

Learning a Graph-based Classifier for Fault Localization

Hao Zhong* & Hong Mei

Department of Computer Science and Engineering, Shanghai Jiao Tong University, China

Abstract Since software emerged, locating software faults has been intensively researched, culminating in various approaches and tools that have been applied in real development. Despite the success of these developments, improved tools are still demanded by programmers. Meanwhile, some programmers are reluctant to use any tools when locating faults in their development. The state-of-the-art situation can be naturally improved by learning how programmers locate faults. The rapid development of open-source software has accumulated many bug fixes. A bug fix is a specific type of comments containing a set of buggy files and their corresponding fixed files, which reveal how programmers repair bugs. Feasibly, an automatic model can learn fault locations from bug fixes, but prior attempts to achieve this vision have been prevented by various technical challenges. For example, most bug fixes are not compilable after checking out, which hinders analyzing bug fixes by most advanced static/dynamic tools. This paper proposes an approach called CLAFa that trains a graph-based fault classifier from bug fixes. CLAFa is built on a recent partial-code tool called GRAPA, which enables the analysis of partial programs by the complete code tool called WALA. Once GRAPA has built a program dependency graph from a bug fix, CLAFa compares the graph from the buggy code with the graph from the fixed code, locates the buggy nodes, and extracts the various graph features of the buggy and clean nodes. Based on the extraction result, CLAFa trains a classifier that combines Adaboost and decision tree learning. The trained CLAFa can predict whether a node of a program dependency graph is buggy or clean. We evaluate CLAFa on thousands of buggy files collected from four open-source projects: Aries, Mahout, Derby, and Cassandra. The f-scores of CLAFa achieves were approximately 80% on all projects.

Keywords Fault classifier, partial code analysis, bug fix analysis

Citation Hao Zhong and Hong Mei. CLAFa. *Sci China Inf Sci*, for review

1 Introduction

As man-made artifacts, software systems can contain bugs that return incorrect values, compromise security, or even crash the systems. In critical applications, bugs can incur huge losses. To improve the quality of software, the software engineering community has devoted much time and attention to bug location (see Section 6 for a detailed survey). Researchers have demonstrated the effectiveness of their proposed approaches, and several tools (*e.g.*, FindBugs [53] and PMD [7]) have been widely applied in modern software development.

Despite the success of their developments, further improvements are greatly demanded by programmers. For example, DiGiuseppe and Jones [30] considered that spectra-based approaches are less effective at locating multiple faults than single faults, as the execution of one fault can hinder the execution of other faults (*e.g.*, crashes). Although possible solutions to this problem have been proposed, their true effectiveness remains in questionable. For example, Abreu *et al.* [12] proposed an approach that locates multiple faults, but evaluated it on only the Siemens benchmark [32], in which faults are manually constructed and each faulty version has exactly one fault. As another example, Wang *et al.* [119] found that the effectiveness of information retrieval (IR)-based approaches is significantly reduced, when bug reports are poorly written or contain misleading descriptions. Johnson *et al.* [57]

* Corresponding author (email: zhonghao@sjtu.edu.cn)

argued that many programmers are dissuaded by the limitations of bug-detection tools, and prefer to rely on their own programming experiences in code debugging. Attempts to improve the state-of-the-art situation are ongoing (see Section 6 for details). In this paper, we explore bug detection by mining programming experiences from real bug fixes. Our approach locates multiple faults with more fidelity than the above dynamic approaches, and with finer accuracy than the above static approaches.

Software programmers often coordinate their development progress with source-code control systems [105]. After repairing bugs, programmers commit their changes to such systems. Meanwhile, the rapid progress of open-source software has promoted bug fixing by open-source communities, and many real bug fixes have accumulated. By examining the commit histories, researchers [117, 132] can identify the commits with repaired bugs, *i.e.*, the bug fixes. Mei and Zhang [80] advocated that big data analysis can revolutionize software automation. Following their visions, we believe that the accumulated data open new research opportunities for locating bugs. Indeed, several researchers [44, 78, 135] have conducted empirical studies on bug fixes to understand how programmers repair bugs and make code changes.

Our insight. Inspired by the state-of-the-art situation and the availability of bug fixes, we contemplated mining a graph-based fault classifier developed from thousands of real bug fixes. Data cleaning is a hot research hotspot in the data mining area [100]. The output quality of many mining algorithms depends on the input quality. If a tool cannot extract accurate code facts from the source files, it cannot produce accurate bug predictions. In the software engineering and programming language communities, source files are often presented as dependency graphs [90], which more easily present facts of the source files than raw texts. Meanwhile, similar software engineering tasks (*e.g.*, predicting bug fixes from commits [117]) can be assisted by classification techniques. Based on the above observations, we make the following proposition: In a proper presentation format (*e.g.*, dependency graphs), one can accurately display related facts of the buggy code, and can train a classifier to locate bugs.

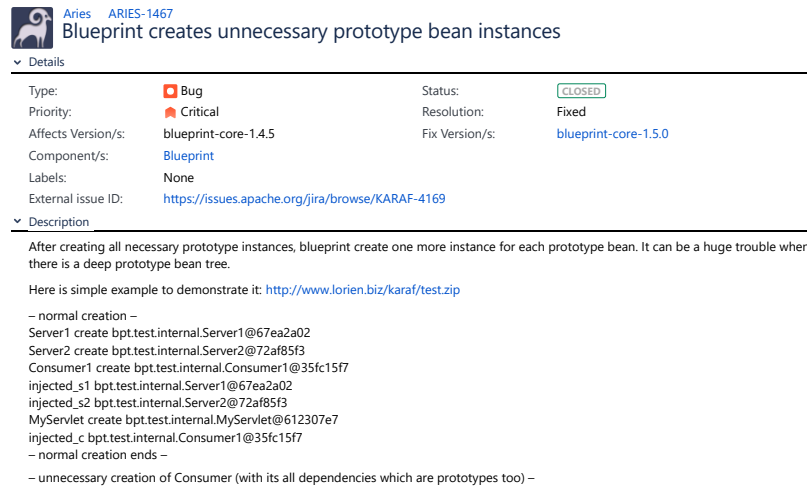
The challenges. To our knowledge, no previous researchers have attempted our proposition despite its benefits, because they are daunted by the following challenges:

Challenge 1. The first challenge is extracting the accurate facts from bug fixes. Most bug fixes are non-compilable, but most static tools (*e.g.*, WALA [9]) need compilable code. An empirical study [118] showed that only 38% of commits are compilable. When checking out the whole project of a commit, researchers must manually repair many compilation errors, before applying complete-code tools in analyses. Although bug signatures have been mined from buggy code [54] (see Section 6 for details), repairing the compilation errors in checked out buggy code requires much efforts; consequently, prior approaches typically analyze only a small number of known buggy samples. For example, Sun and Khoo [111] evaluated their approach only on the Siemens benchmark [55], in which the bugs are manually constructed and known. Li and Ernst [69] mined the bug signatures from only 53 bugs. This a small training set is insufficient for training a reliable classifier. Meanwhile, empirical studies on bug fixes (*e.g.*, [78, 135]) have analyzed only the buggy files and their corresponding fixed files, with partial-code tools (*e.g.*, [38]), which are too imprecise and lightweight [24, 83] to extract accurate facts from bug fixes. Without accurate facts, one cannot train an accurate classifier.

Challenge 2. Even when the facts can be accurately extracted from bug fixes, training a reasonably accurate classifier presents additional challenges. An empirical study [135] showed that faults account for only a small portion of source files; that is, most of the lines in a faulty source file are still clean. From a classifier perspective, the training set is quite imbalanced, because few locations are positive (contain faults). According to Yang and Wu [133], the problem of mining imbalanced data is one of ten challenging problems in data mining research. As most classifiers are designed for balanced data, their effectiveness is reduced when applied on imbalanced data.

Our contributions. To handle the first challenge, we build CLAFa on GRAPA and thereby analyze thousands of real bug fixes. GRAPA [136] extends the inference strategies of PPA [24], and resolves the remaining unknown code names using released binary files that are related to a partial program. GRAPA enables source-file analysis (even in the presence of compilation errors) by the state-of-the-art Java analysis tool WALA [9]. In this way, it builds program dependency graphs (PDGs) from bug fixes:

Definition 1. A PDG is defined as $g = \langle V, E_1, E_2 \rangle$, where V is a set of nodes corresponding to variables or expressions, and $E_1, E_2 \subseteq V \times V$ are two sets of edges. A $\langle s_1, s_2 \rangle \in E_1$ edge denotes a data dependency from s_1 to s_2 , and a $\langle s_1, s_2 \rangle \in E_2$ edge denotes a control dependency from s_1 to s_2 .



Aries ARIES-1467
Blueprint creates unnecessary prototype bean instances

▼ Details

Type:	■ Bug	Status:	CLOSED
Priority:	▲ Critical	Resolution:	Fixed
Affects Version/s:	blueprint-core-1.4.5	Fix Version/s:	blueprint-core-1.5.0
Component/s:	Blueprint		
Labels:	None		
External issue ID:	https://issues.apache.org/jira/browse/KARAF-4169		

▼ Description

After creating all necessary prototype instances, blueprint create one more instance for each prototype bean. It can be a huge trouble when there is a deep prototype bean tree.

Here is simple example to demonstrate it: <http://www.lorien.biz/karaf/test.zip>

```

- normal creation -
Server1 create bpt.test.internal.Server1@67ea2a02
Server2 create bpt.test.internal.Server2@72af85f3
Consumer1 create bpt.test.internal.Consumer1@35fc15f7
injected_s1 bpt.test.internal.Server1@67ea2a02
injected_s2 bpt.test.internal.Server2@72af85f3
MyServlet create bpt.test.internal.MyServlet@612307e7
injected_c bpt.test.internal.Consumer1@35fc15f7
- normal creation ends -
- unnecessary creation of Consumer (with its all dependencies which are prototypes too) -

```

Figure 1 The bug report of ARIES-1467

With GRAPA, researchers [123, 134] have built PDGs from bug fixes to explore open questions in automatic program repair and code change patterns. To handle the first challenge, CLAFa uses GRAPA to build PDGs from bug fixes that are not compilable. To handle the second challenge, we reduce the bias by a cost function that punishes a classifier, if it wrongly predicts a minority instance (Section 3.3.1). The major contributions of this paper are described below:

- We propose the first approach, called CLAFa (Classifying Faults), that predicts faulty nodes in a given PDG. The internal component of CLAFa is a classifier trained on thousands of bug fixes. During the training phase, CLAFa builds PDGs from the bug fixes, and extracts the feature vector of each node via graph analysis. To prepare labeled training data, CLAFa identifies the buggy nodes by comparing the PDGs of the buggy nodes with those of the fixed nodes. During the prediction phase, CLAFa uses its trained classifier to predict whether a node in a PDG is buggy or clean based on its extracted feature vector.
- We evaluated CLAFa on four popular open source projects. In total, we extracted the feature vectors of 19,343 nodes, covering 3,534 buggy methods. To ensure the reliability of our results, we performed within-project and cross-project fault predictions in a time-aware evaluation. CLAFa reasonably predicted the within-project faults (with f-scores 70%), but was less effective in cross-project fault prediction, unless the projects were similar. We further found that the best features can differ among projects, which explains why cross-project fault prediction is less effective than within-project fault prediction. Despite these differences, our results showed that bug reports and their called application programming interfaces (APIs) are more useful for locating faults than other features; also, that the outgoing control-dependent nodes of a node often determine whether the node is buggy or clean. Furthermore, our results revealed the best parameter and classification techniques, but their impacts proved subtle. In summary, our results reveal plenty of scope for improvement, especially by adding more features.

