

Extended Comprehensive Study of Association Measures for Fault Localization

Lucia*, David Lo*, Lingxiao Jiang*, Ferdian Thung, and Aditya Budi

School of Information Systems, Singapore Management University, Singapore

SUMMARY

Spectrum-based fault localization is a promising approach to automatically locate root causes of failures quickly. Two well-known spectrum-based fault localization techniques, Tarantula and Ochiai, employ measures to quantify how likely a program element is a root cause of failures based on profiles of correct and failed program executions. These measures are conceptually similar to *association measures* that have been proposed in statistics, data mining, and other research areas and have been utilized to quantify the relationship strengths between two variables of interest (e.g., the use of a medicine and the cure rate of a disease). In this paper, we view fault localization as a measurement of the relationship strength between the execution of program elements and program failures. We investigate the effectiveness of 40 association measures from the literature on locating bugs and look for existing measures that have similar or better performance than Tarantula and Ochiai. Our empirical evaluations involve programs containing single bug and also multiple bugs. We find that there is no best single measure for all cases. Klosgen and Ochiai outperform other measures for localizing single bug programs. While for localizing multiple-bug programs, Added Value, Odds Ratio, Yule's Q, and Coverage outperform other measures.

KEY WORDS: Association Measures, Fault Localization, Program Spectra

*Correspondence to: 80 Stamford Road, School of Information Systems, Singapore Management University, Singapore 178902. E-mail: lucia.2009@smu.edu.sg, davidlo@smu.edu.sg, lxjiang@smu.edu.sg

1. INTRODUCTION

Software debugging is a difficult and expensive activity to perform. US National Institute of Standards and Technology has reported that software bugs cost US economy 59.5 billion dollars annually [1], and testing and debugging activities account for 30 to 90 percent of labor expended for a project [2]. When a program failure occurs, the execution trace may be long and contain failure-irrelevant information. Often, full execution traces are not even available for debugging if programs fail in the field. Locating the root cause of the failure, which may be far away from the failure point, is non-trivial.

Many approaches have been proposed to help in automating debugging activities, especially in localizing root causes of failures (i.e., *faults*) in programs [3–11]. One family of approaches is spectrum-based fault localization [3,4,12–14], where program traces or abstractions of traces (called *program spectra*) are used to represent program behaviors and the spectra of correct executions and failed executions are compared to identify potentially faulty program elements (e.g., statements, basic blocks, functions, and components). The comparison often employs statistical analysis, and program elements that are observed more often in failed executions than in correct executions (or statistically correlate with failures) may be identified and presented to developers for further inspection.

Spectrum-based fault localization techniques are promising as they are lightweight and have good accuracy. Among spectrum-based fault localization techniques, two well-known techniques are Tarantula [4, 15] and Ochiai [13, 16, 17]. Both approaches compute a suspiciousness score for each program element based on the execution frequencies of the element in correct and failed executions, and rank all program elements according to their scores. Thus, higher suspiciousness scores indicate more likely faulty program elements, and the computation of the scores in effect answers the following question:

What is the strength of association between the executions of a particular program element with failures?

Based on this question, a general solution for fault localization can naturally emerge based on the proliferation of association measures in the literature: measure the strength of association of a program element's executions with failures, and the stronger the association is, the more likely the program element is a fault.

Besides Tarantula and Ochiai which have been used for fault localization, various association measures in the data mining and statistics communities, such as odds ratio [18], Yule's Q [19], and Yule's Y [20] have been proposed to measure the strength of the association of two variables. For example, one might be interested in the association between an application of a particular medical treatment with a recovery from an illness, or in the association between an execution of a business strategy with the revenue change. There are several studies that evaluate the effectiveness of some similarity coefficients, e.g., Jaccard [13, 21], Sorensen-Dice [21], Anderberg [21], Simple Matching [21], Rogers and Tanimoto [21], and Ochiai II [21]. However, there are rich varieties of association measures other than those coefficients that have not been studied for fault localization.

This paper aims to fill this gap by investigating the effectiveness of 40 popular association measures for the purpose of fault localization and comparing them with Tarantula and Ochiai. In particular, we are interested in answering the following research questions (RQs):

- RQ 1.** Are vanilla or off-the-shelf association measures accurate enough in localizing faults?
- RQ 2.** Which off-the-shelf association measures are more accurate to localize faults?
- RQ 3.** What is the relative performance of the off-the-shelf association measures as compared to well-known suspiciousness measures for fault localization, Tarantula and Ochiai in particular?
- RQ 4.** Is the accuracy of the off-the-shelf association measures, Tarantula, and Ochiai different for programs written in C as compared to programs written in Java?
- RQ 5.** What is the effectiveness of the off-the-shelf association measures, Tarantula, and Ochiai in localizing different types of bugs?
- RQ 6.** What is the accuracies of the off-the-shelf association measures, Tarantula, and Ochiai in localizing faulty programs containing multiple bugs?

To answer the above questions, we investigate and compare the accuracies of Tarantula, Ochiai, and the additional 40 association measures on programs from the Siemens test suite [22], and three larger programs from Software Infrastructure Repository (SIR) [23]. The latter includes Space which is written in C, and NanoXml and XmlSecurity, which are written in Java. The programs come along with seeded bugs, test oracles to decide between failures and non-failures, and test cases. We compute various accuracy metrics to evaluate the effectiveness of the association measures in localizing faults to answer the above research questions. We show that a few association measures have better performance than Ochiai and many are better than Tarantula. We also split the programs into those written in C and those written in Java, and analyze the accuracies of the association measures, Tarantula, and Ochiai separately for each of the two sets of programs. Furthermore, we characterize the bugs in the programs and group them into several categories. We then evaluate the effectiveness of the various association measures, Tarantula, and Ochiai, to localize different categories of bugs and also evaluate their accuracies in localizing bugs for programs.

The contributions of this work are as follows:

1. We comprehensively investigate the effectiveness of 40 association measures for fault localization.
2. We highlight a few promising association measures which can outperform Tarantula and Ochiai and those that are comparable with the two well-known spectrum-based fault localization approaches.
3. We provide a partial order of association measures in terms of their accuracies for fault localization.
4. We characterize the effectiveness of the association measures, Ochiai, and Tarantula on programs written in different programming languages.
5. We analyze different kinds of bugs and investigate the effectiveness of the 40 association measures, Ochiai, and Tarantula in localizing each of these categories of bugs.

6. We investigate the accuracies of the 40 association measures, Ochiai, and Tarantula in localizing faulty programs containing multiple bugs.

The structure of this paper is as follows. Section 2 discusses related work. Section 3 discusses the fundamental concepts of spectrum-based fault localization and association measures. Section 4 discusses the particular association measures considered in the paper. Section 5 describes our empirical evaluation and comparison of the association measures. Finally, we conclude our work in Section 6.

2. RELATED WORK

In this section, we describe closely related studies on fault localization and association measures. The survey here is by no means a complete list of all related studies.

2.1. Fault Localization

Recently, there are many studies on fault localization and automated debugging. There are different ways to categorize these studies. Based on the data analyzed by the approaches, fault localization techniques can be classified into *spectrum-based* and *model-based*.

2.1.1. Spectrum-based fault localization. Spectrum-based fault localization techniques often use program spectra, which are program traces or abstractions of program traces that represent program runtime behaviors in certain ways, to correlate program elements (e.g., statements, basic blocks, functions, and components) with program failures (often with the help of statistical analysis).

Many spectrum-based fault localization techniques [4, 5, 12, 24, 25] take as inputs two sets of spectra, one for successful executions and the other for failed executions, and report candidate locations where root causes of program failures (i.e., faults) may reside. Given a failed program spectrum and a set of correct spectra, Renieris and Reiss present a fault localization tool WHITHER [5] that compares the failed execution to the nearest correct execution and reports the most suspicious locations in the program. Liblit et al. propose a technique to search for predicates

whose true evaluation correlates with failures [24]. Chao et al. extend the work by incorporating information on the outcomes of multiple predicate evaluations in a program run in their tool called SOBER [12]. Santelices et al. combine several program spectra to better localize bugs in programs. [25]. All of these techniques need to compare spectra of failed executions with those of successful executions in some way. Evaluating the effectiveness of various association measures can complement all of these techniques by helping to locate the most failure-relevant program elements quickly and improving their performance.

Artzi et al. evaluate the effectiveness of several test generation techniques in generating enough test cases for localizing faults in web applications with the help of Ochiai [26]. Artzi et al. extend their work by proposing a tool named Apollo that can generate test cases to expose failures for web applications and also localize bugs that cause the failures [27]. Their fault localization technique uses Tarantula and a technique that keep information about which program statements are potentially responsible to produce a particular part of an output (e.g., a table in an HTML document). This approach is possible for web applications but may not be applicable to other applications as the output is often not decomposable into parts (it could be a single number) and the number of program statements that are potentially responsible to produce an output are often large. Bandyopadhyay and Ghosh study how the properties of faults affect the effectiveness of Tarantula [28]. Three properties namely accessibility, original state failure condition, and impact are investigated. To answer RQ5, we also investigate the effectiveness of various fault localization techniques on different kinds of bugs. However, we consider a different categorization of faults – our categorization is based on the empirical study performed by Kim et al. on bug fixes [29]. Jiang et al. study the effectiveness of test case prioritization for fault localization using different types of prioritization criteria [30].

Other spectrum-based techniques [8, 31–33] only use failed executions as the input and systematically alter the program structure or program runtime states to locate faults. Zhang et al. [31] search for faulty program predicates by switching the states of program predicates at runtime. Sterling and Olsson use the concept of program chipping [32] to automatically remove parts of a program so that the part that contributes to the failure may become more apparent. While

their tool, ChipperJ, works on syntax trees for Java programs, Gupta et al. [8] work on program dependency graphs and use the intersection of forward and backward program slices to reduce the sizes of failure-relevant code for further inspection. Jeffrey et al. use a value profile based approach to rank program statements according to their likelihood of being faulty [33]. These fault localization techniques do not compare the spectra of failed executions with those of successful executions, and association measures are generally not applicable to them. List of programs that have been used by those past studies in fault localization is shown in Table I.

