inf. Process. Manage.* 24, 5 (1988), 513–523.
- [29] G. Salton and M. J. McGill. 1986. *Introduction to Modern Information Retrieval*. McGraw-Hill, Inc.
- [30] G. Salton, A. Wong, and C. S. Yang. 1975. A Vector Space Model for Automatic Indexing. *Commun. ACM* 18, 11 (1975), 613–620.
- [31] B. Sisman and A.C. Kak. 2013. Assisting Code Search with zAutomatic Query Reformulation for Bug Localization. In *Proc. MSR*. 309–318.
- [32] H. Tang, S. Tan, and X. Cheng. 2009. A Survey on Sentiment Detection of Reviews. *Expert Systems with Applications* 36, 7 (2009).
- [33] Q. Wang, C. Parnin, and A. Orso. 2015. Evaluating the Usefulness of IR-based Fault Localization Techniques. In *Proc. ISSTA*. 1–11.
- [34] S. Wang and D. Lo. 2014. Version History, Similar Report, and Structure: Putting Them Together for Improved Bug Localization. In *Proc. ICPC*. 53–63.
- [35] Shaowei Wang and David Lo. 2016. Amalgam+: Composing Rich Information Sources for Accurate Bug Localization. *Journal of Software: Evolution and Process* 28 (2016), 921–942.
- [36] A. Zeller and R. Hildebrandt. 2002. Simplifying and Isolating Failure-Inducing Input. *IEEE Trans. Softw. Eng.* 28, 2 (2002), 183–200.
- [37] X. Zhang, H. He, N. Gupta, and R. Gupta. 2005. Experimental Evaluation of Using Dynamic Slices for Fault Location. In *Proc. AADeBUG*. 33–42.
- [38] J. Zhou, H. Zhang, and D. Lo. 2012. Where Should the Bugs be Fixed? More Accurate Information Retrieval-Based Bug Localization Based on Bug Reports. In *Proc. ICSE*. 14–24.