2 Motivating Example

Suppose that a programmer Mary needs to repair the reported bug [2] shown in Figure 1. If she locates the program's faults by a spectra-based approach, she must prepare many test cases that trigger both buggy and normal behaviors. CLAFa negates the need to prepare such test cases. Instead, Mary needs only to feed only the reported bug to CLAFa. If using IR-based approaches, Mary must manually analyze the faults in the obtained buggy files. Here she can be further assisted by CLAFa, which can locate the finer faults in buggy files.

Although CLAFa and IR-based approaches have the same inputs, they operate by different techniques. IR-based approaches compare the bug reports and source files to locate similar pairs. Instead of comparing the two types of files, CLAFa builds graphs from the code, and inspects their detailed features. The techniques of CLAFa are similar to those programmers who manually review code for bugs. During a code review, programmers often read the corresponding bug reports to understand the buggy behaviors. Similarly, CLAFa reads the bug reports first

```

1 private Map<String, Object> createInstances(...) {
2   ...
3   for (Map.Entry<String, Recipe> entry : ...) {
4     objects.put(entry.getKey(), entry.getValue().create
5       ());
6   }
7   return objects;
8 }

```

(a) The buggy file of ARIES-1467

```

1 private Map<String, Object> createInstances(...) {
2   ...
3   for (Map.Entry<String, Recipe> entry : ...) {
4     String name = entry.getKey();
5     ComponentMetadata component = blueprintContainer.
6       getComponentDefinitionRegistry().
7       getComponentDefinition(name);
8     boolean prototype = (component instanceof BeanMetadata
9       ) && MetadataUtil.isPrototypeScope((BeanMetadata)
10      component);
11    if (!prototype || names.contains(name)) {
12      objects.put(name, entry.getValue().create());
13    }
14  }
15  return objects;
16 }

```

(b) The fixed file of ARIES-1467

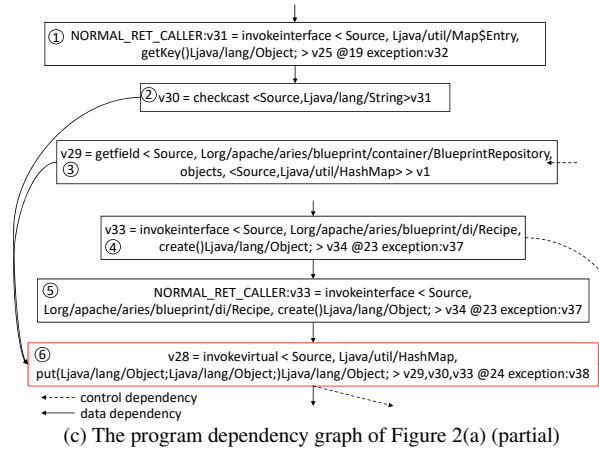


Figure 2 The source code of ARIES-1467

using its trained topic model. In the above example, CLAFa identifies the topics of the bug report in Figure 1 as “blueprint”, “service”, “dependency”, “server”, and “create”. Based on the topics, CLAFa then generates a vector (\vec{b}) presenting the bug report. In this vector, each bit denotes a topic, and “1” denotes that the bug report is related to the topic. After reading a bug report, programmers often look through source files to locate the bug. During this process, they can check many relevant aspects of the problem (*e.g.*, the called APIs). Similarly, CLAFa understands code by building PDGs from source files. Figure 2(a) shows the buggy code, and Figure 2(c) shows the built PDG. For each node of the PDG, CLAFa extracts a vector (\vec{c}) based on the node features (*e.g.*, node names), the local features (*i.e.*, features of the k -deep nodes before and after the node), and the global features (*i.e.*, features of all the nodes before and after the node). The node vector generated by CLAFa combines the \vec{b} and \vec{c} vectors. When input with the vector of the red-boxed node in Figure 2(c), the internal classifier of CLAFa predicts that the node is buggy (see Section 3 for details).

Mary does not need to understand the above details, because CLAFa automatically determines that the red-boxed node in Figure 2(c) is buggy. If she is unfamiliar with PDGs, CLAFa can locate the buggy lines through the WALA interface [8], which allows the location of the code elements corresponding to a given node of a PDG. In this example, the buggy node is mapped to Line 4 in Figure 2(a). Figure 2(b) shows the fixed code, which checks some conditions before calling Line 8 (the buggy line). This example of a bug report illustrates the common use of CLAFa. However, our evaluation results show that even without a bug report, CLAFa achieves reasonably high f-scores when identifying buggy nodes in buggy PDGs (see Section 4.3.6 for details).

Despite the differences, between IR-based approaches and CLAFa, the two approaches can be easily combined. For example, IR-based approaches (*e.g.*, [60]) can reduce the effort of locating buggy files prior to detecting their faults by CLAFa. In addition, CLAFa can be integrated with other tools such as automatic-program-repair tools, which locate faults by spectra-based approaches (*e.g.*, [47]). In each iteration, many candidate patches are generated. If none of these patches pass all test cases, automatic-program-repair approaches must locate new faults in all candidates by spectra-based approaches, which is time-consuming. As CLAFa does not need to execute test

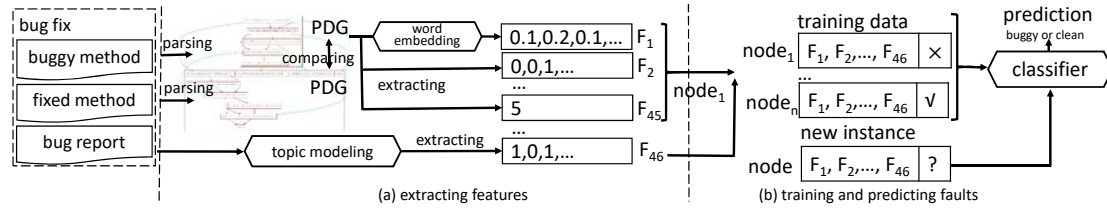


Figure 3 The overview of our approach

cases, it can significantly reduce the time of locating faults.

3 Approach

The term, *fault*, defines different granularities of buggy code. In IR-based approaches, a fault is a buggy file, whereas in spectra-based approaches, it often defines a buggy code line. Although spectra-based approaches can locate faults with finer granularity than IR-based ones, existing studies (*e.g.*, [92]) typically evaluate spectra-based approaches on code-line level, because code lines are more common than other granularities. In our approach, a fault is a buggy node in a PDG. When a code line calls multiple methods, each method invocation is encoded into a node. Consequently, a code line is often encoded into multiple nodes of a PDG, creating a finer granularity. Intuitively, locating finer faults is more challenging than coarser ones.

Our approach involves two stages: model training and fault localization. In the model training stage, a topic model is trained to classify bug reports, and a word embedding model is trained to encode code names into vectors. The faults are then located by a classification model. As the word embedding model and the topic model are learned by unsupervised approaches, they do not need labeled data. Meanwhile, as the classification model is learnt by a supervised approach, it requires labeled data with each node marked as buggy or clean.

Figure 3 shows the overview of our approach. Its major steps are (1) feature extraction (Section 3.2), and (2) training of the classification model and fault prediction (Section 3.3). In Figure 3, the word embedding model and the topic model are already learnt, and the labels for the classification model are prepared as follows. First, CLAFa builds a PDG (S_b) for the buggy method and a PDG (S_f) for its corresponding fixed method, and compares S_b and S_f to locate the modified nodes of S_b . When preparing the training data, it marks any modified node as buggy (\times), and any unmodified node as clean (\checkmark). Our underlying GRAPA tool detects a modified node, if its name or edges are changed. Modifications include the addition and deletion of code lines by programmers. For example, the buggy line in Figure 2(a) is not modified, but its control edges in Figure 2(c) are changed by the code lines added to Figure 2(b). Accordingly, CLAFa detects the modified node in Figure 2(c).

Table 1 lists our extracted features. Note that CLAFa extracts F_{46} from each bug report. For each node, CLAFa extracts a set of code features, and combines them with F_{46} to produce a vector. Once its classifier is trained, CLAFa can predict whether the vector of a node indicates a bug.

3.1 Data Acquisition

The data were acquired from Apache projects, in which bugs are easily identified. Zhong and Su [135] showed that the links of most Apache projects are carefully maintained, meaning that most bug fixes in Apache projects can be identified by extracting and comparing their issue numbers. The issue number is often included in the title of an Apache commit (*e.g.*, “ARIES-960” in Figure 4). In our setting, bug fixes can be identified without complicated techniques (*e.g.*, [117]). We determine that this commit is a bug fix, since its bug report [3] says that “ARIES-960” is a bug. As shown in Figure 4, a commit provides a link to its changed files. The commits are extracted and compared by our extension of the eGit tool [6]. It checks out the buggy files and fixed files of each bug fix.

The bug reports are collected by a web crawler based on XPath [18]. This technical choice reduces the effort of handling different styles of bug reports. Even when composed by the same open source community, bug reports can have subtle differences. For example, although Derby and Cassandra are both from the Apache foundation, the “Detail” section of a Derby report [5] includes an item called “Affects Version/s”, which is absent in a Cassandra report [4]. Our web crawler can extract bug reports of different styles after minor modifications on XPath queries.

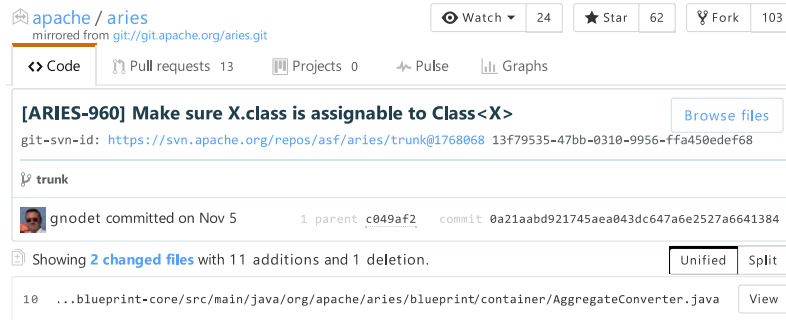


Figure 4 A sample commit

3.2 Feature Extraction

3.2.1 The analysis of bug reports

CLAFa combines the title, description, and comments of each bug report into a bag of words. Stop words (*e.g.*, “is”, “are”, and “would”) are removed, as these words are used in most bug reports and thus have little discriminative power. The remaining words are reduced to their roots by a stemming operation [72]. After that, CLAFa extracts the topics of each bug report using Latent Dirichlet Allocation (LDA) [86].

We choose LDA to classify bug reports, since it is a widely used technique to classify texts [86] including IR-based fault localization (*e.g.*, [87]). LDA treats each topic as a multinomial distribution of a vocabulary of w words. Topics are not predefined, but are learnt by unsupervised labeling. LDA models each document as a mixture of k latent topics, and produces a relevance score between a topic and a document. CLAFa considers that a document contains a topic if the relevance score exceeds 0.3. Based on the LDA results, CLAFa encodes a bug report as a k -dimensional vector, with each bit denoting a topic. If a bug report is related to a topic, the corresponding bit is set to 1; otherwise, it is set to 0.