2.1.2. Model-based fault localization. Compared with spectrum-based techniques, model-based debugging techniques [9, 34–37] are often more accurate, but heavyweight since they are based on more expensive logic reasoning over formal models of programs. Many static and dynamic analysis techniques [6, 38, 39] can be classified as model-based debugging as well. Abreu et al. propose a framework called BARINEL that combines spectrum-based fault localization and model-based debugging to localize single and multiple bugs in programs, and found that the approach is more accurate and heavyweight than spectrum-based fault localization [17]. Although few model-based techniques have employed the concept of failure association, incorporating association measures investigated in this study into program models can be a future direction to improve the performance of model-based debugging techniques.

In our study, we focus on comparisons with two well-known spectrum-based fault localization techniques, namely Tarantula [4, 15] and Ochiai [13, 16, 17]. We evaluate 40 association measures and find promising ones for fault localization.

2.2. Studies On Association Measures

There have been a number of studies proposed in the statistics and data mining community on measures of association between variables since the early 20th century. These include measures such as Yule's Q and Yule's Y [19, 20]. Other measures, such as odds ratio [18], are also commonly considered and utilized in various domains, such as medical [41] and social science [42]. In the data mining community, Agrawal and Srikant have proposed association rule mining which aims to infer

Dataset	LOC	Papers	Dataset	LOC	Papers
SumPowers	27	[35,36]	Power	763	[24]
BinSearch	29	[35,36]	Compress	1590	[24]
BubbleSort	29	[35,36]	Bh	2053	[24]
Hamming	48	[35,36]	Webchess v.0.9.0	2226	[26,27]
Adder	49	[35,36]	Tetris	2403	[40]
Permutation	54	[35,36]	Schoolmate v.1.5.4	4263	[26,27]
Binomia	80	[35,36]	Gzip	5680	[17,30]
Polynom	105	[35,36]	Space	6218	[4,15,17,21]
Tcas	141	[4,5,12,16,17,28,30,33,35–37]	NanoXML	7646	[25,40]
Schedule	292	[4,5,8,12,13,16,17,25,28,30,33]	Phpsysinfo v.2.5.3	7745	[26,27]
Schedule2	301	[4,5,8,12,13,16,17,25,28,30,33]	Li	7761	[24]
Treeadd	385	[24]	Grep	10068	[30,31]
Perimeter	395	[24]	Flex	10459	[30,31]
Print_token2	399	[4,5,8,12,13,16,17,25,28,30,33]	Sed	14427	[17,30]
Tot_info	440	[4,5,12,13,16,17,25,28,30,33]	XMLsecurity	16800	[25]
Print_token	478	[4,5,8,12,13,16,17,25,28,30,33]	Bc-1.06	17042	[31]
Replace	512	[4,5,8,12,13,16,17,28,30,33]	Tar-1.13.25	27137	[31]
Emad	557	[24]	Go	29315	[24]
Rest	617	[24]	Ijpeg	31371	[24]
Bisort	707	[24]	Make	35545	[31]
Health	725	[24]	JABA	37966	[25]
Faqforge v.1.3.2	734	[26,27]			

Table I. Subject Programs Used in Past Fault Localization Studies

associations from two itemsets in a transaction dataset in the early 90s [43]. In that work the metrics of support and confidence for measuring the strength of an association are proposed. Various other metrics, such as interest and collective strength, are proposed later. We describe these measures in detail in Section 4 .

Tan et al. investigate various association measures, compare their properties, and outline the benefits and limitations of each from a computational point of view [44]. The measures are revisited by Geng and Hamilton by including measures for aggregated data summaries [45]. In this paper, we extend their work in the specific domain of fault localization by comparing 40 association measures based on their ability to assign high suspiciousness scores for buggy program elements and low scores for non-buggy ones.

Some of the association measures that we evaluate in this paper, have been studied for fault localizations, e.g., Jaccard [13, 21], Sorensen-Dice [21], Anderberg [21], Simple Matching [21], Rogers and Tanimoto [21], and Ochiai II [21]. In this paper, we revisit the effectiveness of these measures as their effectiveness on various kinds of bugs and programs containing multiple bugs has not been extensively evaluated.

3. CONCEPTS & DEFINITIONS

In this section we formally introduce the problem of spectrum-based fault localization as the computation of association strengths between the executions of various program elements and failures. Also, we describe the concept of dichotomy matrix that is used in the calculation of these association strengths.

3.1. Spectrum-Based Fault Localization

This problem starts with a faulty program, a set of test cases, and a test oracle. The set of test cases are run over the faulty program and observations of how the program runs on each of the test cases are recorded as program spectra. A program spectrum represents certain characteristics of an execution of a program, providing a behavior signature of the execution [46]. The signature of a behavior could be a set of counters, each of which indicates the number of times each program element (e.g., statement, basic block, path, etc.) is executed in one execution [47]. The counters could also simply be 0-1 flags that indicate whether an element is executed. A test oracle is available to label whether a particular output or execution of a test case is correct or wrong. Wrong executions are classified as program failures. The task of a fault localization tool is to find the program elements that are responsible for the failures (i.e., the faults or the root causes) based on the program spectra of both correct and wrong executions.

There have been various spectra proposed in the literature [13, 47]. Different spectra may have different effects on effectiveness of fault localization. The block-hit spectra are at a suitable profiling granularity because all code in the same basic block has the same execution pattern and there is no

Block ID	Program Elements	T15	T16	T17	T18
1	int count; int n; Ele *proc; List *src_queue, *dest_queue; if (prio >= MAXPRIO) /*maxprio=3*/	●	●	●	●
2	{return;}		●	●	●
3	src_queue = prio_queue[prio]; dest_queue = prio_queue[prio+1]; count = src_queue->mem_count; if (count > 1) /* Bug */ /* supposed : count>0 */ {	●	●	●	●
4	n = (int) (count*ratio + 1); proc = find_nth(src_queue, n); if (proc) {		●	●	
5	src_queue = del_ele(src_queue, proc); proc->priority = prio; dest_queue = append_ele(dest_queue, proc); } }		●	●	
Status of Test Case Execution :		Pass	Pass	Pass	Fail

Figure 1. Example of block-hit program spectra

need to instrument individual instructions in a granularity finer than blocks, and because it has been shown in the literature that the instrumentation costs to obtain such spectra are relatively low and it can be used for effective fault localization [13, 16, 17, 21, 47]. The granularity has a balance between reducing instrumentation costs and having sufficient bug-revealing capabilities.

In this paper, we use *block-hit program spectra*, each of which consists of a set of flags to indicate whether each basic block is executed or not in each test case. An example of block-hit program spectra is shown in Figure 1. The first column contains identifiers of basic blocks. The second column contains the statements in the corresponding basic blocks. The other columns indicate whether each basic block is executed in test cases T15, T16, T17, and T18 along with the information whether each of the test cases passes or fails. In this example, ● denotes that a basic block is executed by a test case and an empty cell denotes that the block is not executed by the test case. In the code snippet, a bug lies in the condition of the `if` statement in Block 3, causing Blocks 4–5 to be skipped when the variable `count` is 1. Note that in test cases T16–T18, execution of Block 2 that contains return statement is followed by execution of Block 3. Normally, Block 3 should not be executed, but since this snippet code is being called inside a loop. Thus the traces of these test cases contain execution of Block 2 together with the following blocks.

Symbol	Definition
n	Total number of test cases in the test suite
$n(e)$	Number of test cases that executes a program element e
n_s	Number of test cases that pass
n_f	Number of test cases that fail
$n_s(e)$	Number of test cases that execute e and pass
$n_f(e)$	Number of test cases that execute e and fail

Table II. Some common notations

Based on these spectra, we want to compute the suspiciousness score of each program element following Definition 3.1.

Definition 3.1 (Suspiciousness Score)

Consider a program $P = \{e_1, \dots, e_n\}$ and a set of program spectra $T = T_s \cup T_f$ for P , where P comprises of n elements e_1, \dots, e_n and T comprises of the spectra for correct executions T_s and the spectra for wrong executions T_f . We would like to measure the strength of the association between the executions of each e_i and program failures and assign this strength as the *suspiciousness score* of e_i denoted as $suspiciousness(e)$.

3.1.1. *Tarantula* Jones and Harrold propose Tarantula [4] to rank program elements based on their suspiciousness scores. Intuitively, a program element is more suspicious if it appears in failed executions more frequently than in correct executions. Considering a program P and a test suite T , Table II introduce some common notations which are used in the rest of the paper.

Tarantula's suspiciousness score for a program element e can be computed as follows:

$$suspiciousness(e) = \frac{\frac{n_f(e)}{n_f}}{\frac{n_s(e)}{n_s} + \frac{n_f(e)}{n_f}}$$

Based on block-hit program spectra shown in Figure 1, the suspiciousness score of Block 3 that contains the bug is $\frac{1/1}{3/3+1/1} = 0.5$. Block 1 has the same suspiciousness score. Interestingly, Block 2

receives the highest suspiciousness score: $\frac{1/1}{2/3+1/1} = 0.6$. Following the same calculation, Blocks 4 and 5 are not suspicious. There is no failure that executes these blocks and hence Tarantula returns a suspiciousness score of 0. Assuming developers inspect program elements one by one from the most suspicious to the least, the bug in Block 3 can be found after at most 3 blocks have been inspected.

3.1.2. *Ochiai* Abreu et al. [16] propose *Ochiai* which assigns the suspiciousness score of a program element as follows:

$$\text{suspiciousness}(e) = \frac{n_f(e)}{\sqrt{n_f(n_f(e) + n_s(e))}} = \frac{n_f(e)}{\sqrt{n_f n(e)}}$$

Similar to Tarantula, *Ochiai* considers an element more suspicious if it occurs more frequently in failed executions than in correct executions (the $\sqrt{\frac{n_f(e)}{n(e)}}$ part). Using the same example shown in Figure 1, Blocks 1 and 3 receive a suspiciousness score: $1/\sqrt{3 * (1 + 0)} = 0.14$. Similar to Tarantula, *Ochiai* also returns Block 2 as the most suspicious block: $1/\sqrt{4 * (1 + 0)} = 0.20$, while the remaining blocks are assigned suspiciousness scores of 0. As the case with Tarantula, by employing *Ochiai*, the bug can be found after 3 blocks have been inspected. In this study, we are interested to investigate other association measures that can possibly localize the bug earlier.

