

Bug Report Classification using LSTM and CNN for More Accurate Software Defect Locating

Anonymous Author(s)

ABSTRACT

The PI proposes a two-year project to develop an interactive and practical system that locates software defects automatically. This system will alleviate developers' effort in bug finding and improve productivity. In preliminary work, the PI developed a ranking model that ranks all the source code files for a given bug report. A source file at a higher position in the ranked list is more likely to contain defects than a source file at a lower position. The proposed system is an extension of the PI's preliminary work. In this project, the PI propose to 1) introduce a pre-filtering technique to filter out low-quality bug reports; 2) use artificial neural network to measure semantic relationship between bug reports and source code files; 3) implement an interactive 3-D VR user interface on multiple platforms for the system. The first two objectives aim at making the system be more accurate. The third objective aims at making the system be more usable. The PI propose to recruit two student assistants including one graduate and one undergraduate to help develop the system and perform experimental evaluation. The system will be published and will be used in the software engineering (SE) courses at California Sate University San Marcos to help students learn SE concepts in a lively manner.

This project will result in an interactive and practical system that can be used by the public. To increase system accuracy, this project proposes using artificial neural networks (ANNs) to classify bug reports into either "predictable" or "unpredictable", where "unpredictable" refers to low-quality reports that do not contain sufficient information for locating the bug. The system works on "predictable" reports and keeps silent on "unpredictable" reports. Besides, this project proposes using ANNs to measure the semantic similarity between a bug report and a source file. This semantic similarity is used as a new feature that complements other lexical-similarity features. To make the system be usable, this project proposes to develop an interactive user interface on multiple platforms. Furthermore, this project proposes to develop a supplemental and selectable 3-D VR visualization module that helps developers understand the code change history more efficiently.

This project will publish a practical defect positioning system that can be used in both academic and industry. The research outcome will provide the basis for the PI's future work. The methodology and techniques used in this project will contribute to the software engineering (SE) community in similar work. The funding to this project will help the PI build an SE lab and found an SE research group at California State University San Marcos (CSUSM).

Permission to make digital or hard copies of part or all 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 third-party components of this work must be honored. For all other uses, contact the owner/author(s).

ESEC/FSE 2018, 49 November, 2018, Lake Buena Vista, Florida, United States

© 2018 Copyright held by the owner/author(s).

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM.

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

The research lab will provide students with the necessary computing resources for doing SE research projects. The research group will expose students to new SE techniques and research projects. The funding to this project will also help provide paid positions and new opportunities to our students that have strong desire in doing research projects.

CCS CONCEPTS

• **Computing methodologies** → **Supervised learning by classification**; • **Software and its engineering** → **Software testing and debugging**; • **Information systems** → *Multilingual and cross-lingual retrieval*;

KEYWORDS

Recurrent neural network, long short-term memory, convolutional neural network, bug localization, bug report

ACM Reference Format:

Anonymous Author(s). 2018. Bug Report Classification using LSTM and CNN for More Accurate Software Defect Locating. In *Proceedings of The 26th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE 2018)*. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

1 INTRODUCTION

Software development is complex. During the software development and maintenance process, to ensure product quality, one of the most common tasks of a development team is bug fixing.

A software *bug* or *defect* is a coding mistake that may cause an unintended or unexpected behavior of the software component [5]. Upon discovering an abnormal behavior of the software project, a developer or a user will report it in a document, called a *bug report* or *issue report*. A bug report provides information that could help in fixing a bug, with the overall aim of improving the software quality. A large number of bug reports could be opened during the development life-cycle of a software product. For example, 3,352 bug reports of the Eclipse Platform product were submitted in 2016 alone¹. In a software team, bug reports are extensively used by both managers and developers in decision making during their daily development activities [6].

When a new bug report is received, before being solved, it will be screened and prioritized. This screening process is called bug report triage, which aims at ensuring the bug report quality. During the triage steps, the development team will check whether the bug report contains sufficient information for developers to fix, whether it is a duplicate report, what is its priority, and who should be the appropriate developer to solve it. After the triage, the report will be assigned to a developer to fix. Then the fix will be tested usually by another developer. If the fix is verified, the report will

¹<https://bugs.eclipse.org/bugs/>

be closed. Otherwise, it will be assigned to someone again. This process represents a common defect cycle or bug-fixing cycle.

A developer who is assigned with a bug report usually needs to reproduce the abnormal behavior [19] and perform code reviews [1] in order to find the cause. However, the diversity and uneven quality of bug reports can make this process nontrivial. Essential information is often missing from a bug report [4]. To locate the bug, developers not only need to analyze the bug report using their domain knowledge, but also collect information from peer developers and users to narrow their search. Employing such a manual process in order to find and understand the cause of a bug can be tedious and time-consuming [29].

Bacchelli and Bird [1] surveyed 165 managers and 873 programmers, and reported that finding defects require a high-level understanding of the code and familiarity with the relevant source code files. In the survey, 798 respondents answered that it takes time to review unfamiliar files. While the number of source files in a project is usually large, the number of files that contain the bug is usually very small. For example, the Eclipse JDT bug 424772 was fixed in a patch (commit) with four changed files, while the source code package checked out from this commit contains 10,544 Java files. Therefore, we believe that an automatic approach that ranked the source files with respect to their relevance for the bug report could speed up the bug finding process by narrowing the search to a smaller number of possibly unfamiliar files.

If the bug report is construed as a query and the source code files in the software repository are viewed as a collection of documents, then the problem of finding source files that are relevant to a given bug report can be modeled as a standard task in information retrieval (IR) [23]. In recent years, researchers [17, 22, 30, 33, 34, 40, 47] have proposed various IR models to recommend source code files for bug reports automatically. These models use one or more *features*. Each *feature* is a type of information that measure the relationship between the bug report and the source code file. However, the weight parameters of these models are not learned automatically. Therefore, the feature combinations of these model are sub-optimal. The model scalability is also limited.

To address this issue, in preliminary work, the PI introduced a learning-to-rank model to combine different *features* automatically [43, 44]. The problem of locating software defects is approached as a ranking problem, in which the source code files (documents) are ranked with respect to their *relevance* to a given bug report (query). In this context, *relevance* is equated with the likelihood that a particular source code file contains the cause of the defect described in the bug report. The ranking function is defined as a weighted combination of *features*, where the *features* are automatically trained on previously solved bug reports using a learning-to-rank technique. We performed extensive evaluations on six large open-source Java projects, from which the results showed that the learning-to-rank model outperformed other state-of-the-art approached [42, 44].

Matching bug reports with the relevant source code files is complicated by the lexical gap between natural languages used in the bug reports and the technical terms used in the source code files. To bridge this language mismatch, the PI employed word embeddings, which are vector representations of words, to measure the semantic similarity between bug reports and source code files [45]. These

semantic similarities are used as additional *features* in the ranking function. Results of evaluations on four large open-source Java projects showed that adding these new word-embedding-based semantic features significantly improved the ranking performance [42, 45].

1.0.1 A.1 Objectives. The results of the PI's preliminary work are promising and motivate the PI to continue this work to address the following aspects.