3.2.2 The analysis of source files

Source files in CLAFa are analyzed by GRAPA [136], which extends the state-of-the-art WALA to build PDGs of bug fixes. Whereas WALA does not discriminate between the control dependencies and data dependencies of a PDG, GRAPA extract both types of dependencies. When CLAFa extracts the local and global features of a node (see Table 1) it identifies the nodes before and after the node in the PDG. The node labels are extracted by the Hungarian algorithm [63], which compares PDGs to locate their modified nodes (*i.e.*, buggy nodes).

3.2.3 The analysis of code names

In a built PDG, a node denotes an instruction that can call a method or access a field. For each node, CLAFa extracts code names by parsing its label. As such labels are generated by WALA, they follow specific name patterns. More than one code name can be extracted from a node name. For example, if a node denotes a method invocation, CLAFa extracts its return type, full method name, and parameter type names as its code names. For the buggy node in Figure 2(c), it extracts two code names such as `java/lang/Object` and `java/util/HashMap.put`. Here, we ignore duplicated code names. For each node, CLAFa identifies a main code name: the full method name is the main code name of a method invocation, and the field name is the main code name of a field access. We collect all code names, and use word embedding [82] to train a word model. With the model, CLAFa generates a vector for each code name.

3.2.4 Extracted Features

Table 1 lists our extracted features. For F_1 , CLAFa encodes the code names into vectors, with the support of word embedding [82]. For F_2 , CLAFa identifies the instruction types of nodes. In total, CLAFa identifies 33 node types (*e.g.*, `invokevirtual`, `getfield`, and `putfield`). The type of a node is denoted as one bit in a vector. As mentioned in Section 3.2.3, CLAFa can extract multiple code names from one node. F_3 and F_4 count its API code names and its client code names, respectively. Here, API code names are code names that are declared by third party libraries. For example, when analyzing bug fixes of Aries [1], the prefix `org/apache/aries` determines

Table 1 The features extracted by CLAFa

ID	Node Feature
F_1	Main code name (full method name or field name)
F_2	Node type
F_3	Number of API names
F_4	Number of client code names
F_5	Occurrences of code names in bug reports
ID	Local Feature
F_6	Number of k-depth incoming nodes
F_7	Number of k-depth outgoing nodes
F_8	Similar to F_2 but for k-depth incoming nodes
F_9	Similar to F_2 but for k-depth outgoing nodes
F_{10}	Similar to F_3 but for k-depth incoming nodes
F_{11}	Similar to F_3 but for k-depth outgoing nodes
F_{12}	Similar to F_4 but for k-depth incoming nodes
F_{13}	Similar to F_4 but for k-depth outgoing nodes
F_{14}	Similar to F_5 but for k-depth incoming nodes
F_{15}	Similar to F_5 but for k-depth outgoing nodes
	($F_6 - F_{15}$ are calculated for data-dependent nodes)
$F_{16} - F_{25}$	Similar to $F_6 - F_{15}$ but for control dependency
ID	Global Feature
$F_{26} - F_{45}$	Similar to $F_6 - F_{25}$ but the depth is maximized
ID	Bug Report Feature
F_{46}	Classification result

whether a code name is an API name. Wang *et al.* [119] showed that both IR-based approaches and programmers locate buggy files by searching code names in bug reports. As such code names are useful for locating bugs, we determine F_4 if code names in nodes overlap the code names in bug reports. Local features are derived from node features, and are extracted from the k -depth nodes before and after a current node. For example, when calculating F_6 of the red node in Figure 2(c) and the depth is set to one, we sum up the incoming data-dependent nodes of ②, ③, and ⑤, which have direct data dependencies to ⑥. As ① and ④ are the directly incoming nodes of the three data-dependent nodes, F_6 of ⑥ is 2, when the depth is one. The global features are the local features, whose values are obtained at maximum k -depths. For F_{46} , CLAFa encodes each bug report into a vector based on its topics. Here, the size of the vector is set to one hundred. Whereas other features capture the structure of the code, F_{46} sets the type of a reported bug.

When extracting the features from a bug report, the comments are ignored as they reveal the true locations of faults. For example, some faulty locations are fixed repetitively [43] and can be supplemented with comments added by programmers with true locations. When extracting our features from source files, we also ignore code comments. After these exclusions, no known labels exist in our extracted features. Our evaluation results are encouraging (see Section 4), but we expect that collecting more features in PDGs would improve the effectiveness. We further discuss this issue in Section 5.

3.3 Model Training and Bug Localization

3.3.1 Training the classifier

He and Garcia [50] divide the state-of-the-art solutions for imbalanced learning into sampling approaches and cost-sensitive approaches. Sampling approaches modify an imbalanced data set either by adding instances to minority classes (*e.g.*, [21]) or by removing instances from majority classes (*e.g.*, [71]). Although sampling is

Table 2 Subject.

Project	Single	Multiple	Graph	Fix	%
Aries	37	263	1192	394	76.1%
Mahout	47	253	1573	313	95.8%
Derby	32	268	1981	1134	26.5%
Cassandra	30	270	1468	2536	11.8%
Total	146	1,054	6,214	4,377	27.4%

easily understood, it can lead to overfitting or omission of important concepts [50, 79]. Cost-sensitive approaches (e.g., [112]) impose costs on misclassified instances. Determining whether a node is buggy or clean is a binary classification problem. Accordingly, we define $c(min, maj)$ as the cost of misclassifying a majority class instance as a minority class instance, and $c(maj, min)$ to denote the cost of the contrary case. In traditional classification techniques, both costs are the same. In cost-sensitive approaches, the cost of misclassifying minority instances is usually higher than the cost of misclassifying majority instances, i.e., $c(maj, min) > c(min, maj)$. He and Garcia [50] claimed that cost-sensitive approaches are superior to sampling approaches, because costs can be naturally imposed on imbalanced learning problems [126]. To balance our data, we therefore define the cost of a node (i_t) as follows:

$$cost(i_t) = \sum_i^n \frac{w_i}{n} w_t \frac{|C_m|}{|C|}, i_t \in C_m \quad (1)$$

where w_i denotes the weight of the i th node; w_t denotes the weight of node i_t ; C_m denotes all nodes (buggy and clean) in the m th class; and C denotes all nodes. For simplicity, we assign all w_i as one.

Our classifier comprises two classification techniques: the decision tree learning [39] and AdaBoost [40]. Decision tree is a supervised classification technique that classifies instances by constructing an if-else tree. Each interior node denotes a variable, and each leaf denotes a class. Adaboost is a meta-level learning technique that combines the outputs of weak classifiers into a weighted sum to predict the final output. The training set of CLAFa is a set of labeled data (\vec{f}_i, l_i) , where \vec{f}_i is the feature vector, and l_i is the label of a node. Adaboost repeatedly tunes its weights. We defined by $d_t(i)$ to denote the weight of the i th instance of the training data in the t th iteration. In the next iteration $t + 1$, the weight is updated as follows:

$$d_{t+1}(i) = \frac{d_t(i) \exp(-\alpha_t h_t(\vec{f}_i) l_i)}{z_t} \quad (2)$$

where $\alpha_t = \frac{1}{2} \ln(\frac{1-\varepsilon_t}{\varepsilon_t})$ is the weight updating parameter, $h_t(\vec{f}_i)$ is the prediction on feature vector \vec{f}_i , and z_t is a normalization factor that ensures that the all new weights sum to one. Here, ε_t is the error in the current model over the training set. After imposing cost $cost(i)$ on the i th instance, the above equation is modified to:

$$d_{t+1}(i) = \frac{d_t(i) \exp(-\alpha_t cost(i) h_t(\vec{f}_i) l_i)}{z_t} \quad (3)$$

3.3.2 Locating bugs

For simplicity, a trained classifier is considered as a function ($y = f(\vec{x})$), where \vec{x} denotes the vector of a node and y denotes the prediction. To predict whether a node is buggy, CLAFa extracts its \vec{x} and feeds the vector to its trained classifier. In particular, for a newly reported bug, CLAFa first uses its trained topic model to generate a feature vector (F_{46}). For each source file, it builds a program dependency graph, and for each node of the graph, it then extracts its node features, local features, and global features. Furthermore, CLAFa combines the above features into a vector (\vec{x}). Given \vec{x} as the input, CLAFa predicts whether the node is buggy or clean. In particular, if a prediction value of a node is greater than a threshold (0.5), CLAFa determines that the node is buggy.

Table 3 Overall effectiveness of CLAFa.

Project	Precision	Recall	F-score	The area under ROC
Aries	0.772	0.818	0.787	0.647
Mahout	0.856	0.888	0.871	0.619
Derby	0.902	0.924	0.912	0.647
Cassandra	0.892	0.917	0.903	0.65

4 Evaluation

4.1 Research Question

For a fair comparison, we must align inputs and outputs in a controlled experiment. The inputs of CLAFa differ from both spectra-based approaches (which require test cases), and the outputs of CLAFa differ from IR-based approaches (which cannot faults within source files). The different inputs and outputs preclude a fair comparison of CLAFa and prior spectra-based or IR-based approaches. Instead, we pose the following research questions:

(RQ1) How effectively does CLAFa detect faults in source files (Section 4.3.1)?

(RQ2) How does the local-feature depth affect the effectiveness (Section 4.3.2)?

(RQ3) How effectively does CLAFa detect faults, after learning from other projects (Section 4.3.3)?

(RQ4) What are the best discriminating features of buggy nodes (Section 4.3.4)?

(RQ5) How does CLAFa compare with other classification techniques (Section 4.3.5)?

(RQ6) What is the impact of bug reports (Section 4.3.6)?

4.2 Setup

4.2.1 Dataset

Table 2 shows the subjects of our evaluations (300 randomly selected bug fixes from each project). The “Single” and “Multiple” columns list the number of bug fixes that modify single methods and more than one method, respectively. Column “Graph” list the number of PDGs (over six thousand PDGs in total). Column “Fix” lists the total number of bug fixes. Column “%” is calculated as $\frac{300}{Fix}$.

4.2.2 Metric

Comparing our predictions against the gold standard of the faults (taken as the modified nodes of each bug fix), we classify all the nodes into false negatives (FNs), false positives (FPs), true negatives (TNs), and true positives (TPs). Based on the results, we measured our classifier by the following metrics:

$$precision = \frac{TP}{TP + FP} \quad (4)$$

$$recall = \frac{TP}{TP + FN} \quad (5)$$

$$f - score = \frac{2 \times precision \times recall}{precision + recall} \quad (6)$$

CLAFa classifies nodes into clean and buggy ones. For each type of nodes, we calculate the precisions, recalls, and f-scores, followed by their weighted averages. The weight of each type is the proportion of its nodes over the total number of nodes.