3.2. Dichotomous Association

A common characteristics of the association measures evaluated in this paper is that they are all defined based on dichotomy matrices. The following are the necessary definitions.

Definition 3.2 (Dichotomy)

A *dichotomous outcome* is an outcome whose values could be split into two categories, e.g., wrong or correct, executed or skipped, married or unmarried, etc. A *dichotomous variable* is a variable having a dichotomous outcome. A *dichotomy matrix* is a 2×2 matrix that tries to associate two dichotomous variables in the form of a 2×2 contingency table which records the bivariate frequency distribution of the two variables.

An example of a dichotomy matrix $D(A, B)$ relating variables A and B is shown in Table III. The value c_{00} corresponds to the number of observations in which the value of variable A equals to A_0 and the value of variable B equals to B_0 . The values of the other three entries in the dichotomy matrix are similarly defined.

	$A = A_0$	$A = A_1$
$B = B_0$	c_{00}	c_{01}
$B = B_1$	c_{10}	c_{11}

Table III. An example of a dichotomy matrix. We refer to it as $D(A, B)$.

Based on the concept of dichotomy matrix, we introduce dichotomous association in Definition 3.3.

Definition 3.3 (Dichotomous Association)

A *dichotomous association* is a special form of bivariate association [42] which measures the strength of association between two dichotomous variables, e.g., application of a medical treatment and recovery from the disease, job satisfaction and productivity, and program element execution and program failure. The formulae for calculating dichotomous associations depend on the four entries in dichotomy matrices.

Given a dichotomy matrix relating two variables, two questions are often asked:

1. Is there a (dichotomous) association between the two variables?
2. How strong is the association between the two variables?

A common way to answer these two questions is to define a formula, referred to as an *association measure*, to calculate a score based on the four entries in a dichotomy matrix and consider the association exists (or is strong) if the score is beyond a particular threshold. We define association measure in Definition 3.4.

Definition 3.4 (Association Measure)

An association measure M of two variable A and B is a mathematical function of the four entries of a dichotomy matrix $D(A, B)$, and is denoted as $M(A, B, D(A, B))$ or simply $M(A, B, D)$ if it is clear from the context.

In fault localization, we could produce a dichotomy matrix that relates the executions of a program element with program failures. Consider a program element e . For each test case, a trace is generated when a subject program is executed on the test case. Some traces execute e , others do not. Some traces correspond to failures, others correspond to correct executions. After the test cases are run, a dichotomy matrix as shown in Table IV is produced for every program element e .

	e Executed	e Not Executed
Test Passed	$n_s(e)$	$n_s(\bar{e})$
Test Failed	$n_f(e)$	$n_f(\bar{e})$

Table IV. Dichotomy matrix for fault localization.

The notation \bar{e} means e is not executed, and other notations are defined in Table II. Thus, $n_s(\bar{e})$ is the number of test cases that do not execute e and pass, and $n_f(\bar{e})$ is the number of test cases that do not execute e and fail.

Considering variables E and F to represent a program element being executed and a program failure occurs respectively, we are interested to compute $M(E, F, D_e(E, F))$ (i.e., the association between the execution of e and a failure), where $D_e(E, F)$ represents the dichotomy matrix of the two variables E and F for a program element e . The formulae of the 40 association measures are given in Section 4.

4. ASSOCIATION MEASURES

In this section, we first describe how a dichotomy matrix can be constructed. Next, we present how the 40 association measures can be computed from a dichotomy matrix. Finally, we give an example how we sort program elements for inspection using the association measures.

4.1. Constructing a Dichotomy Matrix

Dichotomy matrix construction requires programs instrumentation that could support collection of program execution traces and a test oracle. In this paper, we instrument all buggy programs in our dataset in basic block level. We manually instrument for C programs and use Cobertura[†] to instrument Java programs. In order to construct dichotomy matrix, we collect the program execution traces which contain information about which basic blocks are executed by each test case. The test oracle would give us information on whether a test case fails or passes. Next, for each basic block, we construct a dichotomy matrix by counting the number of times the basic block is executed by the passing test cases which is denoted as $n_s(e)$, number of times the basic block is not executed by the passing test cases which is denoted as $n_s(\bar{e})$, number of times the basic block is executed by failing test cases which is denoted as $n_f(e)$, and number of times the basic block is not executed by failing test cases which is denoted as $n_f(\bar{e})$. The information in this dichotomy matrix is then used as an input to the association measures to calculate how likely the corresponding basic block contains a bug.

4.2. Association Measures

The 40 association measures that we consider are as follows: ϕ -coefficient [42], odds ratio [18], Yule's Q [19], Yule's Y [20], Kappa [48], J-Measure [49], gini index [50], support [43], confidence [43], Clark and Boswell's Laplace accuracy [51], conviction [52], interest [52], cosine [44], Piatetsky-Shapiro's Leverage [53], certainty factor [54], added value [44], collective strength [55], Jaccard [56], Klogen [57], information gain [58, 59], Coverage [45,

[†]<http://cobertura.sourceforge.net/>

Name	Formula
ϕ -Coefficient (M_1)	$\frac{P(A,B) - P(A)P(B)}{\sqrt{P(A)P(B)(1-P(A))(1-P(B))}}$
Odds ratio (M_2)	$\frac{P(A,B)P(\bar{A},\bar{B})}{P(A,\bar{B})P(\bar{A},B)}$
Yule's Q (M_3)	$\frac{P(A,B)P(\bar{A},\bar{B}) - P(A,\bar{B})P(\bar{A},B)}{P(A,B)P(\bar{A},\bar{B}) + P(A,\bar{B})P(\bar{A},B)} = \frac{\alpha - 1}{\alpha + 1}$
Yule's Y (M_4)	$\frac{\sqrt{P(A,B)P(\bar{A},\bar{B})} - \sqrt{P(A,\bar{B})P(\bar{A},B)}}{\sqrt{P(A,B)P(\bar{A},\bar{B})} + \sqrt{P(A,\bar{B})P(\bar{A},B)}} = \frac{\sqrt{\alpha - 1}}{\sqrt{\alpha + 1}}$
Kappa (M_5)	$\frac{P(A,B) + P(\bar{A},\bar{B}) - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}$
J-Measure (M_6)	$\max(P(A, B) \log\left(\frac{P(B A)}{P(B)}\right) + P(\bar{A}\bar{B}) \log\left(\frac{P(\bar{B} \bar{A})}{P(\bar{B})}\right),$ $P(A, B) \log\left(\frac{P(A B)}{P(A)}\right) + P(\bar{A}\bar{B}) \log\left(\frac{P(\bar{A} \bar{B})}{P(\bar{A})}\right))$
Gini Index (M_7)	$\max(P(A)[P(B A)^2 + P(\bar{B} \bar{A})^2] + P(\bar{A})[P(B \bar{A})^2 +$ $P(\bar{B} \bar{A})^2] - P(B)^2 - P(\bar{B})^2, P(B)[P(A B)^2 + P(\bar{A} B)^2] +$ $P(\bar{B})[P(A \bar{B})^2 + P(\bar{A} \bar{B})^2] - P(A)^2 - P(\bar{A})^2)$
Support (M_8)	$P(A, B)$
Confidence (M_9)	$\max(P(B A), P(A B))$
Laplace (M_{10})	$\max\left(\frac{P(A,B)+1}{P(A)+2}, \frac{P(A,B)+1}{P(B)+2}\right)$
Conviction (M_{11})	$\max\left(\frac{P(A)P(\bar{B})}{P(\bar{A}B)}, \frac{P(B)P(\bar{A})}{P(\bar{B}A)}\right)$
Interest (M_{12})	$\frac{P(A,B)}{P(A)P(B)}$
Piatetsky-Shapiro's (M_{13})	$P(A, B) - P(A)P(B)$
Certainty Factor (M_{14})	$\max\left(\frac{P(B A) - P(B)}{1 - P(B)}, \frac{P(A B) - P(A)}{1 - P(A)}\right)$
Added Value (M_{15})	$\max(P(B A) - P(B), P(A B) - P(A))$
Collective Strength (M_{16})	$\frac{P(A,B) + P(\bar{A},\bar{B})}{P(A)P(B) + P(\bar{A})P(\bar{B})} \times \frac{1 - P(A)P(B) - P(\bar{A})P(\bar{B})}{1 - P(A,B) - P(\bar{A}\bar{B})}$
Jaccard (M_{17})	$\frac{P(A,B)}{P(A) + P(B) - P(A,B)}$
Klosgen (M_{18})	$\sqrt{P(A, B)} \max(P(B A) - P(B), P(A B) - P(A))$
Information Gain (M_{19})	$(-P(B) \log P(B) - P(\bar{B}) \log P(\bar{B})) - (P(A) \times$ $(-P(B A) \log P(B A)) - P(\bar{B} A) \log P(\bar{B} A) - P(\bar{A}) \times$ $(-P(B \bar{A}) \log P(B \bar{A})) - P(\bar{B} \bar{A}) \log P(\bar{B} \bar{A}))$

Table V. Definitions of association measures [Part I]. A and B are the two variables in the dichotomy matrix. $P(A)$ and $P(B)$ correspond to the probabilities of A and B respectively. Other notations follow standard notations in Probability