- (1) **Bug report quality classification:** Bug report quality varies. Some bug reports do not contain sufficient information for developers to locate the defect [4]. Running a ranking system on low-quality reports may obtain misleading results. Therefore, the PI propose to include in the very beginning a pre-filtering step, in which a bug report is classified as either "predictable" or "unpredictable". After filtering out the "unpredictable" reports, the ranking system works only on the "predictable" reports and keeps silent otherwise. The purpose of doing so is to make the ranking results be more trustworthy. This project will use deep neural networks such as Convolutional Neural Network (CNN) [20] and Recurrent Neural Network (RNN) [11] for bug report classification.
- (2) **Document-level similarity:** When calculating the word-embedding-based semantic features, our preliminary work used the Mihalcea's [26] similarity function, in which the similarity between a word w in one document and a bag of words in another document T is computed as the maximum similarity between w and any word w' in T : $sim(w, T) = \max_{w' \in T} sim(w, w')$. However, the similarities between w and other words in T are ignored, which cause information loss. As such, the PI propose to explore alternative methods for aggregating word-level similarities into a document-level similarity. More specifically, this project will use deep neural networks e.g. CNN and RNN to automatically learn feature vectors from bug reports and source code files respectively. Then a new semantic feature can be computed as the cosine similarity between the feature vectors of the bug report and a source file. This new feature can be used alone or as an complementary feature to improve the system performance.
- (3) **User interface:** Despite many theories and algorithms were proposed recently to locate software defects automatically, practical systems with friendly user interfaces are rare, not to mention user studies that evaluate the effectiveness and usefulness of these approaches for developers in real work. Given the promising results of the preliminary work, which provides a self-learning ranking model on the back-end side, the PI propose to develop an interactive and user-friendly user interface on the front-end side. This project will develop a practical system that can be used in the software engineering course at California State University San Marcos (CSUSM) to help students better understand software engineering concepts.

1.1 B. Preliminary Work

The preliminary work of this proposal produced the following results, which provide solid support for the proposed study.

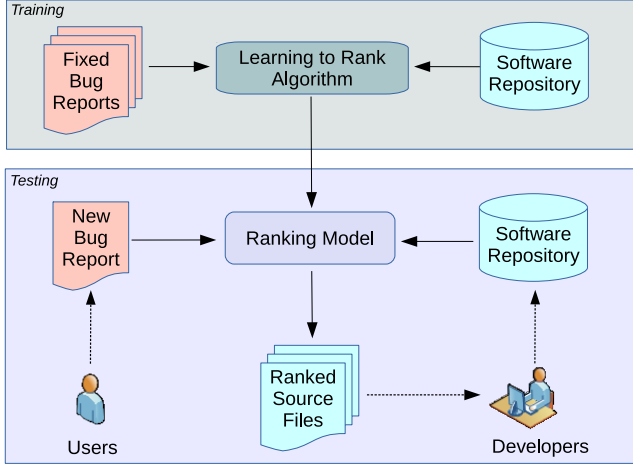


Figure 1: System architecture for training and testing.

A **learning-to-rank system** [43, 44] was designed to rank all the source code files for a given bug report. Figure 1 shows the high-level architecture of this system, in which the **ranking model** is trained to compute a matching score for any bug report r and source code file s combination. The scoring function $f(r, s)$, which is shown in Equation 1, is defined as a weighted sum of k features, where each feature $\phi_i(r, s)$ refers to a measurement of the relation between the source file s and the received bug report r :

$$f(r, s) = \mathbf{w}^T \Phi(r, s) = \sum_{i=1}^k w_i * \phi_i(r, s) \quad (1)$$

Given an arbitrary bug report r as input at test time, the model computes the score $f(r, s)$ for each source file s in the software project and uses this value to rank all the files in descending order. The user is then presented with a ranked list of files, with the expectation that files appearing higher in the list are more likely to be relevant for the bug report i.e., more likely to contain the cause of the bug.

The model parameters w_i are trained on previously solved bug reports using a learning-to-rank algorithm [16]. In this learning framework, the optimization procedure tries to find a set of parameters for which the scoring function ranks the files that are known to be relevant for a bug report at the top of the list for that bug report.

During feature engineering, we designed different types of features including features measuring the text similarity between bug reports and code in different granularity, features measuring the text similarity between bug reports and API references, features measuring the code change history, features about similar fix, and features measuring code complexity based on code-dependency graph [44]. The calculation of text similarity is based on the classic Vector Space Model (VSM) technique, which is widely used for measuring lexical similarity between documents. Later, to bridge the lexical gaps between bug reports and code, we used word embeddings, which are vector representations of words, to create two additional features for measuring their relations [45].

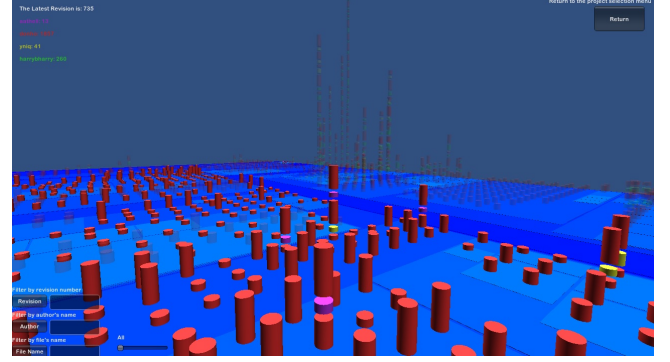


Figure 2: Screenshot of a 3-D software visualization tool

A **3-D software visualization tool** [21] was developed and was used for visualizing software project's *Subversion* (SVN) repository. Figure 2 shows a visualization view. Blue districts refer to folders. A stack of cylinders refers to a source code file. A single small cylinder refers to a revision of a source file. Different colors refer to different developers that contributed to the revisions.

This visualization tool is too old to be used in the proposed study because 1) it works on a very old version control system SVN, which is seldom used now; 2) it visualized all the files of a project in one view, which causes poor scalability and can hardly work on large projects. However, its visualization scheme and the 3-D city metaphor will be used in the proposed study.

1.2 C. Proposed Study

The main goal of the proposed project is to develop an interactive and accurate software defect positioning system that can be used by developers to locate software defects efficiently. This practical system is a continuing work of the PI's preliminary work. To achieve the project goal, this project will pursue several objectives as discussed in Section A.1. More specifically, the PI propose to do the following work.

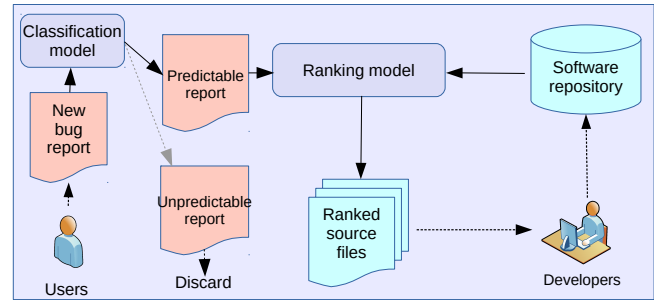


Figure 3: Ranking after pre-filtering.

1.2.1 C.1 Including a pre-filtering step to filter out “unpredictable” bug reports. Inspired by a similar work from Kim et al. [17], the PI propose to include a pre-filtering step in the ranking system. The purpose is to further increase the ranking accuracy by filtering out

“unpredictable” bug reports that do not contain sufficient information for locating the bug. As shown in Figure 3, a new received bug report will first be labeled by a classification model as either “predictable” or “unpredictable”. Then the ranking system will rank all the source code files for a “predictable” report in order to locate the defect. But it will keep silent for “unpredictable” reports as it has no confidence in working correctly for this types of reports.

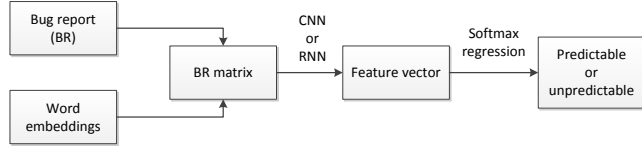


Figure 4: Bug report classification flowchart.

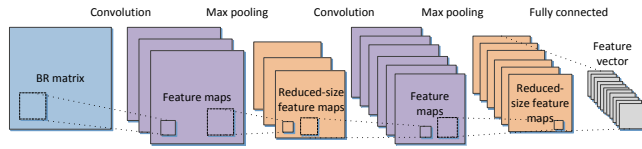


Figure 5: An example of a CNN model.