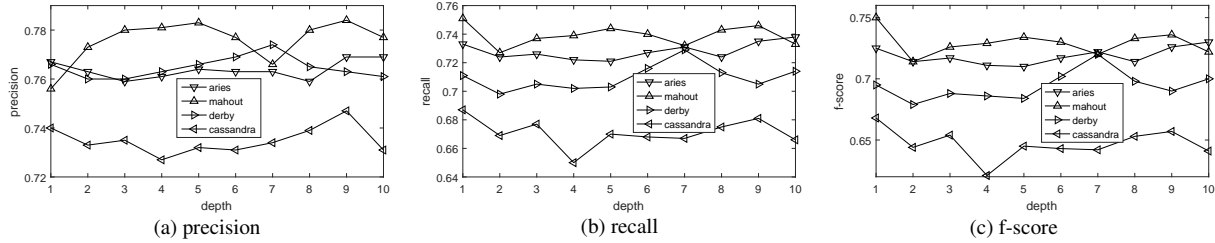


Figure 5 The impacts of our depth parameter

4.3 Empirical Result

4.3.1 RQ1. Effectiveness of CLAFa

Setting This research question explores the effectiveness of CLAFa in bug detection. For this purpose, we apply CLAFa to bug detection for each project listed in Table 2. Di Nucci *et al.* [28] mentioned that the traditional n -fold cross validation ignores the time sequences of the data instances. To avoid this problem, we conduct a time-aware evaluation [27, 48]. In particular, we sort the bug fixes in Table 2 by their issue numbers. After sorting, we use the top 80% bug fixes as the training data, and reserve the remaining bug fixes as the testing data. The effectiveness of CLAFa is evaluated by the precision, recall, f-scores, and the area under the receiver operating curve (ROC).

Results The training data and the testing data consumed approximately 1 gigabyte of storage per project. The classifier was trained after several hours, but once trained, it predicted the state of a node within seconds. Table 3 shows the overall results. Our CLAFa achieved high f-scores for all the projects. Although a fair comparison with spectra-based and IR-based approaches was infeasible, an indirect comparison can provide a reference. The prior studies [30, 31, 73] showed that spectra-based approaches can predict only one fault location effectively. We suspect that when a buggy file has multiple faults their f-scores shall be low. Meanwhile, the f-score of a recent IR-based approach [120] is 0.641. As CLAFa locates finer faults than IR-based approaches or spectra-based approaches, our results show that CLAFa already achieved reasonable accuracies. In summary, our results show that CLAFa is able to predict buggy nodes with reasonable precisions and recalls.

4.3.2 RQ2. Impact of the depth parameter

Setting As mentioned in Section 3.2.3, F_6 to F_{25} are extracted from the k deep incoming or outgoing nodes. To understand the impact of the k -depth parameter, we changed the k -depth from one to ten and investigated its impact on CLAFa. At each step, we recorded the precisions, the recalls, and the f-scores in the four projects, and analyzed the optimized depth.

Results The impacts of varying the k -depth are shown in Figure 5. The k -depth parameter little affected the precisions, recalls and f-scores of the four projects. The three measures followed different trends. For example, the f-scores of Aries and Derby were optimized at a depth of ten, whereas those of Mahout and Cassandra were maximized at unity depth. We suspect that the depth reflects the complexity of bugs in different projects. Most projects yielded favorable results at a depth of nine. Consequently, the depth was set to 9 in other research questions.

4.3.3 RQ3. Learning from other projects

Setting To explore this research question, we use the data of each project as the testing data, and the data from other projects and their combinations as the training data. We then analyze the changes in the f-scores under different settings.

Results Table 4 shows the f-scores during learning from other projects. Each row and column denotes the source of the testing and training data, respectively. For example, when the classifier was trained on Mahout and applied to the prediction of Aries faults, the obtained f-score was 0.445. Comparing the results of Tables 4 and 3, we find that the effectiveness was reduced when the model was learnt from the date of other projects. Column “combination”

Table 4 Learning from other projects.

Project	Aries	Mahout	Derby	Cassandra	Combination
Aries	n/a	0.445	0.431	0.485	0.459
Mahout	0.467	n/a	0.44	0.428	0.439
Derby	0.519	0.439	n/a	0.507	0.479
Cassandra	0.496	0.468	0.457	n/a	0.463

lists the f-scores when all the other projects were used as the training data. For example, when the data of Aries were predicted by the model that was trained on the data of Mahout, Derby, and Cassandra, the f-score was 0.459. Introducing more data for training did not always improve f-scores of the projects. In summary, our results show that predicting the buggy nodes of projects other than the training projects reduces the effectiveness of a trained model. Furthermore, the effectiveness may not be improved, by naively adding more training data. However, some training-testing pairs yielded higher f-scores than others. This suggests that when a project does not have previous bug fixes, CLAFa can be feasibly trained on the data of similar projects.

4.3.4 RQ4. Identification of important features

Setting This research question explores the features that contribute to our high effectiveness. The results will provide insights into the nature of bugs. In particular, we rank the features by their Pearson correlation coefficients [17]. We then compare the top ten features of all projects and selected those with commonality to all projects. This study was not intended to further tune the effectiveness of CLAFa, which can be better achieved by the correlation-based feature selection [45]. Indeed, even the default settings of the classifier yielded satisfactory results in RQ1. We believe that all features are valuable, and removing any one of them will reduce the overall effectiveness, as observed in Section 4.3.6.

Results Table 5 shows the top ten features in the four projects. As shown in Table 1, our features were derived from several basic features (*i.e.*, F_1 to F_7). To improve the presentation, Table 5 shows how features were derived from the features, instead of giving their number. For example, F_8 is written as “ $l, F_2, \pm, \leftarrow^d$ ” meaning that F_8 is derived from F_2 and is a local feature calculated from the incoming-data-dependent nodes. From Table 5, we obtain the following findings:

- 1. The node features alone cannot satisfactorily determine buggy nodes from clean ones.** In total, only two node features ranked among the top features, indicating whether a node is buggy is difficult to determine by inspecting only the node. This result highlights the importance of CLAFa, which accurately extracts local and global features. Without such features, the effectiveness of our approach can be significantly reduced.
- 2. The quality of bug reports is important.** In Table 5, F_5 and F_{46} are related to bug reports. As the effectiveness of IR-based approaches relies on the quality of bug reports [119], the quality of bug reports matters, when locating faults in buggy files. As an extreme case, seven of the top ten features in the Cassandra project were derived from F_5 , but this feature was non-dominant in other projects, confirming that only bug reports alone are insufficient for locating finer faults.
- 3. Calling APIs can introduce bugs.** Although reusing APIs significantly reduces the programming effort, existing empirical studies (*e.g.*, [135]) showed that many bugs are related to wrong API usages. Our results show that API-related features such as F_3 and F_1 are useful for locating faults. In particular, F_3 appeared in the top ten features of Aries and Derby, but was not ranked in Mahout. In fact, six of the ten top features in the Mahout project were derived from F_1 , which can be related to APIs (note that a code name can be an API code name).
- 4. Most of the features are related to outgoing nodes.** In all projects, most of the features were more related to outgoing nodes. As an extreme case, all top ten features of Derby were related to outgoing nodes. After inspecting some bug fixes, we identified two types of bug reports: outsider reports (often submitted by users) and insider reports (usually submitted by the project programmers). Although users have limited knowledge of the implementation details, their reports typically provide error messages for locating faults. As faults often appear before the messages, the outgoing nodes have stronger impacts on the fault localization than the incoming nodes. In contrast, the programmers submitting the insider reports fully understand the implementation details, so their

Table 5 Top ten features, ranked in the order of their importance in the clean versus buggy classification.

Rank	Aries	Mahout	Derby	Cassandra
1	g, F_5 , o, \leftarrow^c	F_{46}	g, F_5 , o, \leftarrow^c	g, F_5 , o, \leftarrow^c
2	l, F_5 , o, \leftarrow^c	g, F_1 , o, \leftarrow^d	g, F_3 , o, \leftarrow^c	g, F_5 , o, \leftarrow^c
3	l, F_7 , o, \leftarrow^c	l, F_1 , o, \leftarrow^d	l, F_5 , o, \leftarrow^c	g, F_5 , i, \leftarrow^c
4	l, F_6 , i, \leftarrow^c	l, F_1 , o, \leftarrow^c	l, F_3 , o, \leftarrow^c	g, F_5 , i, \leftarrow^c
5	g, F_3 , i, \leftarrow^c	g, F_4 , o, \leftarrow^c	l, F_7 , o, \leftarrow^c	l, F_5 , i, \leftarrow^c
6	l, F_3 , i, \leftarrow^c	g, F_1 , i, \leftarrow^c	g, F_2 , o, \leftarrow^c	F_{46}
7	l, F_2 , o, \leftarrow^c	n, F_1	l, F_2 , o, \leftarrow^c	l, F_5 , i, \leftarrow^d
8	g, F_2 , o, \leftarrow^c	l, F_1 , i, \leftarrow^c	g, F_4 , o, \leftarrow^c	n, F_1
9	l, F_3 , o, \leftarrow^c	g, F_5 , o, \leftarrow^c	l, F_4 , o, \leftarrow^c	g, F_4 , o, \leftarrow^c
10	g, F_3 , o, \leftarrow^c	g, F_1 , o, \leftarrow^c	g, F_2 , o, \leftarrow^c	l, F_7 , o, \leftarrow^c

n: node feature; l: local feature; g: global feature.

F_1 : main code name; F_2 : node type; F_3 : number of API code names; F_4 : number of client code names; F_5 : code names mentioned in bug reports; F_6 : in-degree; F_7 : out-degree; F_{46} : categories of bug reports.

i: incoming; o: outgoing.

\leftarrow^c : control dependency; \leftarrow^d : data dependency.

reports typically explain why faults occur. As such explanations often appear before faults, the incoming nodes have stronger impacts on the fault localization than the outgoing nodes. Other important factors, such as code structures and semantics, are also worthy exploring in future work.

5. Most of the features are related to control dependencies. Typically, control dependencies are related to code structures, whereas data dependencies are related to input and output values. For example, various approaches that detect legal call sequences of APIs (*e.g.*, [137]) are more related to control dependencies. Our results indicate that exploring more control-dependency bugs is a promising prospect.

In summary, whether a node is buggy or clean cannot be discerned from the node features alone. The features in bug reports and the called APIs are more effective in fault localization than other types of features. The outgoing control dependencies of a node can also usefully determine whether that node is buggy.

4.3.5 RQ5. Other classification techniques

Setting In this research question, we exchange our classifier with related classifiers. As the Adaboost boosts various classifiers, we eliminate it from this study, and compare the effectiveness of various internal classifiers. Many off-the-shelf classifiers are implemented in the WEKA suite of machine learning algorithms. Among these classifiers, we select the following for comparison with CLAFa:

1. SMO [96] is a sequential optimization algorithm that trains a support vector machine (SVM). An SVM [113] represents a data point as a vector, and the training process searches for one or more hyperplanes that split the data points. The hyperplanes are used for classifying new data points.

2. Naive Bayes [56] is a probabilistic classifier based on Bayes' theorem. During training, a naive Bayes classifier constructs a graph model representing the dependencies among observations, and then estimates the probability distributions of the observations. The graph model is used for predicting new data points.

3. Decision table [62] is a compact model for constructing complex rule sets and their corresponding actions. Given a set of data points, it constructs rules based the distributions of their features. The constructed decision table is used for predicting new data points.