Name	Formula
Coverage (M_{20})	$P(A)$
Accuracy (M_{21})	$P(A, B) + P(\bar{A}, \bar{B})$
Leverage (M_{22})	$P(B A) - P(A)P(B)$
Relative Risk (M_{23})	$P(B A)/P(B \bar{A})$
Interestingness Weighting Dependency (M_{24})	$((\frac{P(A,B)}{P(A)P(B)})^k - 1)P(A, B)^m$ where k,m are coefficients of dependency and generality, respectively, weighting the relative importance of the two factors.
Goodman and Kruskal (M_{25})	$\frac{\sum_i \max_j P(A_i, B_j) + \sum_j \max_i P(A_i, B_j) - \max_i P(A_i) - \max_j P(B_j)}{2 - \max_i P(A_i) - \max_j P(B_j)}$
Normalized Mutual Information (M_{26})	$\sum_i \sum_j P(A_i, B_j) \log_2 \frac{P(A_i, B_j)}{P(A_i)P(B_j)} / \{-\sum_i P(A_i) \log_2 P(A_i)\}$
One-Way Support (M_{27})	$P(B A) \log_2 \frac{P(A,B)}{P(A)P(B)}$
Two-Way Support (M_{28})	$P(A, B) \log_2 \frac{P(A,B)}{P(A)P(B)}$
Two-Way Support Variation (M_{29})	$P(A, B) \log_2 \frac{P(A,B)}{P(A)P(B)} + P(A, \bar{B}) \log_2 \frac{P(A, \bar{B})}{P(A)P(\bar{B})} + P(\bar{A}, B) \log_2 \frac{P(\bar{A}, B)}{P(\bar{A})P(B)} + P(\bar{A}, \bar{B}) \log_2 \frac{P(\bar{A}, \bar{B})}{P(\bar{A})P(\bar{B})}$
Loevinger (M_{30})	$1 - \frac{P(A)P(\bar{B})}{P(A, \bar{B})}$
Sebag-Schoenauer (M_{31})	$\frac{P(A, B)}{P(A, \bar{B})}$
Least Contradiction (M_{32})	$\frac{P(A, B) - P(A, \bar{B})}{P(B)}$
Odd Multiplier (M_{33})	$\frac{P(A, B)P(\bar{B})}{P(B)P(A, \bar{B})}$
Example and Counterexample Rate (M_{34})	$1 - \frac{P(A, \bar{B})}{P(A, B)}$
Zhang (M_{35})	$\frac{P(A, B) - P(A)P(B)}{\max(P(A, B)P(\bar{B}), P(B)P(A, \bar{B}))}$
Sorensen-Dice (M_{36})	$\frac{2P(A, B)}{2P(A, B) + P(\bar{A}, B) + P(A, \bar{B})}$
Anderberg (M_{37})	$\frac{P(A, B)}{P(A, B) + 2(P(\bar{A}, B) + P(A, \bar{B}))}$
Simple-Matching (M_{38})	$P(A, B) + P(\bar{A}, \bar{B})$
Rogers and Tanimoto (M_{39})	$\frac{P(A, B) + P(\bar{A}, \bar{B})}{P(A, B) + P(\bar{A}, \bar{B}) + 2(P(\bar{A}, B) + P(A, \bar{B}))}$
Ochiai II (M_{40})	$\frac{P(A, B) + P(\bar{A}, \bar{B})}{\sqrt{(P(A, B) + P(\bar{A}, \bar{B})) (P(A, B) + P(A, \bar{B})) (P(\bar{A}, \bar{B}) + P(\bar{A}, B)) (P(\bar{A}, \bar{B}) + P(A, \bar{B}))}}$

Table VI. Definitions of association measures [Part II].

60], Accuracy [45, 60], Leverage [45, 60], Relative Risk [45, 60], Interestingness Weighting Dependency [45, 60], Goodman and Kruskal [44, 45, 60, 61], Normalized Mutual Information [44, 45, 60], One-Way Support [45, 60], Two-Way Support [45, 60], Two-Way Support Variation [45, 60], Loevinger [45], Sebag-Schoenauer [45], Least Contradiction [45], Odd Multiplier [45], Example

Name	Range	No	Perfect
ϕ -Coefficient (M_1)	$-1 \dots 0 \dots 1$	0	1
Odds ratio (M_2)	$0 \dots 1 \dots \infty$	1	∞
Yule's Q (M_3)	$-1 \dots 0 \dots 1$	0	1
Yule's Y (M_4)	$-1 \dots 0 \dots 1$	0	1
Kappa (M_5)	$-1 \dots 0 \dots 1$	0	1
J-Measure (M_6)	$0 \dots 1$	0	1
Gini index (M_7)	$0 \dots 1$	0	1
Support (M_8)	$0 \dots 1$	0	1
Confidence (M_9)	$0 \dots 1$	0	1
Laplace (M_{10})	$0 \dots 1$	0	1
Conviction (M_{11})	$0.5 \dots 1 \dots \infty$	1	∞
Interest (M_{12})	$0 \dots 1 \dots \infty$	1	∞
Piatetsky-Shapiro's (M_{13})	$-0.25 \dots 0 \dots 0.25$	0	0.25
Certainty factor (M_{14})	$-1 \dots 0 \dots 1$	0	1
Added Value (M_{15})	$-0.5 \dots 0 \dots 1$	0	1
Collective strength (M_{16})	$0 \dots 1 \dots \infty$	1	∞
Jaccard (M_{17})	$0 \dots 1$	0	1
Klogsen (M_{18})	$(\frac{2}{\sqrt{3}} - 1)^{1/2} [2 - \sqrt{3} - \frac{1}{\sqrt{3}}] \dots 0 \dots \frac{2}{3\sqrt{3}}$	0	$\frac{2}{3\sqrt{3}}$
Information Gain (M_{19})	$0 \dots 1$	0	1

Table VII. Value ranges of the association measures [Part I]. The third and fourth columns give the values corresponding to no association and perfect association respectively. Values of measures with ranges in the format of $a \dots b$ (e.g., J-Measure ($0 \dots 1$), etc.) indicate positive associations with failures. Values of measures with ranges in the format of $a \dots b \dots c$ (e.g., ϕ -coefficient ($-1 \dots 0 \dots 1$), etc.) indicate positive associations with failures (if they are between b and c), or negative associations with failures (if they are between a and b). Values closer to a imply stronger associations with passing executions.

and Counterexample Rate [45], Zhang [45], Sorensen-Dice [62, 63], Anderberg [64], Simple Matching [65], Rogers and Tanimoto [66], and Ochiai II [67].

The mathematical formulae for calculating these 40 association measures are given in Tables V and VI. The formulae are defined in terms of probabilities, instead of frequencies, but we can substitute frequencies recorded in dichotomous matrices for probabilities during actual calculations. The ranges of values that these association measures can take are given in Tables VII and VIII.

Name	Range	No	Perfect
Coverage (M_{20})	$0 \dots 1$	0	1
Accuracy (M_{21})	$0 \dots 1$	0	1
Leverage (M_{22})	$-1 \dots 0 \dots 1$	0	1
Relative Risk (M_{23})	$0 \dots 1 \dots \infty$	1	∞
Interestingness Weighting Dependency (M_{24})	$-1 \dots 0 \dots 1$	0	1
Goodman and Kruskal (M_{25})	$0 \dots 1$	0	1
Normalized Mutual Information (M_{26})	$0 \dots 1$	0	1
One-Way Support (M_{27})	$-1 \dots 0 \dots 1$	0	1
Two-Way Support (M_{28})	$-1 \dots 0 \dots 1$	0	1
Two-Way Support Variation (M_{29})	$0 \dots 1$	0	1
Loevinger (M_{30})	$-1 \dots 0 \dots 1$	0	1
Sebag-Schoenauer (M_{31})	$0 \dots 1 \dots \infty$	1	∞
Least Contradiction (M_{32})	$0 \dots 1$	0	1
Odd Multiplier (M_{33})	$0 \dots 1 \dots \infty$	1	∞
Example and Counterexample Rate (M_{34})	$-\infty \dots 0 \dots 1$	0	1
Zhang (M_{35})	$-1 \dots 0 \dots 1$	0	1
Sorensen-Dice (M_{36})	$0 \dots 1$	0	1
Anderberg (M_{37})	$0 \dots 1$	0	1
Simple-Matching (M_{38})	$0 \dots 1$	0	1
Rogers and Tanimoto (M_{39})	$0 \dots 1$	0	1
Ochiai II (M_{40})	$0 \dots 1$	0	1

Table VIII. Value ranges of the association measures [Part II].

4.3. From Association to Suspiciousness

From Section 3, the suspiciousness score for a program element e with dichotomy matrix D_e is defined as the strength of the association between the executions of e with failures (i.e., $M(E, F, D_e)$). M refers to one of the 40 association measures presented in Section 4.2.

Next, we illustrate how an association measure could be used to rank program elements for inspection. Using the example in Figure 1, we observe that there are 5 blocks: Blocks 1, 2, 3, 4, and 5. The bug resides at the `if` statement in block 3. The suspiciousness score of block 3 is determined by the association strength between the executions of block 3 and failures. For example, by using one of the association measures, e.g., Coverage, blocks 3 and 1 receive the highest suspiciousness score, i.e., 1, followed by block 2 whose suspiciousness score is 0.75. The suspiciousness scores of

blocks 4 and 5 are 0.5. This particular measure can rank the block containing the bug such that no other block receives a score higher than it. However, Tarantula and Ochiai are unable to do so—see Sections 3.1.1 and 3.1.2.

5. EMPIRICAL EVALUATION

In this section we first describe our datasets, followed by our evaluation metrics, and experimental results.

5.1. Datasets

We analyze different programs from Siemens Test Suite [68]. Siemens programs are injected with realistic bugs and often analyzed for fault localization studies [4, 5, 8, 12, 13, 16, 17, 25, 28, 30, 33, 35–37]. We also analyze other three real programs from Software-artifact Infrastructure Repository (SIR) [22] namely: Space [4, 15, 17, 21], NanoXML [25], and XML-Security [25]. Space is written in C, while NanoXML and XML-Security are written in Java. The average lines of code for various versions of Space, NanoXML and XML-Security are 6,218, 4,223, and 21,275 respectively. Siemens test suite was originally used for research in test coverage adequacy and was developed by Siemens Corporation Research. We use the variant provided at www.cc.gatech.edu/aristotle/Tools/subjects/. Each program contains many different versions where each version has one bug. These bugs comprise a wide array of realistic bugs. The Siemens Test Suite comes with 7 programs: `print_tokens`, `print_tokens2`, `replace`, `schedule`, `schedule2`, `tcas`, and `tot_info`. The total number of buggy versions are 132, as shown in Table IX. We manually instrumented the buggy versions at basic block level. Since our instrumentation cannot reach the bugs that reside in variable declarations, we exclude versions that contain this type of bugs, i.e., versions 6, 10, 19, 21 of `tot_info` dataset, version 12 of `replace` dataset, and versions 13, 14, 15, 36, 38 of `tcas` dataset. We exclude versions 4 and 6 of `print_token` because they are identical with the original version. We also exclude version 9 of `schedule2` because there is no test case that results in a failure. Thus, in total, we use 119 buggy versions from the Siemens test suite.