The flowchart of the classification model is shown in Figure 4, where BR stands for bug report. The PI propose to use deep neural networks such as CNN and RNN. First, the content including summary and description of a bug report will be converted into a matrix of real numbers by using word embeddings. Next, CNN or RNN is applied on the BR matrix to learn a feature vector, which will serve as the input to a softmax regression model. The output is a label of either “predictable” or “unpredictable”. The weight parameters of the softmax regression model is learned by using stochastic gradient descent on a dedicated training set with manually labeled bug reports. To create the training data, we will use the ranking results from the preliminary work [45]. For a bug report in the training set, it will be labeled as “predictable” if the preliminary result [45] successfully locate the bug within the top-10 ranked list, and it will be labeled as “unpredictable” otherwise.

Figure 5 shows an example of a CNN model [20] with two convolutional layers. The first convolutional layer uses a number of fix-size filters to perform convolution over the BR matrix and obtain the output of the same number of feature maps. Following is a max pooling layer that reduces the size of the feature maps. The output of the first layer serves as the input to the second convolutional layer. Finally, a fully connected layer produces a feature vector for the input BR matrix. It is noteworthy that when using CNN, all the training and testing BR matrices should have the same size. This can be done by extracting a fix size e.g. 100 words from all the bug reports. Reports with less words can be padded with 0. RNN works on input with variable length.

Recently, Mills et al. [28] proposed an approach, which is based on the Lucene implementation of VSM², to predict query quality for assisting text retrieval tasks in software engineering. We

²<https://lucene.apache.org/>

will perform comparison with this related work and may use techniques from it to improve the proposed system.

1.2.2 C.2 Introducing a new feature to measure the semantic similarity between bug reports and code. In preliminary work, we use Mihalcea’s function to aggregate word similarities into a document similarity. However, as discussed in Section A.1, it might cause information loss. Recently, Huo et al. [13] used CNN to learn feature vectors from bug reports and code for measuring their similarity. But CNN requires fix-length input data and therefore may also cause information loss on code.

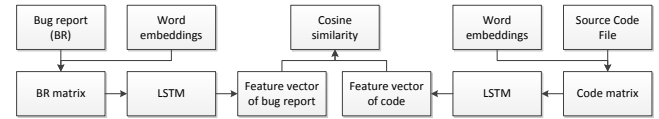


Figure 6: Using LSTM to calculate similarity.

The PI propose to use long short-term memory (LSTM) [11], a special type of RNN that can accommodate bug reports and source code files with variable length. Currently, LSTM networks are the most effective sequence learning models in many applications [9]. We will use the standard LSTM unit defined in [46]. As shown in Figure 6, LSTM will be used to learn feature vectors from bug reports and source code files. The cosine similarity of their vector representations is considered as a measurement of their semantic similarity and will be used as a new feature in the ranking model.

1.2.3 C.5 Evaluation. We will evaluate the software defect positioning system on several large-scale open-source projects that contain a sufficient number (more than 2,000) of source code files and previously fixed bug reports. We will conduct experiments on software projects that are written in Java as well as projects written in C/C++ because these programming languages are widely used in both industry and academic. We will use projects from the Eclipse foundation³ and Apache foundation⁴ because 1) both have many large-scale open-source projects written in Java and C/C++; 2) the source code packages can be easily downloaded from their GIT repositories; 3) their bug reports or issue reports are public accessible. More specifically, the projects used in our preliminary work [44] will be used in this study. Additionally, we will run experiments on more projects such as Apache HTTP Server⁵ written in C, Lucene⁶ (an information retrieval software library) written in Java, and Hadoop⁷ (a software framework for distributed storage and bigdata processing) written in Java. The selection of bug reports for evaluation is based on the same heuristics in [7, 43].

We will run the system to rank all the source code files for a given bug report and compare the result with the actual fix. The evaluation metrics such as *Mean Average Precision* (MAP) [23] used in our preliminary work [44] will also be used in this study. Additionally, we will use *Normalized Discounted Cumulative Gain*

³<https://eclipse.org/>

⁴<https://www.apache.org/>

⁵<https://httpd.apache.org/>

⁶<https://lucene.apache.org/>

⁷<http://hadoop.apache.org/>

(NDCG) [14], which is widely used in evaluating information retrieval models e.g. web search engines, to evaluate our system.

The PI is aware that user study is an effective way to evaluate the effectiveness of the proposed system in developers' real work. However, this is not the main goal of the proposed study. So the PI leave this to future work after the system is published. This study will also provide the basis for our future research on evaluating the effectiveness and usability of the proposed system in assisting teaching software engineering course at CSUSM.

1.3 D. Work Plan

The PI plan to hire one graduate student enrolled in the Master of Science Program in Computer Science (CS) and one undergraduate student enrolled in the Bachelor of Science in Computer Science program at CSUSM to help develop the system. The PI takes the overall responsibility of directing the project and keeps mentoring the students during the development.

In the first year, the graduate student will help implement the bug-report classification model, the LSTM-based semantic similarity feature, and the ranking model running on the server-side. The undergraduate student will implement the client-side programs and the web server that takes charge of the communication between the ranking model and the client-program. Since the ranking model, the database, and the web server work closely, the students will also work together closely during the development.

In the second year, two students will take charge of the maintenance and improvement of the system. Additionally, the graduate student will help run experiments to evaluate the ranking performance of the system on several open-source software projects, analyze the results, and improve the system accordingly. The undergraduate student will help develop a supplemental 3-D VR software visualization module. By the end of this year, the practical system will be published.

Recruitment of students will begin in spring, about three months before they start to work on the project. The PI will broadcast a hiring advertisement on CS-major electronic mailing lists, post a flier on class forums, and make a presentation of the project in the college-wide Frontiers in Science talk. In the coming 2017-2018 academic year, the CS department at CSUSM has a total of 866 undergraduate enrollments and a total of 39 graduate enrollments. Many students have desire to gain hands-on experience by working with faculty members on research projects. The PI will interview interested students. Preference will be given to economically-disadvantaged students, minority students, students making good progress toward their degree, and students who have demonstrated interest in this project.

1.4 E. Broader Impacts

The proposed study will result in a practical software system that alleviates developers' effort in bug finding and improves productivity. The system will be published. Any software developers can use and test it on their own projects. The system will be used in the software engineering courses at CSUSM to help students learn software engineering concepts in a lively manner. This project will expose students to an up-to-date software engineering research

topic as well as some cutting-edge techniques including artificial neural networks and virtual reality.

The research outcome of this project will provide the basis for the PI's future research in performing user study and collecting user feedback to evaluate the effectiveness of the system in both academic and industry. The experience learned from this project will be very helpful for PI's research in code recommendation and automatic programming. The methodology and techniques used in this project will contribute to the software engineering research community.

The PI graduated with a Ph.D. degree from Ohio University at May 2016 and joined the CS department at CSUSM as an tenure-track faculty at August 2016. The funding to the proposed study will help the PI obtain important computing resources for building a software engineering (SE) research lab at CSUSM to continue his research. This fund, which supports two student assistants, will also help the PI found an SE research group at CSUSM. The research group will work on adapting new techniques to solve SE tasks, applying up-to-date SE research outcome to assist teaching SE courses, and engaging students in doing SE research projects.

The funding will have a significant impact on the quality of education for students here at CSUSM. The funds requested for student assistants will allow some of our economically-disadvantaged students, many of whom work in retail and on campus dining, to have paid positions during the academic year and summer.

Furthermore, funding of the proposed research program will create new opportunities for our students that have strong desire to participate in research projects. The research opportunities created from the funding of the proposed research will directly contribute to providing talented undergraduates at CSUSM with the research experience necessary to be competitive applicants for top-tier Ph.D. programs.