4. Logistic [67] is a multinomial logistic regression model. During training, the regression model minimizes the estimated prediction errors. Given a set of independent variables, it predicts the probabilities of the categories of dependent variables.

These classifiers were selected because they are widely applied in software engineering papers (*e.g.*, [60, 117]). In this study, all classifiers were operated with their default settings. For different thresholds, we compared ROC and precision-recall (PR) curves of the five classifiers.

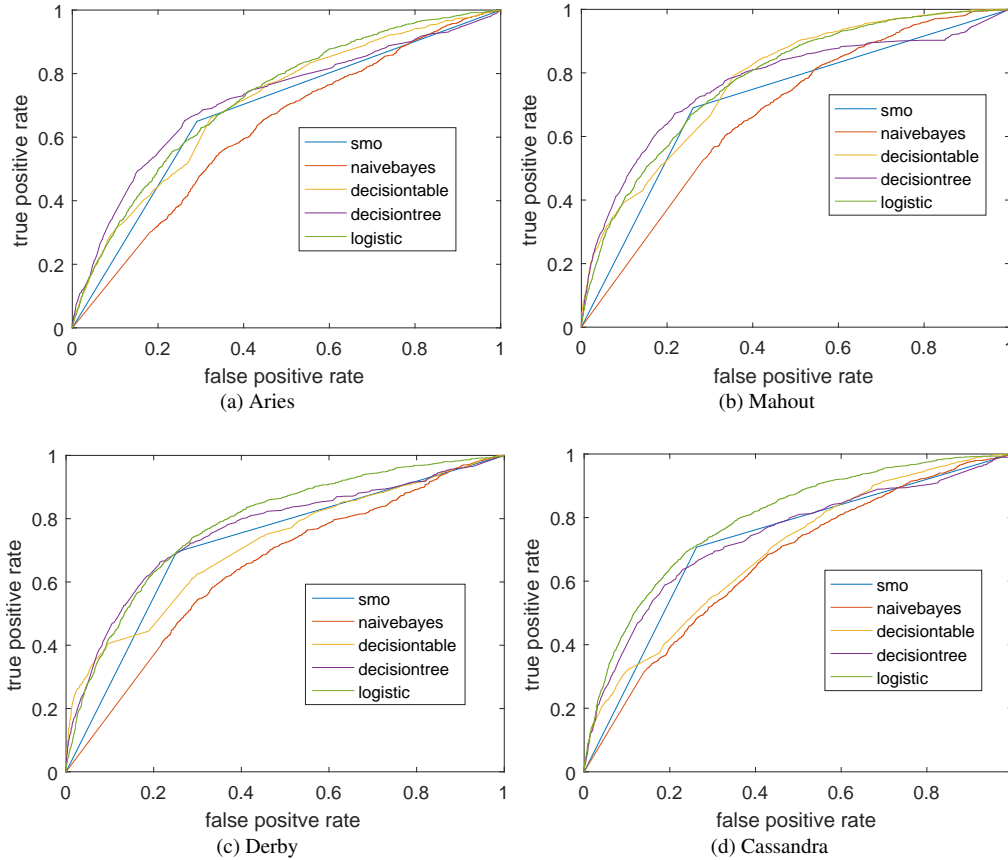


Figure 6 ROC

Results Figures 6 and 7 show the ROC curves and the PR curves, respectively. As mentioned in Section 3.3.1, the internal classifier of CLAFa is a decision tree. The findings are summarized below:

1. The logistic classifier outperformed the other classifiers on Cassandra. The top ten features of Cassandra differed from those of the other three projects (see Table 5). We suspect that the features in this project are intrinsically more suited to logistic regression than to other classifiers. When programmers apply CLAFa to bug detections in their own projects, they can tune its internal classifier to suite the nature of their debug features.

2. CLAFa was the optimum classifier on the other three projects. The left sides of our ROC and PR curves were above the left sides of the other classifiers. Although the curves intertwine, Fawcett's analysis [36] implies that the shape of our curves is optimal in the case of highly imbalanced data.

3. There remains much space for improvement. As the ROC areas are not maximized, all classifiers were sub-optimal. The ROC curves are intertwined, suggesting that multiple classifiers can be interpolated by existing techniques (*e.g.*, [37]). However, as the curves closely resemble each other, we suspect that the combining the classifiers will little improve the results. This issue is further discussed in Section 5.

Figures 6 and 7 show that the differences among several classifiers (*e.g.*, decision tree and logistic) were minor, and the classifiers could not be distinguished by their f-scores and AUC values. Consequently, we answer this research question by analyzing the ROC and PR graphs.

Ghotra *et al.* [42] evaluated the effectiveness of buggy-file prediction by different classifiers. Although their research goal differed from ours, some of their findings are consistent with ours. For example, they reported that logistic (SL) was one of the best classifiers, and smo was among the bottom classifiers. Other findings of Ghotra *et al.* [42] differed from ours, owing to the different settings in the two studies. Whereas we compared individual classifiers, they compared both individual classifiers and their combinations. They reported that some combinations are significantly better than individual ones (*e.g.*, Bag+J48).

In summary, CLAFa is the best classifier for Aries, Derby, and Mahout, and Logistic was the best classifier for Cassandra. However, the differences were minor, indicating that tuning classifiers may not achieve significant

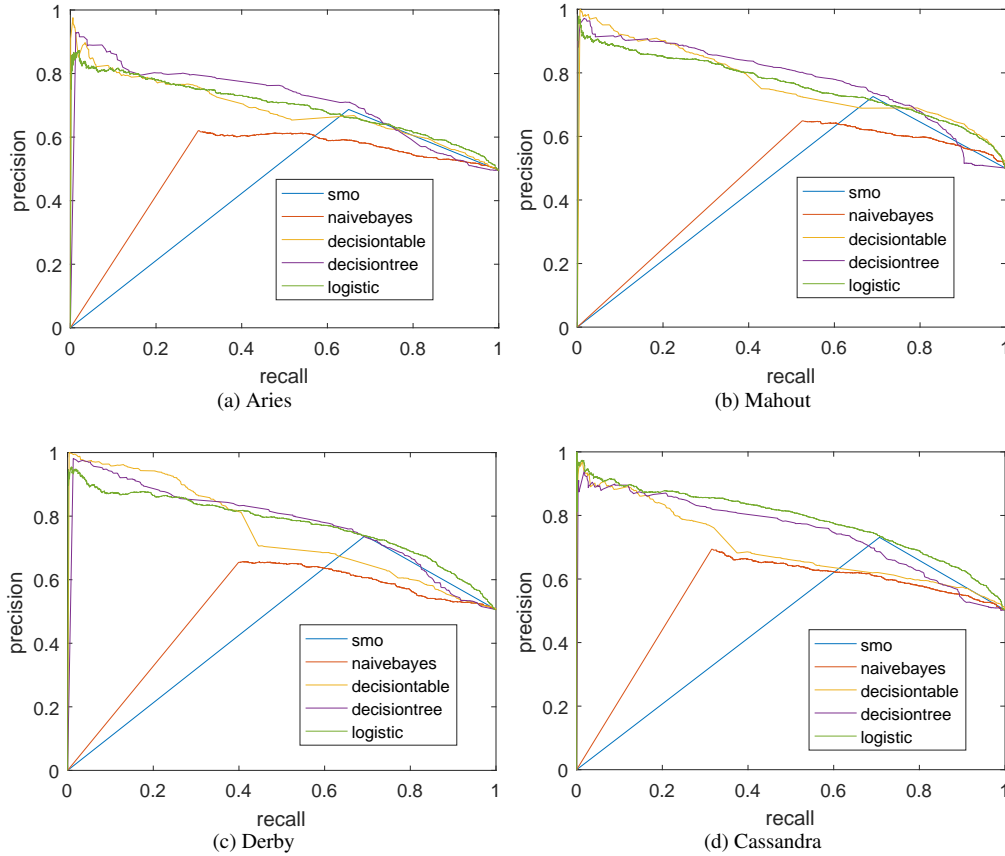


Figure 7 Precision-recall graph

improvement. This finding is unsurprising, because Hall *et al.* [46] also reported that simple models such as logistic achieved even better results than complicated models such as SVM. Instead, we expect that exploring more features is effective for better results.

4.3.6 RQ6. Impact of bug reports

Setting We supplement the graph features with the information from bug reports, which were used in existing approaches (*e.g.*, [103]). In Section 4.3.4, our results showed that the features from bug reports are important in fault localization, which somewhat overshadows the significance of CLAFa. In this research question, we reveal the true effectiveness of the graph features, removing all features from bug reports. Specifically, we removed F_{46} , F_5 from Table 1, and all features that were derived from F_5 . We then analyzed the precisions, recalls, and f-scores to understand the classification potential of the graph features.

Results Table 6 shows the results. Remove the features from bug reports reduced all the precision, recall, and f-score. Wang *et al.* [119] reported that the overall effectiveness of their compared IR-based approach [138] was severely degraded, when the bug reports lacked code entity names. CLAFa achieved reasonable high f-score even without inputs from bug reports (Table 6).

Although Table 5 showed that more top features were related to bug reports in Cassandra than other projects, the f-score on this project was less reduced by removing the bug-report features than f-scores of the other project. Note that Table 5 does not present the weight of each top feature. The results of learning classification model by code features alone were consistent (see Table 6), but the impacts of adding more features are apparently complicated.

Table 6 The results without bug reports.

Project	Precision	Recall	F-score
Aries	0.691	0.664	0.651
Mahout	0.71	0.679	0.668
Derby	0.735	0.672	0.652
Cassandra	0.706	0.657	0.631

4.3.7 Threats to validity

The threats to internal validity include the experimental bias of labeling faults in buggy files. For example, when programmers fix a bug, they somewhat also modify clean code lines (*e.g.*, refactoring). Another threat is missed bug fixes in commits. Although such cases can be rare, they will negatively impact on the results. The threat could be reduced by introducing manual inspection in future work. Meanwhile, the threats to external validity include the selected subjects, which are all open-source projects. This threat could be reduced by replicating our studies on commercial projects.

5 Discussion and Future Work

Exploring more basic features. As shown in Table 1, our basic features are limited. Previous researchers explored other data sources such as stack traces [128], version histories [110], and commit activities [61]. They also proposed useful metrics [61] for predicting buggy code. When we introduce more features in future work, we plan to borrow ideas from these approaches. To handle the complexity of features, we must allow features with variable lengths. For example, mined specifications (*e.g.*, [137]) are widely used in the detection of bugs related to wrong API sequences. To detect such bugs, the API call sequences can be encoded into feature vectors. As client code methods can call different APIs, feature vectors of the encoded API call sequences are length-variable. To fully leverage such features in future, we plan to tune our classification technique or introduce more advanced techniques.

Definition of bugs and its measures. Outside Apache projects, programmers can fail to maintain the links between the bug reports and their fixes [16]. Even if programmers maintain the links, they can define bugs in different ways, and their definitions can be disputed by researchers [15, 51]. These inconsistent definition can lead to different labels of faulty nodes. Here, the effectiveness of CLAFa was evaluated only on the definitions of programmers. Evaluate CLAFa under other definitions would help to assess the rationality of these definitions. A reasonable definition shall be consistent, and should therefore yield better results than a dubious definition. We plan to explore this issue in future work.