Dataset	LOC	Language	Num. of Faulty Version	Num. of Test Cases
print.token	478	C	5	4130
print.token2	399	C	10	4115
replace	512	C	31	5542
schedule	292	C	9	2650
schedule2	301	C	9	2710
tcas	141	C	36	1608
tot.info	440	C	19	1051
space	6,218	C	35	13,585
NanoXML v1	3,497	Java	6	214
NanoXML v2	4,007	Java	7	214
NanoXML v3	4,608	Java	9	216
NanoXML v5	4,782	Java	8	216
XML security v1	21,613	Java	6	92
XML security v2	22,318	Java	6	94
XML security v3	19,895	Java	4	84

Table IX. Dataset descriptions

Space is an interpreter for Array Definition Language (ADL) used by European Space Agency. We analyze all 35 faulty versions of Space downloaded from SIR. NanoXML is a utility for parsing XML. SIR contains 5 versions of NanoXML. We exclude NanoXML_v4 because there is no buggy version. For each version, SIR provides a few bugs. In total there are 32 buggy versions for NanoXML_v1, NanoXML_v2, NanoXML_v3, and NanoXML_v5, and we analyze 30 of them. We exclude two buggy versions because there is no test case that results in a failure. XML-Security is a Java library that supports digital signature and encryption. SIR contains 3 versions of XML-Security. Again for each version, SIR provides a few bugs. In total there are 52 buggy versions for *XMLSec.v1*, *XMLSec.v2*, and *XMLSec.v3*. We exclude 16 buggy versions because there is no test case that results in a failure. Thus, the total number of buggy versions that we analyze for the 3 programs are 81.

Table IX provides the details on the number of lines, the programming language in which the program is written, the number of faulty versions, and the number of test cases, of each subject programs (programs in the Siemens test suite and the other 3 medium-size programs).

5.2. Evaluation Metrics

We use 40 association measures, Tarantula, and Ochiai to rank program elements based on their suspiciousness. Developers could then use this list to investigate the more suspicious program elements first. We assume that developers would be able to identify the buggy program element when they inspect it.

We consider two commonly used evaluation metrics: average percentage of code inspected to find all bugs, and proportion of bugs found when a given proportion of the code is inspected. We describe these two metrics in the following paragraphs.

5.2.1. Percentage of Code Inspected. We evaluate the performance of the measures by the number of elements that are ranked as high or higher than the program element containing the fault/the bug. For a suspiciousness score to be effective, buggy program elements should have a relatively larger value of suspiciousness scores than the non-buggy elements.

When a buggy program element has the same suspiciousness score with several other elements, the largest rank of the elements that has this suspiciousness score is used as the rank of the buggy element. For example, consider the case where the two highest suspiciousness scores are 0.92 and 0.91, two elements have suspiciousness scores of 0.92, and 3 elements have suspiciousness scores of 0.91. If a buggy element is given the suspiciousness score of 0.91, then the rank of this buggy element is 5, instead of 3. Since we do not know how the programmer will traverse elements that have the same suspiciousness score, we use the worst case scenario where the programmer inspects all elements having the same score. In the situation when a buggy version contains multiple bugs or one single bug that involves multiple program elements, we use the largest rank among the buggy elements as it is the worst case rank that represents scenarios when programmers would like to find all the buggy elements by using the given list of most suspicious program elements.

Ranking program elements by using suspiciousness score is useful to evaluate the accuracy of a fault-localization approach. However, it may not be representative enough to know how programmers really locate the bugs with their expertise, as this ranking approach does not provide

the context of the bugs to help programmers in understanding the root cause of the bugs, as the study by Parnin and Orso has shown [40]. It is an interesting future study to augment suspiciousness ranking with some additional information to more effectively guide programmers to locate all buggy program elements.

Suspiciousness measures that rank buggy elements first are more effective than those that rank them last. We then use this rank to compute the percentage of program elements that need to be inspected to find the buggy elements by the formula:

$$\frac{\textit{largest rank among the buggy elements}}{\textit{total elements}}$$

In our experiment, we thus choose to use basic block as the granularity of the elements and the above percentage is applied as the first accuracy criterion of the association measures. The accuracy of a particular measure to localize bug in all buggy versions of our datasets is then evaluated by calculating the overall mean of all percentages of code inspected for the measure which is the average of the percentages of code inspected of the measure for all buggy versions in our datasets. The smaller the overall mean, the more accurate the measure in localizing bug. However, the overall mean might not necessarily reflect that the measure would localize bugs with the same accuracy for all buggy versions. A measure could have a good accuracy in localizing one bug, but might not have a good accuracy in localizing other bugs. Thus, we also calculate the overall standard deviation of a measure to evaluate the variance of percentages of code inspected of the measure for all buggy versions. The smaller the overall standard deviation, the better the overall mean reflects the accuracy of the measure because of less variation of percentage of code inspected for all buggy versions.

5.2.2. Proportion of Bugs Localized. Next, we calculate the proportion of bugs that could be localized assuming that the developers are only willing to investigate a given proportion of code. To compute this measure we vary the proportion of code that the developers are willing to inspect and for each proportion we compute the proportion of bugs that can be localized.

5.3. Evaluation Results

We describe the accuracy of the association measures in comparison with Ochiai and Tarantula in the following paragraphs.

5.3.1. Percentage of Code Inspected. The overall means and standard deviations of percentage of code inspected of the 40 association measures along with those of Ochiai and Tarantula for all subject programs are shown in Table X. The smallest mean is 22.42% which is achieved by Klosgen. Ochiai and Tarantula achieve 22.66% and 24.66% respectively. Other measures that have similar accuracy to Klosgen (M_{18}) and Ochiai (i.e., having a mean in the range of 22.5% and 23.5%) are Added Value (M_{15}), ϕ -Coefficient (M_1), Normalized Mutual Information (M_{26}), and Two-way Support (M_{28}). There are 13 measures that have similar accuracy as Tarantula (i.e., having a mean between 23.5% to 25.5%), i.e., Interestingness Weighting Dependency (M_{24}), J-Measure (M_6), Information Gain (M_{19}), Two-Way Support Variation (M_{29}), Example and Counterexample Rate (M_{34}), Kappa (M_5), Confidence (M_9), Odd Multiplier (M_{33}), Sebag (M_{31}), Interest (M_{12}), Zhang (M_{35}), One-Way Support (M_{27}), Jaccard (M_{17}), Sorensen-Dice (M_{36}), and Anderberg (M_{37}). We notice that most of the association measures have similar standard deviations of around 20% (ranging from 19% to 27%), except for Goodman and Kruskal (M_{25}) which has the highest standard deviation among all measures i.e., 37%.

In our evaluation dataset, we have eight C programs and seven Java programs, the accuracy of an association measure in localizing bugs for different programs might be different. We are interested to evaluate the accuracy of the association measures in localizing bugs for each program. Thus, we also calculate the mean and standard deviation of the measures for buggy versions in each program. Tables XI & XII show the detail of the accuracy values (i.e., mean and standard deviation) of each measure for each C program, while Tables XIII & XIV show the detail for each of Java program. Based on the results, the measures have different accuracy for different program and different programming languages. The measures generally require more than 10% of percentage of code to be inspected in order to localize bugs for three C programs (i.e., schedule2, tcas, and tot_info) and

Association Measures	Mean	StdDev	Association Measures	Mean	StdDev
Klogsen(M_{18})	22.42%	24.58%	Sorensen-Dice(M_{36})	24.86%	24.88%
Ochiai	22.66%	23.81%	Anderberg(M_{37})	24.86%	24.90%
Collective Strength(M_{16})	22.97%	24.20%	Ochiai II(M_{40})	25.78%	26.94%
Added Value(M_{15})	23.03%	24.93%	Gini Index(M_7)	25.83%	26.84%
ϕ -Coefficient(M_1)	23.12%	24.54%	Leverage(M_{22})	27.12%	26.16%
Normalized Mutual Information(M_{26})	23.22%	24.04%	Least Contradiction(M_{32})	31.71%	27.70%
Two-Way Support(M_{28})	23.37%	24.93%	Rogers and Tanimoto(M_{39})	31.76%	27.73%
Interestingness Weighting Dependency(M_{24})	23.74%	24.94%	Accuracy(M_{21})	31.98%	27.69%
J-Measure(M_6)	24.23%	25.80%	Simple-Matching(M_{38})	31.98%	27.69%
Information Gain(M_{19})	24.37%	25.82%	Odds Ratio(M_2)	40.07%	19.80%
Two-Way Support Variation(M_{29})	24.37%	25.82%	Yule's Q (M_3)	40.09%	19.82%
Example and Counterexample Rate(M_{34})	24.49%	24.59%	Conviction(M_{11})	40.09%	19.70%
Kappa(M_5)	24.51%	24.76%	Certainty Factor(M_{14})	40.10%	19.70%
Confidence(M_9)	24.52%	24.74%	Yule's Y (M_4)	40.10%	19.81%
Sebag(M_{31})	24.63%	24.67%	Goodman and Kruskal(M_{25})	43.04%	37.55%
Odd Multiplier(M_{33})	24.63%	24.68%	Relative Risk(M_{23})	43.30%	20.71%
Zhang(M_{35})	24.64%	24.66%	Support(M_8)	43.64%	20.32%
Interest(M_{12})	24.64%	24.68%	Laplace(M_{10})	43.64%	20.30%
One-Way Support(M_{27})	24.65%	24.76%	Coverage(M_{20})	45.09%	23.07%
Tarantula	24.66%	24.66%	Piatetsky-Shapiro's(M_{13})	56.02%	24.49%
Jaccard(M_{17})	24.72%	24.94%	Loevinger(M_{30})	56.31%	24.79%

Table X. Overall mean and standard deviation (in parentheses) of accuracy values (smaller the better)

the Java programs. For the other two C programs (i.e., schedule and space), more than half of the measures could guide a debugger to localize bugs in these programs by inspecting less than 10% of code.