2 EVALUATION

In this section, we describe an extensive set of experiments that are intended to determine the utility of the new document similarity measures based on word embeddings in the context of bug localization. This is an information retrieval task in which queries are bug reports and the system is trained to identify relevant, buggy files.

2.1 Text Pre-processing

There are three types of text documents used in the experimental evaluations in this section: 1) the Eclipse API reference, developer guides, Java API reference, and Java tutorials that are used to train the word embeddings; 2) the bug reports; and 3) the source code files. When creating bag-of-words for these documents, we use the same pre-processing steps on all of them: we remove punctuation and numerical numbers, then split the text by whitespace.

The tokenization however is done differently for each category of documents. In the one-vocabulary setting, compound words in the bug reports and the source code files are split based on capital letters. For example, "WorkbenchWindow" is split into "Workbench" and "Window", while its original form is also reserved. We then apply the Porter stemmer on all words/tokens.

Table 1: Benchmark Projects: *Eclipse** refers to *Eclipse Platform UI*.

Project	Time Range	# of bug reports used for testing	# of bug reports used for training	# of bug reports used for tuning	total
Birt	2005-06-14 – 2013-12-19	583	500	1,500	2,583
Eclipse*	2001-10-10 – 2014-01-17	1,656	500	1,500	3,656
JDT	2001-10-10 – 2014-01-14	632	500	1,500	2,632
SWT	2002-02-19 – 2014-01-17	817	500	1,500	2,817

In the two-vocabulary setting, a code token such as a method name “clear” is marked as “@clear@” so that it can be distinguished from the adjective “clear”. Then we stem only the natural language words. We also split compound natural language words. In order to separate code tokens from natural language words in the training corpus, we wrote a dedicated HTML parser to recognize and mark the code tokens. For bug reports, we mark words that are not in an English dictionary as code tokens. For source code files, all tokens except those in the comments are marked as code tokens. Inside the comments, words that are not in an English dictionary are also marked as code tokens.

2.2 Corpus for Training Word Embeddings

To train the shared embeddings, we created a corpus from documents in the following Eclipse repositories: the Platform API Reference, the JDT API Reference, the Birt API Reference, the Java SE 7 API Reference, the Java tutorials, the Platform Plug-in Developer Guide, the Workbench User Guide, the Plug-in Development Environment Guide, and the JDT Plug-in Developer Guide. The number of documents and words/tokens in each repository are shown in Table 2. All documents are downloaded from their official website⁸⁹. Code tokens in these documents are usually placed between special HTML tags such as `<code>` or emphasized with different fonts.

Table 2: Documents for training word embeddings.

Data sources	Documents	Words/Tokens
Platform API Reference	3,731	1,406,768
JDT API Reference	785	390,013
Birt API Reference	1,428	405,910
Java SE 7 API Reference	4,024	2,840,492
The Java Tutorials	1,282	1,024,358
Platform Plug-in Developer Guide	343	182,831
Workbench User Guide	426	120,734
Plug-in Development Environment Guide	269	90,356
JDT Plug-in Developer Guide	164	64,980
Total	12,452	6,526,442

Table 3: The vocabulary size.

Word embeddings trained on:	Vocabulary size
one-vocabulary setting	21,848
two-vocabulary setting	25,676

⁸<http://docs.oracle.com/javase/7/docs>

⁹<http://www.eclipse.org/documentation>

Table 4: Number of word pairs.

Approach	# of word pairs
One-vocabulary embeddings	238,612,932
Two-vocabulary embeddings	329,615,650
SEWordSim [37]	5,636,534
SWordNet [41]	1,382,246

To learn the shared embeddings, we used the Skip-gram model, modified such that it works in the training scenarios described in Section ?? Table 3 shows the number of words in each vocabulary setting. Table 4 compares the number of word pairs used to train word embeddings in the one- and two-vocabulary settings with the number of word pairs used in two related approaches. Thus, when word embeddings are trained on the one-vocabulary setting, the vocabulary size is 21,848, which leads to 238,612,932 word pairs during training. This number is over 40 times the number of word pairs in SEWordSim [37], and is more than 172 times the number of word pairs in SWordNet [41].

2.3 Benchmark Datasets

We perform evaluations on the fine-grained benchmark dataset from [?]. Specifically, we use four open-source Java projects: Birt¹⁰, Eclipse Platform UI¹¹, JDT¹², and SWT¹³. For each of the 10,000 bug reports in this dataset, we checkout a before-fixed version of the source code, within which we rank all the source code files for the specific bug report.

Since the training corpus for word embeddings (shown in Table 2) contains only Java SE 7 documents, for testing we use only bug reports that were created for Eclipse versions starting with 3.8, which is when Eclipse started to add Java SE 7 support. The Birt, JDT, and SWT projects are all Eclipse Foundation projects, and also support Java SE 7 after the Eclipse 3.8 release. Overall, we collect for testing 583, 1656, 632, and 817 bug reports from Birt, Eclipse Platform UI, JDT, and SWT, respectively. Older bug reports that were reported for versions before release 3.8 are used for training and tuning the learning-to-rank systems.

Table 1 shows the number of bug reports from each project used in the evaluation. The methodology used to collect the bug reports is discussed at length in [?]. Here we split the bug reports into a testing set, a training set, and a tuning set. Taking Eclipse Platform UI for example, the newest 1,656 bug reports, which were reported

¹⁰<https://www.eclipse.org/birt/>

¹¹<http://projects.eclipse.org/projects/eclipse.platform.ui>

¹²<http://www.eclipse.org/jdt/>

¹³<http://www.eclipse.org/swt/>

Table 5: MAP and MRR for the 5 ranking systems.

Project	Metric	LR+WE ¹ $\phi_1-\phi_8$	LR+WE ² $\phi_1-\phi_8$	LR $\phi_1-\phi_6$	WE ¹ $\phi_7-\phi_8$	WE ² $\phi_7-\phi_8$
Eclipse Platform UI	MAP	0.40	0.40	0.37	0.26	0.26
	MRR	0.46	0.46	0.44	0.31	0.31
JDT	MAP	0.42	0.42	0.35	0.22	0.23
	MRR	0.51	0.52	0.43	0.27	0.29
SWT	MAP	0.38	0.38	0.36	0.25	0.25
	MRR	0.45	0.45	0.43	0.30	0.30
Birt	MAP	0.21	0.21	0.19	0.13	0.13
	MRR	0.27	0.27	0.24	0.17	0.17

starting with Eclipse 3.8, are used for testing. The older 500 bug reports in the training set are used for learning the weight parameters of the ranking function in Equation ??, using the SVM^{rank} package [15, 16]. The oldest 1,500 bug reports are used for tuning the hyper-parameters of the Skip-gram model and the SVM^{rank} model, by repeatedly training on 500 and testing on 1000 bug reports. To summarize, we tune the hyper-parameters of the Skip-gram model and the SVM^{rank} model on the tuning dataset, then train the weight vector used in the ranking function on the training dataset, and finally test and report the ranking performance on the testing dataset. After tuning, the Skip-gram model was trained to learn embeddings of size 100, with a context window of size 10, a minimal word count of 5, and a negative sampling of 25 words.

2.4 Results and Analysis

We ran extensive experiments for the bug localization task, in order to answer the following research questions:

- RQ1*: Do word embeddings help improve the ranking performance, when added to an existing strong baseline?
- RQ2*: Do word embeddings trained on different corpora change the ranking performance?
- RQ3*: Do the word embedding training heuristics improve the ranking performance, when added to the vanilla Skip-gram model?
- RQ4*: Do the modified text similarity functions improve the ranking performance, when compared with the original similarity function in [26]?

We use the Mean Average Precision (MAP) [23], which is the mean of the average precision values for all queries, and the Mean Reciprocal Rank (MRR) [38], which is the harmonic mean of ranks of the first relevant documents, as the evaluation metrics. MAP and MRR are standard evaluation metrics in IR, and were used previously in related work on bug localization [34, 47?].