Integrating with automatic program repair. Automatic program repair (*e.g.*, [125]) has recently become a research hotspot, but remains limited in scope (*e.g.*, [99]). Automatic program repair typical invokes many iterations of spectra-based fault localization to locate the buggy lines of mutated candidates. This time-consuming process might partially explain why only limited number of bug fixes by the prior approaches. CLAFa reduces the time of fault localization by removing the need to execute test cases. In addition, spectra-based approaches cannot effectively to locate multiple faults. CLAFa resolves this problem by eliminating the interference among multiple bugs. In future work, we plan to integrate CLAFa with automatic program repair. An automatic program-repair approach must know the nodes to be modified, but modified nodes do not necessarily indicate bugs. To assist bug discrimination by the prior approach, we intend to locate the bug-causing nodes in an extended version of CLAFa.

Tuning CLAFa. Although CLAFa achieved high f-scores in our evaluation analysis, the results were obtained under the default settings of classifiers. The effectiveness of CLAFa could be further improved by fine-tuning. He and Garcia [50] recommended cost-sensitive approaches for handling imbalanced data, but sampling approaches might be more effective for our specific application. The literature is replete with approaches the automatically tune the parameters of a classifier [108, 114–116]. In future work, we plan to tune CLAFa using these approaches.

Combining with other fault-localization approaches. Bug finding bugs has been widely researched. Even the 100-plus papers cited in Section 6 did not cover all approaches; moreover, many new approaches are proposed

each year. As mentioned in Section 4.1, a fair controlled experiment is precluded because the inputs and outputs of CLAFa differ from those of other approaches. Despite the large number of available approaches, many programmers prefer to their own debugging experience [57]. To handle the problem, a feasible way is to combine existing approaches, instead of arguing the best approach. For example, Le *et al.* [66] combined IR-based and spectra-based approaches to improve the fault-localization results. In future work, we will explore the combination of CLAFa and other approaches in a real development context.

Words and topics outside the vocabulary. Our present evaluation trained our topic model and word embedding on all subjects, which excludes words or topics outside the vocabulary. In real usage, a bug report can contain new words or topics that never appeared in the previous histories, which confound the trained models. This problem could be feasibly solved by to re-training the models on the new bug reports. Alternatively, the new words and topics could be replaced with similar ones in the vocabulary. In future work, we plan to explore both approaches.

6 Related Work

Spectra-based fault localization. A typical spectra-based approach calculates the suspicious scores of buggy lines based on passed and failed tests. These approaches assume that if a code line often appears in failed tests, it probably contains a bug. This research topic has been intensively studied, and some early approaches (*e.g.*, [109]) were published in 1980s [130]. Various research metrics (*e.g.*, [58, 85, 129]) can calculate suspicious values of code lines, and their effectiveness has been explored in various empirical studies (*e.g.*, [10, 11]). Besides these explicitly defined metrics for suspicious buggy code, researchers have explored the automated mining of metrics. For example, Wong and Qi [131] learnt a neural network for buggy code prediction. Approaches other than metrics have also improved the state-of-the-art situation. Hao *et al.* [47] proposed an approach that reduces the number of test cases while minimizing the impacts on fault localization. Mao *et al.* [75] remove irrelevant code by slicing before calculating suspicious buggy code. However, most of the published papers assume that one bug in each buggy file. Although some approaches (*e.g.*, [12, 29]) locate multiple faults, these multi-faults are typically assumed as independent [130]. Abreu *et al.* [12] proposed an approach that locates multiple faults, but their approach is evaluated on only the Siemens benchmark [32], on which faults are manually constructed and each faulty version contains exactly one fault. Although Abreu *et al.* [12] combined multiple versions to simulate multiple faults, the combination again assumed independent multiple faults, which do not represent true faults. Gao and Wong [41] proposed another approach that locates multiple faults, but their evaluation also generated multiple faults from single fault versions. Pearson *et al.* [92] argued that artificial faults cannot confirm the true effectiveness of spectra-based localizations. In particular, their collected traces may not reflect the interferences among multiple faults observed in empirical studies (*e.g.*, [26, 30]). As the problem lies in the underlying assumption, many researchers (*e.g.*, [26, 30]) consider that spectra-based approaches cannot insufficiently to locate multiple faults, although at least one fault can top a suspicious list. Perez *et al.* [93] claimed that over 82% of bug fixes are single-fault fixes. In their definition, a fix is single-faulted, if all tests affect at least one changed component of the fix [93]. They considered that single faults can modify multiple locations, which is inconsistent with our definitions. Moreover, the best results from spectra-based approaches are typically obtained only after many high quality test cases. Just *et al.* [59] showed that if the test cases are of superior quality, the time of generating correct patches can be significantly reduced. Spectra-based approaches have been adapted with fewer test inputs [20, 94], but our approach focuses on a precise static analysis rather than test coverage. The multi-fault detection by our approach, which needs no test inputs, can complement spectra-based approaches.

IR-based fault localization. Typical IR-based approaches locate buggy files by comparing the bug reports with the source files. The bug reports are assumed to share similarity with their corresponding buggy files. The similarity is detected by various machine learning techniques such as TF-IDF [103], LDA [74], naive Bayes [60], and SVM [122, 138]. The early papers handled source files such as natural language texts, but more recent papers (*e.g.*, [106]) have extracted the code structures (*i.e.*, class, methods, variables, and comments) from source files to improve fault-localization accuracy. Besides natural language descriptions of bug reports, other data sources such as stack traces [128], version histories [110], and their combination [121] have been explored. Wang *et al.* [119] argued that the similarity lies in the code names that appear in both bug reports and source files. CLAFa is not an

IR-based approach, as it does not calculate the similarity between bug reports and source files, although it extracts features from bug reports. Being built on accurate source analysis, CLAFa can determine the buggy locations in source files, completing the IR-based approaches.

Model-based fault localization. A model can define either legal usages (specs) or illegal usages (bug signatures). Ammons *et al.* [14] mined automata for APIs. Pandita *et al.* [91] refined their approach, while other researchers [23, 88, 89, 137] mined graphs for specs. Robillard *et al.* [104] showed that automata and graphs are equivalent. These types of model-based fault localization can be reduced to the grammar inference problem. Method-pair extraction, proposed by Li and Zhou [70], has been improved in more complicated contexts [107]. The approach of Engler *et al.* [34], which extracts frequent call sequences, has also been improved by advanced techniques [102, 124]. Furthermore, mined sequences have been encoded as temporal logic [68, 76]. This research line can be reduced to sequence mining [13]. All of the above approaches are based on call sequences. Ernst *et al.* [35] inferred invariants to define the variable rules. More informative specs can be obtained by combining the invariants with sequences [65], and spec-mining has been enriched in various test cases [25, 98]. Marc and David [19] mined performance models from runtime traces. Zhong and Mei [134] empirically analyzed several open questions (*e.g.*, the best formats of specs). Meanwhile, bug signatures have been mined in sequences [54] and graphs [22, 69, 95, 111, 139]. Most of these approaches analyzed traces [54, 69, 111, 139]; only a few approaches [69] analyzed the buggy source files. Hsu *et al.* [54] and Sun and Khoo [111] applied a frequency-based mining approach; Li and Ernst [69] extracted subgraphs; and Cheng *et al.* [22] applied discriminative graph mining. These approaches differ from our approach by applying different techniques and taking different inputs from our approach.

Bug prediction. Studies on bug prediction have applied well-known metrics [33] or self-defined metrics [77]. In a bug prediction, a bug is assumed to cause a metric violation. However, defining the metrics is a challenging task. Nagappan *et al.* [84] predicted module-level bugs using several metrics, but none of the metrics yielded high effectiveness on all projects. CLAFa is not built on known metrics. Bug prediction can also be based on commit activities [48, 61], assuming that bugs can be identified by specific commit activities (*e.g.*, bursting). Rahman *et al.* [101] showed that bug predicting by activities is not substantially better than counting the number of commits touching a source file. CLAFa differs from the above approaches in both its underlying assumptions and technical details. In addition, most of these approaches predict faults on courser levels (*e.g.*, modules and classes) [46].

Classification in software engineering. Classification is an intensively studied technique in data mining, and commonly assists software engineering tasks such as requirement elicitation [49], development [127], testing [64], debugging [97], maintenance [52], and software reuse [81]. CLAFa is designed mainly for debugging and maintenance. Because it accurately analyzes source code, CLAFa predicts the faults within source files, complementing existing approaches (*e.g.*, [117]).

7 Conclusion

Open-source communities accumulate many real bug fixes. Locating code faults using knowledge mined from bug these fixes, and is desired, but is rendered challenging by various technical limitations. This paper proposed an approach called CLAFa that combines graph analysis and classification to locate multiple faults in source files. First, it builds program dependency graphs from the bug fixes and compares them to detect buggy nodes. The buggy nodes are used to label the faults. CLAFa extracts the graph features of each node, and denotes them in a vector. A classification model is then trained on the labeled data. We have evaluated CLAFa on thousands of real bug fixes extracted from four popular open-source projects. Our results confirmed that CLAFa can locate multiple faults with reasonably high accuracy.

Acknowledgments

We appreciate the anonymous reviewers for their constructive comments. This work is sponsored by the National Key R&D Program of China No. 2018YFC0830500, the National Nature Science Foundation of China No. 61572313, and the grant of Science and Technology Commission of Shanghai Municipality No. 15DZ1100305.

Conflict of interest The authors declare that they have no conflict of interest.