We also perform statistical tests for each pair of measures including Tarantula and Ochiai using Wilcoxon signed rank test [69] at 0.05 statistical significance threshold to see if some measures are statistically significantly better than others. We use this statistical test because it does not assume that the data follows normal distribution. We visualize the statistical significance relationship as a partial order in Figure 2. A link from A to B in the partial order denotes that A is statistically significantly better than B. The relationships expressed in the partial order are transitive: if A outperforms B, and B outperforms C, then A outperforms C too. To reduce the number of links in the partial order, we omit links that could be captured by this transitivity property.

Association Measures	Programs							
	print_token	print_token2	schedule	schedule2	replace	tcas	tot_info	space
ϕ -Coefficient	14%(19%)	13%(17%)	<u>9%(15%)</u>	47%(24%)	11%(14%)	50%(27%)	19%(15%)	<u>*3%(7%)</u>
OddsRatio	38%(10%)	44%(2%)	49%(28%)	53%(7%)	33%(13%)	62%(13%)	49%(20%)	18%(6%)
Yule's Q	38%(10%)	44%(3%)	48%(28%)	53%(7%)	33%(13%)	62%(13%)	49%(20%)	18%(6%)
Yule's Y	38%(9%)	44%(3%)	48%(28%)	53%(7%)	33%(13%)	62%(13%)	49%(20%)	18%(6%)
Kappa	22%(22%)	16%(20%)	<u>7%(6%)</u>	51%(24%)	14%(16%)	51%(27%)	21%(16%)	<u>4%(8%)</u>
J-Measure	13%(22%)	11%(13%)	24%(19%)	49%(23%)	11%(15%)	52%(30%)	18%(15%)	<u>5%(10%)</u>
Gini Index	16%(28%)	16%(22%)	12%(12%)	58%(26%)	11%(15%)	53%(30%)	24%(19%)	<u>6%(12%)</u>
Support	38%(10%)	45%(3%)	50%(28%)	56%(8%)	33%(12%)	65%(13%)	51%(20%)	20%(8%)
Confidence	23%(22%)	13%(18%)	<u>7%(6%)</u>	51%(24%)	14%(16%)	51%(27%)	23%(16%)	<u>5%(8%)</u>
Laplace	38%(10%)	45%(3%)	50%(28%)	56%(8%)	34%(12%)	65%(13%)	50%(20%)	20%(8%)
Conviction	38%(10%)	44%(3%)	48%(28%)	53%(7%)	33%(12%)	62%(13%)	49%(20%)	18%(6%)
Interest	23%(22%)	17%(20%)	<u>7%(6%)</u>	51%(24%)	14%(16%)	51%(27%)	22%(16%)	<u>5%(8%)</u>
Piatetsky-Shapiro's	61%(32%)	77%(16%)	71%(29%)	68%(28%)	60%(27%)	65%(25%)	63%(28%)	47%(13%)
Certainty Factor	38%(10%)	44%(3%)	48%(28%)	53%(7%)	33%(13%)	62%(13%)	49%(20%)	18%(6%)
Added Value	14%(23%)	14%(19%)	<u>7%(10%)</u>	50%(25%)	10%(14%)	50%(27%)	18%(15%)	<u>4%(8%)</u>
Collective Strength	13%(15%)	13%(16%)	11%(20%)	44%(23%)	11%(13%)	50%(27%)	18%(15%)	<u>*3%(7%)</u>
Jaccard	21%(21%)	12%(18%)	<u>7%(6%)</u>	50%(24%)	13%(15%)	51%(27%)	21%(15%)	<u>*3%(7%)</u>
Klogsen	12%(19%)	11%(15%)	<u>8%(12%)</u>	48%(24%)	<u>*9%(12%)</u>	50%(28%)	*17%(13%)	<u>4%(8%)</u>
Information Gain	13%(22%)	11%(14%)	14%(19%)	49%(23%)	11%(15%)	52%(30%)	19%(16%)	<u>5%(10%)</u>
Coverage	58%(26%)	63%(20%)	69%(30%)	*29%(20%)	50%(22%)	*45%(23%)	51%(28%)	36%(19%)
Accuracy	33%(24%)	28%(22%)	<u>10%(5%)</u>	60%(24%)	26%(21%)	57%(29%)	37%(29%)	<u>9%(16%)</u>
Leverage	24%(23%)	20%(22%)	<u>7%(4%)</u>	55%(26%)	16%(18%)	53%(28%)	28%(23%)	<u>8%(12%)</u>
Relative Risk	38%(10%)	45%(3%)	50%(28%)	56%(8%)	34%(13%)	66%(12%)	50%(20%)	18%(6%)
Int. Weighting Dependency	17%(23%)	15%(20%)	<u>8%(9%)</u>	50%(24%)	12%(15%)	50%(27%)	21%(16%)	<u>4%(8%)</u>
Goodman and Kruskal	20%(36%)	30%(44%)	67%(38%)	83%(9%)	40%(40%)	64%(33%)	66%(41%)	12%(22%)
Normalized Mutual Info.	<u>10%(12%)</u>	<u>9%(11%)</u>	16%(27%)	41%(20%)	11%(15%)	49%(27%)	*17%(13%)	<u>4%(8%)</u>
One-Way Support	23%(22%)	17%(20%)	<u>7%(6%)</u>	51%(24%)	14%(16%)	51%(27%)	22%(16%)	<u>5%(8%)</u>
Two-Way Support	14%(22%)	14%(19%)	<u>8%(10%)</u>	50%(25%)	11%(14%)	50%(27%)	20%(15%)	<u>*3%(8%)</u>
Two-Way Support Variation	13%(22%)	11%(14%)	14%(19%)	49%(23%)	11%(15%)	52%(30%)	19%(16%)	<u>5%(10%)</u>
Loevinger	68%(26%)	78%(21%)	72%(29%)	38%(22%)	64%(23%)	52%(22%)	70%(26%)	59%(24%)
Sebag	23%(22%)	17%(20%)	<u>7%(6%)</u>	51%(24%)	14%(16%)	51%(27%)	22%(16%)	<u>5%(8%)</u>
Least Contradiction	30%(27%)	28%(22%)	<u>10%(5%)</u>	59%(24%)	26%(21%)	57%(29%)	36%(29%)	<u>9%(16%)</u>
Odd Multiplier	23%(22%)	17%(20%)	<u>7%(6%)</u>	51%(24%)	14%(16%)	51%(27%)	22%(16%)	<u>5%(8%)</u>
Example and Counter.	23%(22%)	17%(20%)	<u>7%(6%)</u>	51%(24%)	14%(16%)	50%(27%)	22%(16%)	<u>5%(8%)</u>
Zhang	23%(22%)	17%(20%)	<u>7%(6%)</u>	51%(24%)	14%(16%)	51%(27%)	22%(16%)	<u>5%(8%)</u>

Table XI. Detailed means and standard deviations (in parentheses) of percentages of code inspected to find all bugs in the C programs [Part I]. The star (*) marks the measure(s) with the lowest mean and the underline marks measures that have a mean within 10%.

Association Measures	Programs							
	print_token	print_token2	schedule	schedule2	replace	tcas	tot_info	space
Sorensen-Dice	21%(21%)	16%(20%)	<u>7%(6%)</u>	50%(24%)	14%(15%)	51%(27%)	20%(14%)	<u>*3%(7%)</u>
Anderberg	21%(21%)	16%(20%)	<u>7%(6%)</u>	51%(24%)	14%(15%)	51%(27%)	20%(14%)	<u>*3%(7%)</u>
Simple-Matching	33%(24%)	28%(22%)	<u>10%(5%)</u>	60%(24%)	26%(21%)	57%(29%)	37%(29%)	<u>9%(16%)</u>
Rogers and Tanimoto	30%(27%)	28%(22%)	<u>10%(5%)</u>	60%(24%)	26%(21%)	57%(29%)	36%(29%)	<u>9%(16%)</u>
Ochiai II	19%(27%)	16%(21%)	<u>*6%(4%)</u>	57%(26%)	13%(15%)	52%(29%)	28%(27%)	<u>4%(10%)</u>
Tarantula	23%(22%)	17%(20%)	<u>7%(6%)</u>	51%(24%)	14%(16%)	51%(27%)	22%(16%)	<u>5%(8%)</u>
Ochiai	<u>*8%(8%)</u>	<u>9%(10%)</u>	13%(28%)	42%(22%)	10%(13%)	48%(24%)	*17%(12%)	<u>*3%(8%)</u>

Table XII. Detailed means and standard deviations (in parentheses) of percentages of code inspected to find all bugs in the C programs [Part II]. The star (*) marks the measure(s) with the lowest mean and the underline marks measures that have a mean within 10%.