2.4.1 RQ1: Do word embeddings help improve the ranking performance? The results shown in Table 5 compare the **LR** system introduced in [?] with a number of systems that use word embeddings in the one- and two-vocabulary settings, as follows: **LR+WE¹** refers to combining the one-vocabulary word-embedding-based features with the six features of the LR system from [?], **LR+WE²** refers to combining the two-vocabulary word-embedding-based features with the LR system, **WE¹** refers to using only the one-vocabulary word-embedding-based features, and **WE²** refers to using only the two-vocabulary word-embedding-based features. The parameter vector of each ranking system is learned automatically.

The results show that the new word-embedding-based similarity features, when used as additional features, improve the performance of the **LR** system. The results of both **LR+WE¹** and **LR+WE²** show that the new features help achieve 8.1%, 20%, 5.6%, and 16.7% relative improvements in terms of MAP over the original **LR** approach, for Eclipse Platform UI, JDT, SWT, and Birt respectively. In [?], **LR** was reported to outperform other state-of-the-art bug localization models such as the VSM-based BugLocator from Zhou et al. [47] and the LDA-based BugScout from Nguyen et al. [30].

Another observation is that using word embeddings trained on one-vocabulary and using word embeddings trained on two-vocabulary achieve almost the same results. By looking at a sample of API documents and code, we discovered that class names, method names, and variable names are used with a consistent meaning throughout. For example, developers use *Window* to name a class that is used to create a window instance, and use *open* to name a method that performs an open action. Therefore, we believe the two-vocabulary setting will be more useful when word embeddings are trained on both software engineering (SE) and natural language (NL) corpora (e.g. Wikipedia), especially in situations in which a word has NL meanings that do not align well with its SE meanings. For example, since *eclipse* is used in NL mostly with the astronomical sense, it makes sense for *eclipse* to be semantically more similar with *light* than *ide*. However, in SE, we want *eclipse* to be more similar to *ide* and *platform* than to *total*, *color*, or *light*. By training separate embeddings for *eclipse* in NL contexts (i.e. *eclipse_NL*) vs. *eclipse* in SE contexts (i.e. *eclipse_SE*), the expectation is that, in an SE setting, the *eclipse_SE* embedding would be more similar with the *ide_SE* embedding than the *total_SE* or *color_SE* embeddings.

Kochhar et al. [18] reported from an empirical study that the localized bug reports, which explicitly mention the relevant file names, “statistically significantly and substantially” impact the bug localization results. They suggested that there is no need to run automatic bug localization techniques on these bug reports. Therefore, we separate the testing bug reports for each project into two subsets T1 (easy) and T2 (difficult). Bug reports in T1 mention either the relevant file names or their top-level public class names, whereas T2 contains the other bug reports. Note that, although bug reports in T1 make it easy for the programmer to find a relevant buggy file, there may be other relevant files associated with the same bug report that could be more difficult to identify, as shown in the statistics from Table 6.

Table 6 shows the MAP and MRR results on T1 and T2. Because **LR+WE¹** and **LR+WE²** are comparable on the test bug reports, here we compare only **LR+WE¹** with **LR**. The results show that both **LR+WE¹** and **LR** achieve much better performance on bug reports in T1 than T2 for all projects. This confirms the conclusions of the empirical study from Kochhar et al. [18]. The results in Table 6 also show that overall using word embeddings helps on both T1 and T2. One exception is Birt, where the use of word embeddings hurts performance on the 27 easy bugs in T1, a result that deserves further analysis in future work.

To summarize, we showed that using word embeddings to create additional semantic similarity features helps improve the ranking performance of a state-of-the-art approach to bug localization. However, separating the code tokens from the natural language words in two vocabularies when training word embeddings on the

Table 6: Results on easy (T1) vs. difficult (T2) bug reports, together with # of bug reports (size) and average # of relevant files per bug report (avg).

		T1		T2	
		LR+WE ¹	LR	LR+WE ¹	LR
Eclipse	Size/Avg	322/2.11		1,334/2.89	
	MAP	0.80	0.78	0.30	0.27
	MRR	0.89	0.87	0.36	0.33
JDT	Size/Avg	84/2.60		548/2.74	
	MAP	0.79	0.75	0.36	0.29
	MRR	0.90	0.87	0.45	0.37
SWT	Size/Avg	376/2.35		441/2.57	
	MAP	0.57	0.55	0.22	0.21
	MRR	0.66	0.65	0.27	0.26
Birt	Size/Avg	27/2.48		556/2.24	
	MAP	0.48	0.54	0.20	0.17
	MRR	0.62	0.69	0.25	0.22

SE corpus did not improve the performance. In future work, we plan to investigate the utility of the two-vocabulary setting when training with both SE and NL corpora.

2.4.2 RQ2: Do word embeddings trained on a different corpora change the ranking performance? To test the impact of the training corpus, we train word embeddings in the one-vocabulary setting using the Wiki data dumps¹⁴, and redo the ranking experiment.

Table 7: The size of the different corpora.

Corpus	Vocabulary	Words/Tokens
Eclipse and Java	21,848	6,526,442
Wiki	2,098,556	3,581,771,341

Table 8: Comparison of the LR+WE¹ results when using word embeddings trained on different corpora.

Corpus	Metric	Eclipse Platform UI	JDT	SWT	Birt
Eclipse and Java documents	MAP	0.40	0.42	0.38	0.21
	MRR	0.46	0.51	0.45	0.27
Wiki	MAP	0.40	0.41	0.38	0.21
	MRR	0.46	0.51	0.45	0.27

The advantage of using the Wiki corpus is its large size for training. Table 7 shows the size of the Wiki corpus. The number of words/tokens in the Wiki corpus is 548 times the number in our corpus, while its vocabulary size is 96 times the vocabulary size of our corpus. Theoretically, the larger the size of the training corpus the better the word embeddings. On the other hand, the advantage of the smaller training corpus in Table 2 is that its vocabulary is

¹⁴<https://dumps.wikimedia.org/enwiki/>

close to the vocabulary used in the queries (bug reports) and the documents (source code files).

Table 8 shows the ranking performance by using the Wiki embeddings. Results show that the project specific embeddings achieve almost the same MAP and MRR for all projects as the Wiki embeddings. We believe one reason for the good performance of the Wiki embeddings is the pre-processing decision to split compound words such as WorkbenchWindow that do not appear in the Wiki vocabulary into their components words Workbench and Window, which belong to the Wiki vocabulary. Correspondingly, Table 9 below shows the results of evaluating just the word-embeddings features (WE¹) on the Eclipse project with the two types of embeddings, with and without splitting compound words. As expected, the project-specific embeddings have better performance than the Wiki-trained embeddings when compound words are not split; the comparison is reversed when splitting is used. Overall, each cor-

Table 9: Project-specific vs. Wikipedia embeddings performance of WE¹ features, with and without splitting compound words.

Project	Metric	No Split	Split
Eclipse/Java	MAP	0.254	0.260
	MRR	0.307	0.310
Wikipedia	MAP	0.248	0.288
	MRR	0.300	0.346

pus has its own advantages: while the embeddings trained on the project-specific corpus may better capture specific SE meanings, the embeddings trained on Wikipedia may benefit from the substantially larger amount of training examples. Given the complementary advantages, in future work we plan to investigate training strategies that exploit both types of corpora.

2.4.3 RQ3: Do the word embedding training heuristics improve the ranking performance? Table 10 shows the results of using the original Skip-gram model without applying the heuristic techniques discussed in Sections ?? and ?. It shows that both the enhanced and the original Skip-gram model achieve the same results most of the time. These results appear to indicate that increasing the number of training pairs for word embeddings will not lead to further improvements in ranking performance, which is compatible with the results of using the Wiki corpus vs. the much smaller project-specific corpora.