References

- 1 Apache Aries. <http://aries.apache.org>.
- 2 ARIES-1612. <https://issues.apache.org/jira/browse/ARIES-1612>.
- 3 ARIES-960. <https://issues.apache.org/jira/browse/ARIES-960>.
- 4 CASSANDRA-2044. <https://issues.apache.org/jira/browse/CASSANDRA-2044>.
- 5 DERBY-5396. <https://issues.apache.org/jira/browse/DERBY-5396>.
- 6 eGit. <http://www.eclipse.org/egit/>.
- 7 PMD. <https://pmd.github.io>.
- 8 The usage guide of WALA. <http://wala.sourceforge.net/wiki/index.php/UserGuide:MappingToSourceCode>.
- 9 WALA. <http://wala.sf.net>.
- 10 R. Abreu, P. Zoetewej, et al. An evaluation of similarity coefficients for software fault localization. In *Proc. PRDC*, pages 39–46, 2006.
- 11 R. Abreu, P. Zoetewej, R. Golsteijn, and A. J. Van Gemund. A practical evaluation of spectrum-based fault localization. *Journal of Systems and Software*, 82(11):1780–1792, 2009.
- 12 R. Abreu, P. Zoetewej, and A. J. Van Gemund. Spectrum-based multiple fault localization. In *Proc. ASE*, pages 88–99, 2009.
- 13 R. Agrawal and R. Srikant. Mining sequential patterns. In *Proc. ICDE*, pages 3–14, 1995.
- 14 G. Ammons, R. Bodík, and J. Larus. Mining specifications. In *Proc. 29th POPL*, pages 4–16, 2002.
- 15 G. Antoniol, K. Ayari, M. Di Penta, F. Khomh, and Y.-G. Guéhéneuc. Is it a bug or an enhancement?: a text-based approach to classify change requests. In *Proc. CASCON*, page 23, 2008.
- 16 A. Bachmann, C. Bird, F. Rahman, P. Devanbu, and A. Bernstein. The missing links: bugs and bug-fix commits. In *Proc. ESEC/FSE*, pages 97–106, 2010.
- 17 J. Benesty, J. Chen, Y. Huang, and I. Cohen. *Pearson correlation coefficient*, pages 1–4. 2009.
- 18 A. Berglund, S. Boag, D. Chamberlin, M. F. Fernandez, M. Kay, J. Robie, and J. Siméon. Xml path language (xpath). *World Wide Web Consortium (W3C)*, 2003.
- 19 M. Brünink and D. S. Rosenblum. Mining performance specifications. In *Proc. ESEC/FSE*, pages 39–49, 2016.
- 20 J. Campos, R. Abreu, G. Fraser, and M. d’Amorim. Entropy-based test generation for improved fault localization. In *Proc. ASE*, pages 257–267, 2013.
- 21 N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer. SMOTE: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, 16:321–357, 2002.
- 22 H. Cheng, D. Lo, Y. Zhou, X. Wang, and X. Yan. Identifying bug signatures using discriminative graph mining. In *Proc. ISSTA*, pages 141–152, 2009.
- 23 J. C. Corbett, M. B. Dwyer, J. Hatcliff, S. Laubach, C. S. Păsăreanu, R. Bby, and H. Zheng. Bandera: Extracting finite-state models from Java source code. In *Proc. 22nd ICSE*, pages 439–448, 2000.
- 24 B. Dagenais and L. J. Hendren. Enabling static analysis for partial Java programs. In *Proc. OOPSLA*, pages 313–328, 2008.
- 25 N. K. C. M. S. H. Dallmeier, Valentin and A. Zeller. Generating test cases for specification mining. In *Proc. ISSTA*, pages 85–96, 2010.
- 26 V. Debroy and W. E. Wong. Insights on fault interference for programs with multiple bugs. In *Proc. ISSRE*, pages 165–174, 2009.
- 27 D. Di Nucci, F. Palomba, G. De Rosa, G. Bavota, R. Oliveto, and A. De Lucia. A developer centered bug prediction model. *IEEE Transactions on Software Engineering*, 44(1):5–24, 2017.
- 28 D. Di Nucci, F. Palomba, D. A. Tamburri, A. Serebrenik, and A. De Lucia. Detecting code smells using machine learning techniques: are we there yet? In *Proc. SANER*, pages 612–621, 2018.
- 29 W. Dickinson, D. Leon, and A. Podgurski. Finding failures by cluster analysis of execution profiles. In *Proc. ICSE*, pages 339–348, 2001.
- 30 N. DiGiuseppe and J. A. Jones. On the influence of multiple faults on coverage-based fault localization. In *Proc. ISSTA*, pages 210–220, 2011.
- 31 N. DiGiuseppe and J. A. Jones. Fault density, fault types, and spectra-based fault localization. *Empirical Software Engineering*, 20(4):928–967, 2015.
- 32 H. Do, S. Elbaum, and G. Rothermel. Supporting controlled experimentation with testing techniques: An infrastructure and its potential impact. *Empirical Software Engineering*, 10(4):405–435, 2005.
- 33 K. El Emam, W. Melo, and J. C. Machado. The prediction of faulty classes using object-oriented design metrics. *Journal of Systems and Software*, 56(1):63–75, 2001.
- 34 D. Engler, D. Chen, and A. Chou. Bugs as inconsistent behavior: A general approach to inferring errors in systems code. In *Proc. 18th SOSP*, pages 57–72, 2001.
- 35 M. Ernst, J. Perkins, P. Guo, S. McCamant, C. Pacheco, M. Tschantz, and C. Xiao. The Daikon system for dynamic detection of likely invariants. *Science of Computer Programming*, 69(1-3):35–45, 2007.
- 36 T. Fawcett. An introduction to roc analysis. *Pattern recognition letters*, 27(8):861–874, 2006.
- 37 P. A. Flach and S. Wu. Repairing concavities in roc curves. In *In Proc. IJCAI*, pages 702–707, 2005.
- 38 B. Fluri, M. Wursch, M. Pinzger, and H. C. Gall. Change distilling: Tree differencing for fine-grained source code change extraction. *Transactions on Software Engineering*, 33(11):725–743, 2007.
- 39 E. Frank. *Pruning Decision Trees and Lists*. PhD thesis, Department of Computer Science, University of Waikato, 2000.
- 40 Y. Freund and R. E. Schapire. Experiments with a new boosting algorithm. In *Proc. ICML*, pages 148–156, San Francisco, 1996.

- 41 R. Gao and W. E. Wong. Mseer-an advanced technique for locating multiple bugs in parallel. *IEEE Transactions on Software Engineering*, pages 1–1, 2018.
- 42 B. Ghotra, S. McIntosh, and A. E. Hassan. Revisiting the impact of classification techniques on the performance of defect prediction models. In *Proc. ICSE*, pages 789–800, 2015.
- 43 Z. Gu, E. T. Barr, D. J. Hamilton, and Z. Su. Has the bug really been fixed? In *Proc. 32nd ICSE*, pages 55–64, 2010.
- 44 P. J. Guo, T. Zimmermann, N. Nagappan, and B. Murphy. Characterizing and predicting which bugs get fixed: an empirical study of microsoft windows. In *Proc. ICSE*, pages 495–504, 2010.
- 45 M. A. Hall. Correlation-based feature selection for machine learning. 1999.
- 46 T. Hall, S. Beecham, D. Bowes, D. Gray, and S. Counsell. A systematic literature review on fault prediction performance in software engineering. *IEEE Transactions on Software Engineering*, 38(6):1276–1304, 2012.
- 47 D. Hao, T. Xie, L. Zhang, X. Wang, J. Sun, and H. Mei. Test input reduction for result inspection to facilitate fault localization. *Automated software engineering*, 17(1):5–31, 2010.
- 48 A. E. Hassan. Predicting faults using the complexity of code changes. In *Proc. ICSE*, pages 78–88, 2009.
- 49 J. H. Hayes, A. Dekhtyar, and J. Osborne. Improving requirements tracing via information retrieval. In *Proc. RE*, pages 138–147, 2003.
- 50 H. He and E. A. Garcia. Learning from imbalanced data. *IEEE Transactions on Knowledge and Data Engineering*, 21(9):1263–1284, 2009.
- 51 K. Herzig, S. Just, and A. Zeller. It’s not a bug, it’s a feature: how misclassification impacts bug prediction. In *Proc. ICSE*, pages 392–401, 2013.
- 52 A. Hindle, D. M. German, and R. Holt. What do large commits tell us?: a taxonomical study of large commits. In *Proc. 4th MSR*, pages 99–108, 2008.
- 53 D. Hovemeyer and W. Pugh. Finding bugs is easy. In *Proc. OOPSLA*, pages 132–136, 2004.
- 54 H.-Y. Hsu, J. A. Jones, and A. Orso. Rapid: Identifying bug signatures to support debugging activities. In *Proc. ASE*, pages 439–442, 2008.
- 55 M. Hutchins, H. Foster, T. Goradia, and T. Ostrand. Experiments of the effectiveness of dataflow-and controlflow-based test adequacy criteria. In *Proc. ICSE*, pages 191–200, 1994.
- 56 G. H. John and P. Langley. Estimating continuous distributions in bayesian classifiers. In *Proc. UAI*, pages 338–345, 1995.
- 57 B. Johnson, Y. Song, E. Murphy-Hill, and R. Bowdidge. Why don’t software developers use static analysis tools to find bugs? In *Proc. ICSE*, pages 672–681, 2013.
- 58 J. A. Jones, M. J. Harrold, and J. Stasko. Visualization of test information to assist fault localization. In *Proc. ICSE*, pages 467–477, 2002.
- 59 R. Just, C. Parnin, I. Drosos, and M. D. Ernst. Comparing developer-provided to user-provided tests for fault localization and automated program repair. In *Proc. ISSTA*, pages 287–297, 2018.
- 60 D. Kim, Y. Tao, S. Kim, and A. Zeller. Where should we fix this bug? a two-phase recommendation model. *IEEE Transactions on Software Engineering*, 39(11):1597–1610, 2013.
- 61 S. Kim, T. Zimmermann, E. J. Whitehead Jr, and A. Zeller. Predicting faults from cached history. In *Proc. 29th ICSE*, pages 489–498, 2007.
- 62 R. Kohavi. The power of decision tables. In *Proc. ECML*, pages 174–189, 1995.
- 63 H. W. Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955.
- 64 M. Last, M. Friedman, and A. Kandel. The data mining approach to automated software testing. In *Proc. KDD*, pages 388–396, 2003.
- 65 T. Le, X. Le, D. Lo, and I. Beschastnikh. Synergizing specification miners through model fissions and fusions. In *Proc. ASE*, pages 115–125, 2015.
- 66 T.-D. B. Le, R. J. Oentaryo, and D. Lo. Information retrieval and spectrum based bug localization: better together. In *Proc. ESEC/FSE*, pages 579–590, 2015.
- 67 S. le Cessie and J. van Houwelingen. Ridge estimators in logistic regression. *Applied Statistics*, 41(1):191–201, 1992.
- 68 C. Lemieux, D. Park, and I. Beschastnikh. General LTL specification mining. In *Proc. ASE*, pages 81–92, 2015.
- 69 J. Li and M. D. Ernst. CBCD: Cloned buggy code detector. In *Proc. ICSE*, pages 310–320, 2012.
- 70 Z. Li and Y. Zhou. PR-Miner: Automatically extracting implicit programming rules and detecting violations in large software code. In *Proc. ESEC/FSE*, pages 306–315, 2005.
- 71 X.-Y. Liu, J. Wu, and Z.-H. Zhou. Exploratory undersampling for class-imbalance learning. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 39(2):539–550, 2009.
- 72 J. B. Lovins. *Development of a stemming algorithm*. MIT Information Processing Group, Electronic Systems Laboratory Cambridge, 1968.
- 73 L. Lucia, D. Lo, L. Jiang, F. Thung, and A. Budi. Extended comprehensive study of association measures for fault localization. *Journal of Software: Evolution and Process*, 26(2):172–219, 2014.
- 74 S. K. Lukins, N. A. Kraft, and L. H. Etzkorn. Bug localization using latent dirichlet allocation. *Information and Software Technology*, 52(9):972–990, 2010.
- 75 X. Mao, Y. Lei, Z. Dai, Y. Qi, and C. Wang. Slice-based statistical fault localization. *Journal of Systems and Software*, 89:51–62, 2014.
- 76 S. Maoz and J. O. Ringert. GR(1) synthesis for LTL specification patterns. In *Proc. ESEC/FSE*, pages 96–106, 2015.
- 77 A. Marcus, D. Poshyvanyk, and R. Ferenc. Using the conceptual cohesion of classes for fault prediction in object-oriented systems. *IEEE Transactions on Software Engineering*, 34(2):287–300, 2008.
- 78 M. Martinez and M. Monperrus. Mining software repair models for reasoning on the search space of automated program fixing. *Empirical Software Engineering*, 20(1):176–205, 2013.
- 79 D. Mease, A. J. Wyner, and A. Buja. Boosted classification trees and class probability/quantile estimation. *Journal of Machine Learning Research*, 8:409–439, 2007.