Association Measures	Programs							
	Nano_v1	Nano_v2	Nano_v3	Nano_v5	XML-sec_v1	XML-sec_v2	XML-sec_v3	
ϕ -Coefficient	*21%(28%)	32%(25%)	*31%(30%)	20%(15%)	29%(0%)	31%(0%)	40%(0%)	
OddsRatio	45%(18%)	43%(18%)	44%(19%)	23%(10%)	29%(0%)	31%(0%)	40%(0%)	
Yule's Q	44%(18%)	43%(18%)	44%(19%)	25%(14%)	29%(0%)	31%(0%)	40%(0%)	
Yule's Y	44%(18%)	43%(18%)	44%(19%)	25%(14%)	29%(0%)	31%(0%)	40%(0%)	
Kappa	*21%(28%)	32%(25%)	*31%(30%)	20%(15%)	29%(0%)	31%(0%)	40%(0%)	
J-Measure	25%(30%)	31%(26%)	*31%(30%)	26%(17%)	29%(0%)	31%(0%)	40%(0%)	
Gini Index	25%(30%)	32%(25%)	*31%(30%)	26%(19%)	29%(0%)	31%(0%)	40%(0%)	
Support	67%(11%)	51%(12%)	55%(10%)	41%(18%)	29%(0%)	31%(0%)	40%(0%)	
Confidence	*21%(28%)	26%(29%)	*31%(30%)	20%(15%)	29%(0%)	31%(0%)	40%(0%)	
Laplace	67%(11%)	51%(12%)	55%(10%)	41%(18%)	29%(0%)	31%(0%)	40%(0%)	
Conviction	44%(18%)	42%(18%)	44%(19%)	25%(14%)	29%(0%)	31%(0%)	40%(0%)	
Interest	*21%(28%)	*26%(29%)	*31%(30%)	20%(15%)	29%(0%)	31%(0%)	40%(0%)	
Piatetsky-Shapiro's	52%(20%)	46%(16%)	36%(13%)	46%(22%)	29%(0%)	31%(0%)	40%(0%)	
Certainty Factor	44%(18%)	42%(18%)	44%(19%)	25%(14%)	29%(0%)	31%(0%)	40%(0%)	
Added Value	*21%(28%)	*26%(29%)	*31%(30%)	20%(15%)	29%(0%)	31%(0%)	40%(0%)	
Collective Strength	*21%(28%)	31%(26%)	*31%(30%)	20%(15%)	29%(0%)	31%(0%)	40%(0%)	
Jaccard	*21%(28%)	53%(9%)	*31%(30%)	19%(15%)	29%(0%)	31%(0%)	40%(0%)	
Klosgen	*21%(28%)	28%(28%)	*31%(30%)	19%(14%)	29%(0%)	31%(0%)	40%(0%)	
Information Gain	25%(30%)	32%(25%)	*31%(30%)	26%(18%)	29%(0%)	31%(0%)	40%(0%)	
Coverage	59%(21%)	47%(19%)	34%(15%)	40%(19%)	29%(0%)	31%(0%)	40%(0%)	
Accuracy	25%(30%)	38%(27%)	32%(30%)	24%(21%)	29%(0%)	31%(0%)	40%(0%)	
Leverage	25%(30%)	28%(28%)	*31%(30%)	24%(17%)	29%(0%)	31%(0%)	40%(0%)	

Table XIII. Detailed means and standard deviations (in parentheses) of percentages of code inspected to find all bugs in the Java programs [Part I]. The star (*) marks the measure(s) with the lowest mean and the underline marks measures that have a mean within 10%.

Association Measures	Programs						
	Nano.v1	Nano.v2	Nano.v3	Nano.v5	XML-sec.v1	XML-sec.v2	XML-sec.v3
Relative Risk	67%(11%)	51%(12%)	55%(10%)	40%(16%)	29%(0%)	31%(0%)	40%(0%)
Int. Weighting Dependency	*21%(28%)	27%(28%)	*31%(30%)	20%(15%)	29%(0%)	31%(0%)	40%(0%)
Goodman and Kruskal	29%(35%)	38%(27%)	40%(30%)	42%(30%)	29%(0%)	31%(0%)	40%(0%)
Normalized Mutual Info.	*21%(28%)	43%(18%)	*31%(30%)	21%(14%)	29%(0%)	31%(0%)	40%(0%)
One-Way Support	*21%(28%)	*26%(29%)	*31%(30%)	21%(18%)	29%(0%)	31%(0%)	40%(0%)
Two-Way Support	*21%(28%)	30%(27%)	*31%(30%)	20%(13%)	29%(0%)	31%(0%)	40%(0%)
Two-Way Support Variation	25%(30%)	32%(25%)	31%(30%)	26%(18%)	29%(0%)	31%(0%)	40%(0%)
Loevinger	63%(20%)	50%(21%)	35%(15%)	45%(21%)	29%(0%)	31%(0%)	40%(0%)
Sebag	*21%(28%)	*26%(29%)	*31%(30%)	20%(15%)	29%(0%)	31%(0%)	40%(0%)
Least Contradiction	25%(30%)	38%(27%)	32%(30%)	24%(21%)	29%(0%)	31%(0%)	40%(0%)
Odd Multiplier	*21%(28%)	*26%(29%)	*31%(30%)	20%(15%)	29%(0%)	31%(0%)	40%(0%)
Example and counter.	*21%(28%)	*26%(29%)	*31%(30%)	20%(15%)	29%(0%)	31%(0%)	40%(0%)
Zhang	*21%(28%)	*26%(29%)	*31%(30%)	20%(15%)	29%(0%)	31%(0%)	40%(0%)
Sorensen-Dice	*21%(28%)	53%(9%)	42%(35%)	19%(15%)	29%(0%)	31%(0%)	40%(0%)
Anderberg	*21%(28%)	53%(9%)	42%(35%)	19%(15%)	29%(0%)	31%(0%)	40%(0%)
Simple-Matching	25%(30%)	38%(27%)	43%(35%)	24%(21%)	29%(0%)	31%(0%)	40%(0%)
Rogers and Tanimoto	25%(30%)	38%(27%)	43%(35%)	24%(21%)	29%(0%)	31%(0%)	40%(0%)
Ochiai II	25%(30%)	32%(25%)	42%(35%)	24%(18%)	29%(0%)	31%(0%)	40%(0%)
Tarantula	*21%(28%)	*26%(29%)	*31%(30%)	20%(15%)	29%(0%)	31%(0%)	40%(0%)
Ochiai	*21%(28%)	53%(10%)	*31%(30%)	*18%(15%)	29%(0%)	31%(0%)	40%(0%)

Table XIV. Detailed means and standard deviations (in parentheses) of percentages of code inspected to find all bugs in the Java programs [Part II]. The star (*) marks the measure(s) with the lowest mean and the underline marks measures that have a mean within 10%.

It is interesting to note that Klosgen (M_{18}) and Normalized Mutual Information (M_{26}) perform comparably with Ochiai. Klosgen (M_{18}), Normalized Mutual Information (M_{26}), and Ochiai are significantly better than Tarantula. It is also interesting to note that 12 other measures also perform significantly better than Tarantula. These are: ϕ -coefficient (M_1), Added Value (M_{15}), Collective Strength (M_{16}), J-Measure (M_6), Information Gain (M_{19}), Two-way Support (M_{28}), Two-way Support Variation (M_{29}), Interestingness Weighting Dependency (M_{24}), Gini Index (M_7), Example and Counterexample Rate (M_{34}), Kappa (M_5), and Ochiai II (M_{40}). Measures that perform comparably with Tarantula are Confidence (M_9), Interest (M_{12}), One-way Support (M_{27}), Sebag (M_{31}), Odds Multiplier (M_{33}), and Zhang (M_{35}). Also it could be noted that Piatetsky-Shapiro's (M_{13}), and Loevinger (M_{30}) perform worse than other measures for fault localization.

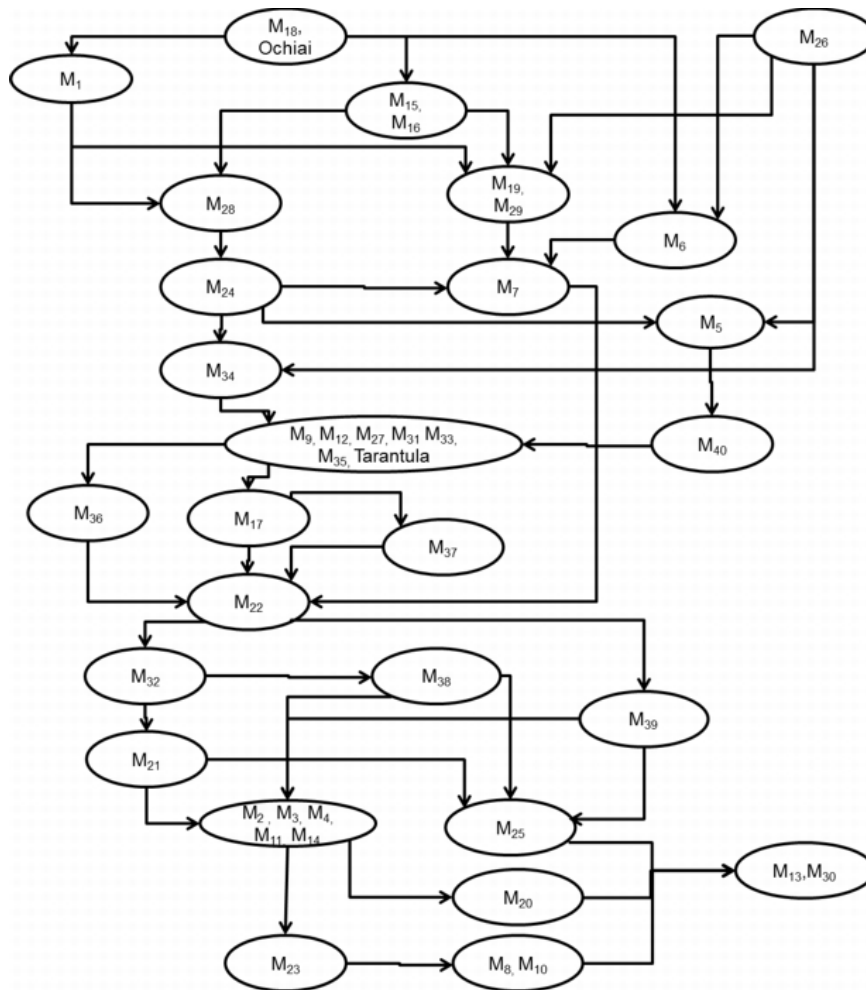


Figure 2. Accuracy partial order

5.3.2. *Proportion of Bugs Localized.* We also plot the curve showing the proportion of code that are investigated (x-axis) vs. the proportion of bugs localized (y-axis) for all dataset. We split the large graphs into several smaller graphs so that measures that have similar accuracies would be grouped together, as shown in Figures 3, 4, 5, 6, 7, and 8. For each graph, we compare a number of association measures with Tarantula and Ochiai.

Association measures included in Figure 3 perform better than Ochiai and Tarantula or as good as Ochiai. When only 10% of program elements are inspected, Tarantula and Ochiai could localize 40% and 47% of the bugs. Klogen (M_{18}) and Added Value (M_{15}) could localize more bugs than Tarantula and Ochiai—they could localize 49% and 50% of the bugs respectively. Two-way Support (M_{28}) could localize the same proportion of bugs as Ochiai.