2.4.4 RQ4: Do the modified text similarity functions improve the ranking performance? Table 11 below compares the new text similarity functions shown in Equation ?? with the original text similarity function from Mihalcea et al. [26], shown in Equation ?. In WE¹_{ori}, the new features ϕ_7 and ϕ_8 are calculated using the one-vocabulary word embeddings and the original *idf*-weighted text similarity function. The results of LR+WE¹ and LR are copied from Table 5, for which ϕ_7 and ϕ_8 are calculated using the new text similarity functions.

Results show that the new text similarity features lead to better performance than using the original text similarity function. The new features obtain a 20% relative improvement in terms of MAP

Table 10: LR+WE¹ results obtained using the enhanced vs. the original Skip-gram model.

Project	Metric	LR $\phi_1-\phi_8$	Enhanced Skip-gram $\phi_1-\phi_6$	Original Skip-gram $\phi_1-\phi_8$
Eclipse Platform UI	MAP	0.37	0.40	0.40
	MRR	0.44	0.46	0.46
JDT	MAP	0.35	0.42	0.42
	MRR	0.43	0.51	0.51
SWT	MAP	0.36	0.38	0.37
	MRR	0.43	0.45	0.44
Birt	MAP	0.19	0.21	0.21
	MRR	0.24	0.27	0.27

Table 11: Comparison between the new text similarity function (LR+WE¹) and the original similarity function (LR+WE¹_{ori}).

Project	Metric	LR $\phi_1-\phi_8$	LR+WE ¹ $\phi_1-\phi_6$	LR+WE ¹ _{ori} $\phi_7-\phi_8$
Eclipse Platform UI	MAP	0.37	0.40	0.37
	MRR	0.44	0.46	0.43
JDT	MAP	0.35	0.42	0.36
	MRR	0.43	0.51	0.45
SWT	MAP	0.36	0.38	0.37
	MRR	0.43	0.45	0.44
Birt	MAP	0.19	0.21	0.20
	MRR	0.24	0.27	0.25

over the LR approach, while features calculated based on the original text similarity function achieve only a 3% relative improvement.

3 EVALUATION OF WORD EMBEDDINGS FOR API RECOMMENDATION

To assess the generality of using document similarities based on word embeddings for information retrieval in software engineering, we evaluate the new similarity functions on the problem of linking API documents to Java questions posted on the community question answering (cQA) website Stack Overflow (SO). The SO website enables users to ask and answer computer programming questions, and also to vote on the quality of questions and answers posted on the website. In the Question-to-API (Q2API) linking task, the aim is to build a system that takes as input a user’s question in order to identify API documents that have a non-trivial semantic overlap with the (as yet unknown) correct answer. We see such a system as being especially useful when users ask new questions, for which they would have to wait until other users post their answers. Recommending relevant API documents to the user may help the user find the answer on their own, possibly even before the correct answer is posted on the website. To the best of our knowledge, the Q2API task for cQA websites has not been addressed before.

In order to create a benchmark dataset, we first extracted all questions that were tagged with the keyword ‘java’, using the datadump archive available on the Stack Exchange website. Of the 1,493,883 extracted questions, we used a script to automatically select only the questions satisfying the following criteria:

- (1) The question score is larger than 20, which means that more than 20 people have voted this question as “useful”.
- (2) The question has answers of which one was checked as the “correct” answer by the user who asked the question.
- (3) The “correct” answer has a score that is larger than 10, which means that more than 10 people gave a positive vote to this answer.
- (4) The “correct” answer contains at least one link to an API document in the official Java SE API online reference (versions 6 or 7).

This resulted in a set of high quality 604 questions, whose correct answers contain links to Java API documents. We randomly selected 150 questions and asked two proficient Java programmers to label the corresponding API links as *helpful* or *not helpful*. The remaining 454 questions were used as a (noisy) training dataset. Out of the 150 randomly sampled questions, the 111 questions that were labeled by both annotators as having *helpful* API links were used for testing. The two annotators were allowed to look at the correct answer in order to determine the semantic overlap with the API document.

Although we allow API links to both versions 6 and 7, we train the word embeddings in the one-vocabulary setting, using only the Java SE 7 API documentations and tutorials. There are 5,306 documents in total, containing 3,864,850 word tokens.

We use the Vector Space Model (VSM) as the baseline ranking system. Given a question T , for each API document S we calculate the VSM similarity as feature $\phi_1(T, S)$ and the asymmetric semantic similarities that are based on word embeddings as features $\phi_2(T, S)$ and $\phi_3(T, S)$. In the VSM+WE system, the file score of each API document is calculated as the weighted sum of these three features, as shown in Equation ???. During training on the 454 questions, the objective of the learning-to-rank system is to find weights such that, for each training question, the relevant (helpful) API documents are ranked at the top. During evaluation on the 111 questions in the test dataset, we rank all the Java API documents S for each question T in descending order of their ranking score $f(T, S)$.

Table 12: Results on the Q2API task.

Approach	MAP	MRR
VSM	0.11	0.12
VSM+WE	0.35	0.39

Table 12 shows the MAP and MRR performance of the baseline VSM system that uses only the VSM similarity feature, vs. the performance of the VSM+WE system that also uses the two semantic similarity features. The results in this table indicate that the document similarity features based on word embeddings lead to substantial improvements in performance. As such, these results can serve as an additional empirical validation of the utility of word

embeddings for information retrieval tasks in software engineering.

We note that these results are by no means the best results that we expect for this task, especially since the new features were added to a rather simple VSM baseline. For example, instead of treating SO questions only as bags of undifferentiated words, the questions could additionally be parsed in order to identify code tokens or code-like words that are then disambiguated and mapped to the corresponding API entities [35? ?]. Given that, like VSM, these techniques are highly lexicalized, we expect their performance to improve if used in combination with additional features based on word embeddings.

4 RELATED WORK

Related work on word embedding in NLP was discussed in Section ?? In this section we discuss other methods for computing word similarities in software engineering and related approaches for bridging the lexical gap in software engineering tasks.

4.1 Word Similarities in SE

To the best of our knowledge, word embedding techniques have not been applied before to solve information retrieval tasks in SE. However, researchers [12, 39, 41] have proposed methods to infer semantically related software terms, and have built software-specific word similarity databases [36, 37].

Tian et al. [36, 37] introduce a software-specific word similarity database called SEWordSim that was trained on StackOverflow questions and answers. They represent words in a high-dimensional space in which every element within the vector representation of word w_i is the Positive Pointwise Mutual Information (PPMI) between w_i and another word w_j in the vocabulary. Because the vector space dimension equals the vocabulary size, the scalability of their vector representation is limited by the size of the vocabulary. When the size of the training corpus grows, the growing vector dimension will lead to both larger time and space complexities. Recent studies [3, 27] also showed that this kind of traditional count-based language models were outperformed by the neural-network-based low-dimensional word embedding models on a wide range of word similarity tasks.

Howard et al. [12] and Yang et al. [41] infer semantically related words directly from comment-code, comment-comment, or code-code pairs without creating the distributional vector representations. They first need to map a line of comment (or code) to another line of comment (or code), and then infer word pairs from these line pairs. Similarly, Wang et al. [39] infer word similarities from tags in FreeCode. The main drawback of these approaches is that they rely solely on code, comments, and tags. More general free-text contents are ignored. Many semantically related words (e.g. “placeholder” and “view”) are not in the source code but in the free-text contents of project documents (e.g. the Eclipse user guide, developer guide, and API document shown in Figure ?? to Figure ??). However, these types of documents are not exploited in these approaches.

More importantly, all the above approaches did not explain how word similarities can be used to estimate document similarities. They reported user studies in which human subjects were recruited

to evaluate whether the word similarities are accurate. However, these subjective evaluations do not tell whether and how word similarities can be used in solving IR tasks in SE.