- 80 H. MEI and L. Zhang. Can big data bring a breakthrough for software automation? *Science China Information Sciences*, 61:056101, 2018.
- 81 T. Menzies and J. S. Di Stefano. More success and failure factors in software reuse. *IEEE Transactions on Software Engineering*, 29(5):474, 2003.
- 82 T. Mikolov, K. Chen, G. Corrado, and J. Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- 83 A. Mishne, S. Shoham, and E. Yahav. Typestate-based semantic code search over partial programs. In *Proc. OOPSLA*, pages 997–1016, 2012.
- 84 N. Nagappan, T. Ball, and A. Zeller. Mining metrics to predict component failures. In *Proc. ICSE*, pages 452–461, 2006.
- 85 L. Naish, H. J. Lee, and K. Ramamohanarao. A model for spectra-based software diagnosis. *ACM Transactions on Software Engineering and Methodology*, 20(3):11, 2011.
- 86 D. Newman, A. Asuncion, P. Smyth, and M. Welling. Distributed algorithms for topic models. *Journal of Machine Learning Research*, 10:1801–1828, 2009.
- 87 A. T. Nguyen, T. T. Nguyen, J. Al-Kofahi, H. V. Nguyen, and T. N. Nguyen. A topic-based approach for narrowing the search space of buggy files from a bug report. In *Proc. ASE*, pages 263–272, 2011.
- 88 H. V. Nguyen, H. A. Nguyen, A. T. Nguyen, and T. N. Nguyen. Mining interprocedural, data-oriented usage patterns in JavaScript web applications. In *Proc. ICSE*, pages 791–802, 2014.
- 89 T. T. Nguyen, H. A. Nguyen, N. H. Pham, J. M. Al-Kofahi, and T. N. Nguyen. Graph-based mining of multiple object usage patterns. In *Proc. ESEC/FSE*, pages 383–392, 2009.
- 90 K. J. Ottenstein and L. M. Ottenstein. The program dependence graph in a software development environment. *ACM Sigplan Notices*, 19(5):177–184, 1984.
- 91 R. Pandita, X. Xiao, H. Zhong, T. Xie, S. Oney, and A. Paradkar. Inferring method specifications from natural language API descriptions. In *Proc. 34th ICSE*, pages 815–825, 2012.
- 92 S. Pearson, J. Campos, R. Just, G. Fraser, R. Abreu, M. D. Ernst, D. Pang, and B. Keller. Evaluating and improving fault localization. In *Proc. ICSE*, pages 609–620, 2017.
- 93 A. Perez, R. Abreu, and M. d’Amorim. Prevalence of single-fault fixes and its impact on fault localization. In *Proc. ICST*, pages 12–22, 2017.
- 94 A. Perez, R. Abreu, and A. van Deursen. A test-suite diagnosability metric for spectrum-based fault localization approaches. In *Proc. ICSE*, pages 654–664, 2017.
- 95 N. H. Pham, T. T. Nguyen, H. A. Nguyen, X. Wang, A. T. Nguyen, and T. N. Nguyen. Detecting recurring and similar software vulnerabilities. In *Proc. ICSE*, volume 2, pages 227–230, 2010.
- 96 J. C. Platt. Fast training of support vector machines using sequential minimal optimization. *Advances in kernel methods*, pages 185–208, 1999.
- 97 A. Podgurski, D. Leon, P. Francis, W. Masri, M. Minch, J. Sun, and B. Wang. Automated support for classifying software failure reports. In *Proc. 25th ICSE*, pages 465–475, 2003.
- 98 M. Pradel and T. R. Gross. Leveraging test generation and specification mining for automated bug detection without false positives. In *Proc. ICSE*, pages 288–298, 2012.
- 99 Y. Qi, X. Mao, Y. Lei, Z. Dai, and C. Wang. The strength of random search on automated program repair. In *Proc. 36th ICSE*, pages 254–265, 2014.
- 100 E. Rahm and H. H. Do. Data cleaning: Problems and current approaches. *IEEE Data Engineering Bulletin*, 23(4):3–13, 2000.
- 101 F. Rahman, D. Posnett, A. Hindle, E. Barr, and P. Devanbu. Bugcache for inspections: hit or miss? In *Proc. ESEC/FSE*, pages 322–331, 2011.
- 102 M. Ramanathan, A. Grama, and S. Jagannathan. Path-sensitive inference of function precedence protocols. In *Proc. 29th ICSE*, pages 240–250, 2007.
- 103 S. Rao and A. Kak. Retrieval from software libraries for bug localization: a comparative study of generic and composite text models. In *Proc. MSR*, pages 43–52, 2011.
- 104 M. P. Robillard, E. Bodden, D. Kawrykow, M. Mezini, and T. Ratchford. Automated API property inference techniques. *IEEE Transactions on Software Engineering*, 39(5):613–637, 2013.
- 105 M. J. Rochkind. The source code control system. *IEEE transactions on Software Engineering*, (4):364–370, 1975.
- 106 R. K. Saha, M. Lease, S. Khurshid, and D. E. Perry. Improving bug localization using structured information retrieval. In *Proc. ASE*, pages 345–355, 2013.
- 107 A. Saied, O. Benomar, H. Abdeen, and H. Sahraoui. Mining multi-level API usage patterns. In *Proc. SANER*, pages 23–32, 2015.
- 108 F. Sarro, S. Di Martino, F. Ferrucci, and C. Gravino. A further analysis on the use of genetic algorithm to configure support vector machines for inter-release fault prediction. In *Proc. SAC*, pages 1215–1220, 2012.
- 109 E. Shapiro. *Algorithmic program debugging*. PhD thesis, Yale University, 1983.
- 110 B. Sisman and A. C. Kak. Incorporating version histories in information retrieval based bug localization. In *Proc. MSR*, pages 50–59, 2012.
- 111 C. Sun and S.-C. Khoo. Mining succinct predicated bug signatures. In *Proc. ESEC/FSE*, pages 576–586, 2013.
- 112 Y. Sun, M. S. Kamel, A. K. Wong, and Y. Wang. Cost-sensitive boosting for classification of imbalanced data. *Pattern Recognition*, 40(12):3358–3378, 2007.
- 113 J. A. Suykens and J. Vandewalle. Least squares support vector machine classifiers. *Neural processing letters*, 9(3):293–300, 1999.
- 114 C. Tantithamthavorn, S. McIntosh, A. E. Hassan, and K. Matsumoto. Automated parameter optimization of classification techniques for defect prediction models. In *Proc. ICSE*, pages 321–332, 2016.
- 115 C. Tantithamthavorn, S. Mcintosh, A. E. Hassan, and K. Matsumoto. The impact of automated parameter optimization on defect

- prediction models. *IEEE Transactions on Software Engineering*, PP(99):1–1, 2018.
- 116 C. Thornton, F. Hutter, H. H. Hoos, and K. Leyton-Brown. Auto-WEKA: Combined selection and hyperparameter optimization of classification algorithms. In *Proc. KDD*, pages 847–855, 2013.
- 117 Y. Tian, J. Lawall, and D. Lo. Identifying linux bug fixing patches. In *Proc. 34th ICSE*, pages 386–396, 2012.
- 118 M. Tufano, F. Palomba, G. Bavota, M. Di Penta, R. Oliveto, A. De Lucia, and D. Poshyvanyk. There and back again: Can you compile that snapshot? *Journal of Software: Evolution and Process*, 2016.
- 119 Q. Wang, C. Parnin, and A. Orso. Evaluating the usefulness of IR-based fault localization techniques. In *Proc. ISSTA*, pages 1–11, 2015.
- 120 S. Wang, T. Liu, and L. Tan. Automatically learning semantic features for defect prediction. In *Proc. ICSE*, pages 297–308, 2016.
- 121 S. Wang and D. Lo. Amalgam+: Composing rich information sources for accurate bug localization. *Journal of Software: Evolution and Process*, 28(10):921–942, 2016.
- 122 S. Wang, D. Lo, and J. Lawall. Compositional vector space models for improved bug localization. In *Proc. ICSME*, pages 171–180, 2014.
- 123 Y. Wang, N. Meng, and H. Zhong. An empirical study of multi-entity changes in real bug fixes. In *Proc. ICSME*, page to appear, 2018.
- 124 A. Wasylkowski, A. Zeller, and C. Lindig. Detecting object usage anomalies. In *Proc. ESEC/FSE*, pages 35–44, 2007.
- 125 W. Weimer, T. Nguyen, C. Le Goues, and S. Forrest. Automatically finding patches using genetic programming. In *Proc. ICSE*, pages 364–374, 2009.
- 126 G. M. Weiss. Mining with rarity: a unifying framework. *ACM SIGKDD Explorations Newsletter*, 6(1):7–19, 2004.
- 127 C. C. Williams and J. K. Hollingsworth. Automatic mining of source code repositories to improve bug finding techniques. *IEEE Transactions on Software Engineering*, 31(6):466–480, 2005.
- 128 C. Wong, Y. Xiong, H. Zhang, D. Hao, L. Zhang, and H. Mei. Boosting bug-report-oriented fault localization with segmentation and stack-trace analysis. In *Proc. ICSME*, pages 181–190, 2014.
- 129 W. E. Wong, V. Debroy, and D. Xu. Towards better fault localization: A crosstab-based statistical approach. *IEEE Transactions on Systems, Man, and Cybernetics*, 42(3):378–396, 2012.
- 130 W. E. Wong, R. Gao, Y. Li, R. Abreu, and F. Wotawa. A survey on software fault localization. *IEEE Transactions on Software Engineering*, 42(8):707–740, 2016.
- 131 W. E. Wong and Y. Qi. Bp neural network-based effective fault localization. *International Journal of Software Engineering and Knowledge Engineering*, 19(04):573–597, 2009.
- 132 R. Wu, H. Zhang, S. Kim, and S.-C. Cheung. Relink: recovering links between bugs and changes. In *Proc. ESEC/FSE*, pages 15–25, 2011.
- 133 Q. Yang and X. Wu. 10 challenging problems in data mining research. *International Journal of Information Technology & Decision Making*, 5(04):597–604, 2006.
- 134 H. Zhong and N. Meng. Towards reusing hints from past fixes -an exploratory study on thousands of real samples. *Empirical Software Engineering*, 23(5):2521–2549, 2018.
- 135 H. Zhong and Z. Su. An empirical study on real bug fixes. In *Proc. ICSE*, pages 913–923, 2015.
- 136 H. Zhong and X. Wang. Boosting complete-code tools for partial program. In *Proc. ASE*, pages 671–681, 2017.
- 137 H. Zhong, L. Zhang, T. Xie, and H. Mei. Inferring resource specifications from natural language API documentation. In *Proc. ASE*, pages 307–318, 2009.
- 138 J. Zhou, H. Zhang, and D. Lo. Where should the bugs be fixed? more accurate information retrieval-based bug localization based on bug reports. In *Proc. ICSE*, pages 14–24, 2012.
- 139 Z. Zuo, S.-C. Khoo, and C. Sun. Efficient predicated bug signature mining via hierarchical instrumentation. In *Proc. ISSTA*, pages 215–224, 2014.