4.2 Bridging the Lexical Gap to Support Software Engineering Tasks

Text retrieval techniques have been shown to help in various SE tasks [10, 24]. However, the system performance is usually suboptimal due to the lexical gap between user queries and code [25]. To bridge the lexical gap, a number of approaches [2, 8, 25? ? ?] have been recently proposed that exploit information from API documentations. These approaches extract API entities referenced in code, and use the corresponding documentations to enhance the ranking results.

Specifically, McMillan et al. [25] measure the lexical similarity between the user query and API entities, then rank higher the code that uses the API entities with higher similarity scores. Bajracharya et al. [2] augment the code with tokens from other code segments that use the same API entries. Ye et al. [?] concatenate the descriptions of all API entries used in the code, and directly measure the lexical similarity between the query and the concatenated document. The main drawback of these approaches is that they consider only the API entities used in the code. The documentations of other API entities are not used. Figure ?? shows the Eclipse bug 384108. Figure ?? shows its relevant file *PartServiceImpl.java*. Figure ?? shows the description of an API entry *IPageLayout*. Although *IPageLayout* is not used in *PartServiceImpl.java*, its API descriptions contains useful information that can help bridge the lexical gap by mapping the term “view” in bug 384108 with the term “placeholder” in *PartServiceImpl.java*. Therefore, to bridge the lexical gap, we should consider not only the descriptions of the API entities used in the code but also all API documents and project documents (e.g. the user guide shown in Figure ?? and the developer guide in Figure ??) that are available.

Latent Semantic Indexing (LSI) and Latent Dirichlet Allocation (LDA) have been used in the area of feature location and bug localization. Poshyanyk et al. [31, 32] use LSI to reduce the dimension of the term-document matrix, represent code and queries as vectors, and estimate the similarity between code and queries using the cosine similarity between their vector representations. Similarly, Nguyen et al. [30] and Lukins et al. [22] use LDA to represent code and queries as topic distribution vectors. Rao et al. [33] compare various IR techniques on bug localization, and report that traditional IR techniques such as VSM and Unigram Model (UM) outperform the more sophisticated LSI and LDA techniques. These approaches create vector representations for documents instead of words and estimate query-code similarity based on the cosine similarity between their vectors. McMillan et al. [?] introduced a LSI-based approach for measuring program similarity, and showed that their model achieve higher precision than a LSA-based approach [?] in detecting similar applications. All these works neither measure word similarities nor try to bridge the lexical gap between code and queries.

5 FUTURE WORK

This is the future-work section.

6 CONCLUSION

This is the conclusion section.

REFERENCES

- [1] Alberto Bacchelli and Christian Bird. 2013. Expectations, Outcomes, and Challenges of Modern Code Review. In *Proceedings of the 2013 International Conference on Software Engineering (ICSE '13)*. Piscataway, NJ, USA, 712–721. <http://dl.acm.org/citation.cfm?id=2486788.2486882>
- [2] Sushil K. Bajracharya, Joel Osher, and Cristina V. Lopes. 2010. Leveraging Usage Similarity for Effective Retrieval of Examples in Code Repositories. In *Proceedings of the Eighteenth ACM SIGSOFT International Symposium on Foundations of Software Engineering (FSE '10)*. New York, NY, USA, 157–166.
- [3] Marco Baroni, Georgiana Dinu, and Germán Kruszewski. 2014. Don't count, predict! A systematic comparison of context-counting vs. context-predicting semantic vectors. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Baltimore, Maryland, 238–247. <http://www.aclweb.org/anthology/P14-1023>
- [4] Silvia Breu, Rahul Premraj, Jonathan Sillito, and Thomas Zimmermann. 2010. Information Needs in Bug Reports: Improving Cooperation Between Developers and Users. In *Proceedings of the 2010 ACM Conference on Computer Supported Cooperative Work (CSCW '10)*. New York, NY, USA, 301–310.
- [5] Bernd Bruegge and Allen H. Dutoit. 2009. *Object-Oriented Software Engineering Using UML, Patterns, and Java* (3rd ed.). Prentice Hall Press, Upper Saddle River, NJ, USA.
- [6] Raymond P. L. Buse and Thomas Zimmermann. 2012. Information Needs for Software Development Analytics. In *Proceedings of the 2012 International Conference on Software Engineering (ICSE 2012)*. Piscataway, NJ, USA, 987–996.
- [7] Valentin Dallmeier and Thomas Zimmermann. 2007. Extraction of Bug Localization Benchmarks from History. In *Proceedings of the Twenty-second IEEE/ACM International Conference on Automated Software Engineering (ASE '07)*. New York, NY, USA, 433–436.
- [8] Tathagata Dasgupta, Mark Grechanik, Evan Moritz, Bogdan Dit, and Denys Poshyvanyk. 2013. Enhancing Software Traceability by Automatically Expanding Corpora with Relevant Documentation. In *Proceedings of the 2013 IEEE International Conference on Software Maintenance (ICSM '13)*. Washington, DC, USA, 320–329.
- [9] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. 2016. *Deep Learning*. MIT Press. <http://www.deeplearningbook.org>.
- [10] Sonia Haiduc, Gabriele Bavota, Andrian Marcus, Rocco Oliveto, Andrea De Lucia, and Tim Menzies. 2013. Automatic Query Reformulations for Text Retrieval in Software Engineering. In *Proceedings of the 2013 International Conference on Software Engineering (ICSE '13)*. IEEE Press, Piscataway, NJ, USA, 842–851. <http://dl.acm.org/citation.cfm?id=2486788.2486898>
- [11] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long Short-Term Memory. *Neural Comput.* 9, 8 (Nov. 1997), 1735–1780. <https://doi.org/10.1162/neco.1997.9.8.1735>
- [12] Matthew J. Howard, Samir Gupta, Lori Pollock, and K. Vijay-Shanker. 2013. Automatically Mining Software-based, Semantically-similar Words from Comment-code Mappings. In *Proceedings of the 10th Working Conference on Mining Software Repositories (MSR '13)*. IEEE Press, Piscataway, NJ, USA, 377–386. <http://dl.acm.org/citation.cfm?id=2487085.2487155>
- [13] Xuan Huo, Ming Li, and Zhi-Hua Zhou. 2016. Learning Unified Features from Natural and Programming Languages for Locating Buggy Source Code. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence (IJCAI'16)*. AAAI Press, 1606–1612. <http://dl.acm.org/citation.cfm?id=3060832.3060845>
- [14] Kalervo Järvelin and Jaana Kekäläinen. 2002. Cumulated Gain-based Evaluation of IR Techniques. *ACM Trans. Inf. Syst.* 20, 4 (Oct. 2002), 422–446. <https://doi.org/10.1145/582415.582418>
- [15] Thorsten Joachims. 2002. Optimizing Search Engines Using Clickthrough Data. In *Proceedings of the Eighth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '02)*. New York, NY, USA, 133–142.
- [16] Thorsten Joachims. 2006. Training Linear SVMs in Linear Time. In *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '06)*. New York, NY, USA, 217–226.
- [17] Dongsun Kim, Yida Tao, Sungjun Kim, and Andreas Zeller. 2013. Where Should We Fix This Bug? A Two-Phase Recommendation Model. *IEEE Trans. Softw. Eng.* 39, 11 (Nov. 2013), 1597–1610.
- [18] Pavneet Singh Kochhar, Yuan Tian, and David Lo. 2014. Potential Biases in Bug Localization: Do They Matter?. In *Proceedings of the 29th ACM/IEEE International Conference on Automated Software Engineering (ASE '14)*. ACM, New York, NY, USA, 803–814. <https://doi.org/10.1145/2642937.2642997>
- [19] Thomas D. LaToza and Brad A. Myers. 2010. Hard-to-answer Questions About Code. In *Evaluation and Usability of Programming Languages and Tools (PLATEAU '10)*. New York, NY, USA, Article 8, 6 pages.
- [20] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. 1989. Backpropagation Applied to Handwritten Zip Code Recognition. *Neural Comput.* 1, 4 (Dec. 1989), 541–551. <https://doi.org/10.1162/neco.1989.1.4.541>
- [21] C. Liu, X. Ye, and E. Ye. 2014. Source Code Revision History Visualization Tools: Do They Work and What Would it Take to Put Them to Work? *IEEE Access* 2 (2014), 404–426. <https://doi.org/10.1109/ACCESS.2014.2322102>
- [22] Stacy K. Lukins, Nicholas A. Kraft, and Letha H. Etzkorn. 2010. Bug Localization Using Latent Dirichlet Allocation. *Inf. Softw. Technol.* 52, 9 (Sept. 2010), 972–990.
- [23] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. 2008. *Introduction to Information Retrieval*. Cambridge University Press, New York, NY, USA.
- [24] A Marcus and G Antoniol. 2012. On the use of text retrieval techniques in software engineering. In *Proceedings of 34th IEEE/ACM International Conference on Software Engineering, Technical Briefing*.
- [25] C. McMillan, M. Grechanik, D. Poshyvanyk, Chen Fu, and Qing Xie. 2012. Exemplar: A Source Code Search Engine for Finding Highly Relevant Applications. *Software Engineering, IEEE Transactions on* 38, 5 (Sept 2012), 1069–1087. <https://doi.org/10.1109/TSE.2011.84>
- [26] Rada Mihalcea, Courtney Corley, and Carlo Strapparava. 2006. Corpus-based and knowledge-based measures of text semantic similarity. In *Proceedings of the 21st national conference on Artificial intelligence (AAAI'06)*. AAAI Press, 775–780.
- [27] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient Estimation of Word Representations in Vector Space. In *Proc. of Workshop at ICLR '13*.
- [28] Chris Mills, Gabriele Bavota, Sonia Haiduc, Rocco Oliveto, Andrian Marcus, and Andrea De Lucia. 2017. Predicting Query Quality for Applications of Text Retrieval to Software Engineering Tasks. *ACM Trans. Softw. Eng. Methodol.* 26, 1, Article 3 (May 2017), 45 pages. <https://doi.org/10.1145/3078841>
- [29] Emerson Murphy-Hill, Thomas Zimmermann, Christian Bird, and Nachiappan Nagappan. 2013. The Design of Bug Fixes. In *Proceedings of the 2013 International Conference on Software Engineering (ICSE '13)*. Piscataway, NJ, USA, 332–341. <http://dl.acm.org/citation.cfm?id=2486788.2486833>
- [30] Anh Tuan Nguyen, Tung Thanh Nguyen, Jafar Al-Kofahi, Hung Viet Nguyen, and Tien N. Nguyen. 2011. A Topic-based Approach for Narrowing the Search Space of Buggy Files from a Bug Report. In *Proceedings of the 2011 26th IEEE/ACM International Conference on Automated Software Engineering (ASE '11)*. Washington, DC, USA, 263–272. <https://doi.org/10.1109/ASE.2011.6100062>
- [31] Denys Poshyvanyk, Yann-Gael Gueheneuc, Andrian Marcus, Giuliano Antoniol, and Vaclav Rajlich. 2007. Feature Location Using Probabilistic Ranking of Methods Based on Execution Scenarios and Information Retrieval. *IEEE Trans. Softw. Eng.* 33, 6 (June 2007), 420–432.
- [32] Denys Poshyvanyk, Andrian Marcus, Vaclav Rajlich, Yann-Gael Gueheneuc, and Giuliano Antoniol. 2006. Combining Probabilistic Ranking and Latent Semantic Indexing for Feature Identification. In *Proceedings of the 14th IEEE International Conference on Program Comprehension (ICPC '06)*. Washington, DC, USA, 137–148.
- [33] 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 (MSR '11)*. New York, NY, USA, 43–52.
- [34] R.K. Saha, M. Lease, S. Khurshid, and D.E. Perry. 2013. Improving bug localization using structured information retrieval. In *Automated Software Engineering (ASE), 2013 IEEE/ACM 28th International Conference on*. 345–355.
- [35] Siddharth Subramanian, Laura Inozemtseva, and Reid Holmes. 2014. Live API Documentation. In *Proceedings of the 36th International Conference on Software Engineering (ICSE 2014)*. ACM, New York, NY, USA, 643–652. <https://doi.org/10.1145/2568225.2568313>
- [36] Yuan Tian, D. Lo, and J. Lawall. 2014. Automated construction of a software-specific word similarity database. In *Software Maintenance, Reengineering and Reverse Engineering (CSMR-WCRE), 2014 Software Evolution Week - IEEE Conference on*. 44–53. <https://doi.org/10.1109/CSMR-WCRE.2014.6747213>
- [37] Yuan Tian, David Lo, and Julia Lawall. 2014. SEWordSim: Software-specific Word Similarity Database. In *Companion Proceedings of the 36th International Conference on Software Engineering (ICSE Companion 2014)*. ACM, New York, NY, USA, 568–571. <https://doi.org/10.1145/2591062.2591071>
- [38] Ellen M. Voorhees. 1999. The TREC-8 Question Answering Track Report. In *In Proceedings of TREC-8*. 77–82.
- [39] Shaowei Wang, David Lo, and Lingxiao Jiang. 2012. Inferring Semantically Related Software Terms and Their Taxonomy by Leveraging Collaborative Tagging. In *Proceedings of the 2012 IEEE International Conference on Software Maintenance (ICSM '12)*. IEEE Computer Society, Washington, DC, USA, 604–607. <https://doi.org/10.1109/ICSM.2012.6405332>
- [40] Chu-Pan Wong, Yingfei Xiong, Hongyu Zhang, Dan Hao, Lu Zhang, and Hong Mei. 2014. Boosting Bug-Report-Oriented Fault Localization with Segmentation and Stack-Trace Analysis. In *Proc. ICSME'14, To Appear*.

- [41] Jinqiu Yang and Lin Tan. 2012. Inferring semantically related words from software context. In *Mining Software Repositories (MSR), 2012 9th IEEE Working Conference on*. 161–170. <https://doi.org/10.1109/MSR.2012.6224276>
- [42] Xin Ye. 2016. *Automated Software Defect Localization*. Ph.D. Dissertation. Ohio University, Athens, OH. <https://etd.ohiolink.edu/>.
- [43] Xin Ye, Razvan Bunescu, and Chang Liu. 2014. Learning to Rank Relevant Files for Bug Reports Using Domain Knowledge. In *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering (FSE 2014)*. New York, NY, USA, 689–699. <http://dl.acm.org/citation.cfm?id=2337223>. 2337226
- [44] Xin Ye, Razvan Bunescu, and Chang 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. <https://doi.org/10.1109/TSE.2015.2479232>
- [45] Xin Ye, Hui Shen, Xiao Ma, Razvan Bunescu, and Chang Liu. 2016. From Word Embeddings to Document Similarities for Improved Information Retrieval in Software Engineering. In *Proceedings of the 38th International Conference on Software Engineering (ICSE '16)*. ACM, New York, NY, USA, 404–415. <https://doi.org/10.1145/2884781.2884862>
- [46] Wojciech Zaremba, Ilya Sutskever, and Oriol Vinyals. 2014. Recurrent Neural Network Regularization. *CoRR* abs/1409.2329 (2014). <http://dblp.uni-trier.de/db/journals/corr/corr1409.html#ZarembaSV14>
- [47] 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 (ICSE 2012)*. Piscataway, NJ, USA, 14–24. <http://dl.acm.org/citation.cfm?id=2337223.2337226>