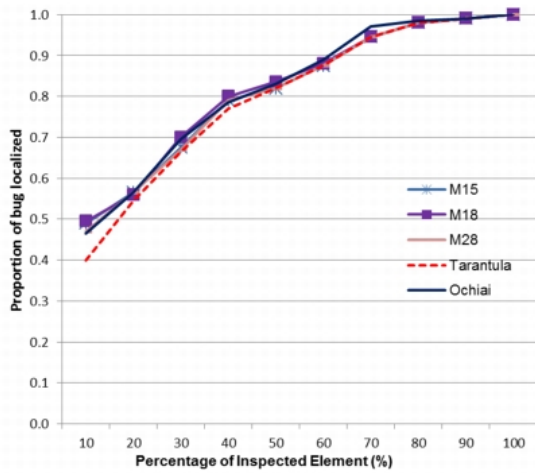


Figure 3. Comparing M_{15} , M_{18} , and M_{28} with Ochiai and Tarantula for all datasets

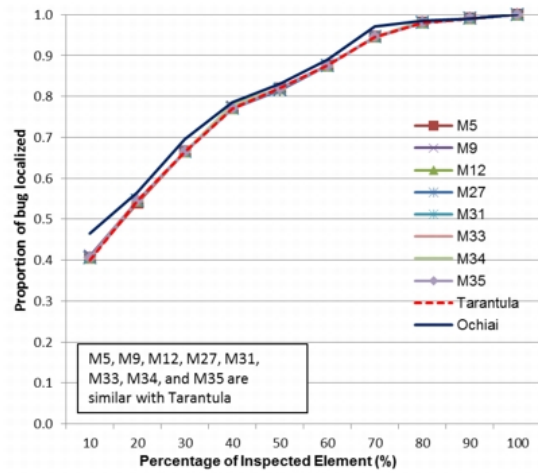


Figure 4. Comparing M_5 , M_9 , M_{12} , M_{27} , M_{31} , M_{33} - M_{35} with Ochiai and Tarantula for all datasets

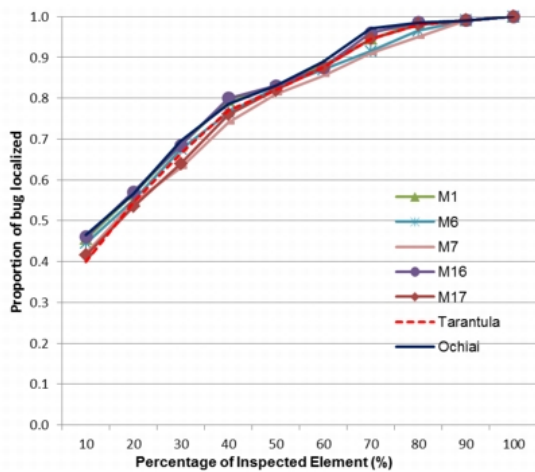


Figure 5. Comparing M_1 , M_6 , M_7 , M_{16} , M_{17} with Ochiai and Tarantula for all datasets

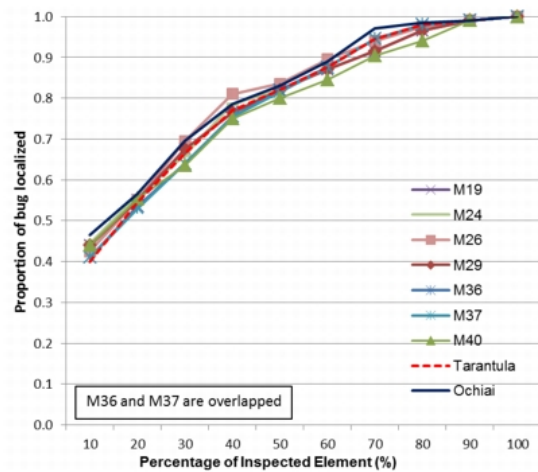


Figure 6. Comparing M_{19} , M_{24} , M_{26} , M_{29} , M_{36} , M_{37} , M_{40} with Ochiai and Tarantula for all datasets

Figures 5 and 6 show association measures that perform better than Tarantula but not as good as Ochiai when only 10% of program elements are inspected. The measures are ϕ -coefficient (M_1), J-Measure (M_6), Gini Index (M_7), Collective Strength (M_{16}), Jaccard (M_{17}), Information Gain (M_{19}), Interestingness Weighting Dependency (M_{24}), Normalized Mutual Information (M_{26}), Loevinger (M_{29}), Sorensen-Dice (M_{36}), Anderberg (M_{37}), and Ochiai II (M_{40}).

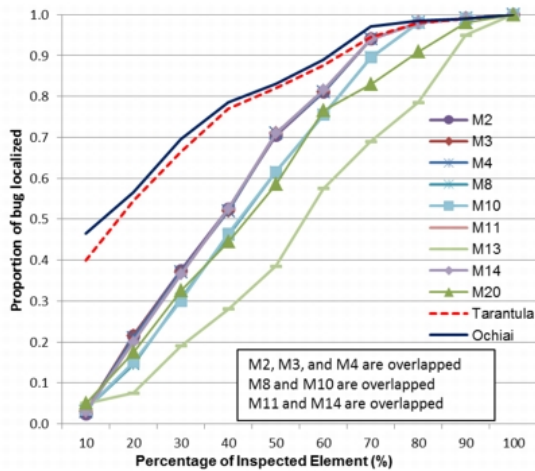


Figure 7. Comparing $M_2 - M_4$, M_8 , M_{10} , M_{11} , M_{13} , M_{14} , M_{20} with Ochiai and Tarantula for all datasets

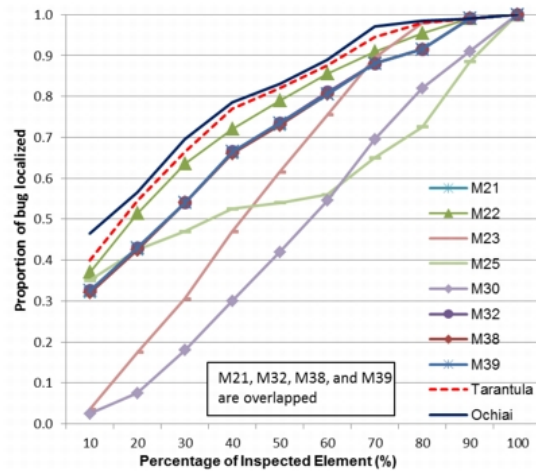


Figure 8. Comparing $M_{21} - M_{23}$, M_{25} , M_{30} , M_{32} , M_{38} , M_{39} with Ochiai and Tarantula for all datasets

Measures that perform similar to Tarantula are shown in Figure 4. They localize 41% of the bugs when 10% of program elements are inspected. The association measures included in Figures 7 and 8 perform worse than Tarantula and Ochiai.

5.4. Effectiveness for Various Programming Languages.

We are interested to evaluate the accuracy of measures for different programming language (C and Java). We find that for different programming language, measures could perform differently. Based on the first accuracy criterion (i.e., percentage of code inspected), we generate two partial orders for the C and Java programs separately to evaluate which measures perform well in each of the programming languages. We highlight measures that are at the top of the partial orders (i.e., no other measures perform statistically significantly better than them). For C programs, Ochiai and Klosgen (M_{18}) are measures that are at the top of the partial order. For Java programs, a number of measures are at the top of the partial order namely Klosgen (M_{18}), Ochiai, Confidence (M_9), Interest (M_{12}), Added Value (M_{15}), Sebag (M_{31}), Example and Counterexample Rate (M_{34}), Zhang (M_{35}), Tarantula, Interestingness Weighting Dependency (M_{24}), and Normalized Mutual Information (M_{26}).

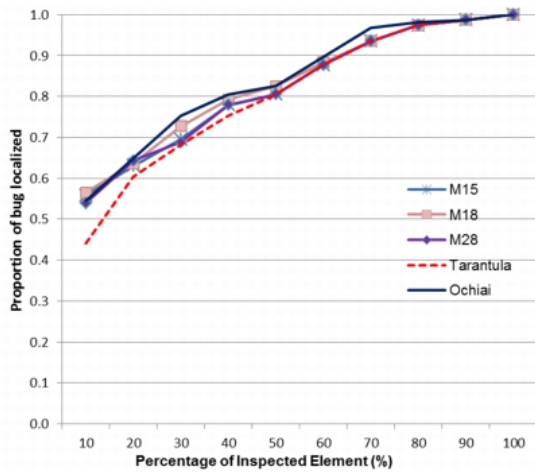


Figure 9. Comparing M_{15} , M_{18} , and M_{28} with Ochiai and Tarantula for C programs

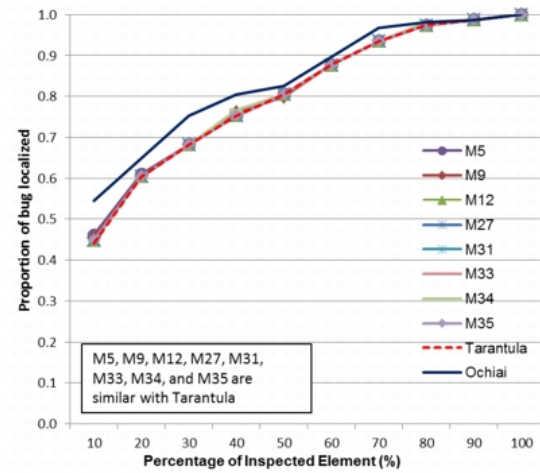


Figure 10. Comparing M_5 , M_9 , M_{12} , M_{27} , M_{31} , M_{33} , M_{34} , and M_{35} with Ochiai and Tarantula for C programs

Next, we evaluate the measures based on the second accuracy criterion (i.e., proportion of bugs found). For localizing bugs in C programs, Tarantula and Ochiai could localize 44% and 55% of total bugs within 10% of inspected program elements. Ochiai performs better than Tarantula. Figure 9 show measures that perform better than Ochiai and Tarantula or similar with Ochiai when only 10% of program elements are inspected. Added Value (M_{15}) and Klosgen (M_{18}) could localize 56% and 56% of the bugs respectively. Two-way Support (M_{28}) could localize slightly lower than Ochiai, 54%.

Figures 11 and 12 show measures that perform better than Tarantula but not as good as Ochiai. ϕ -Coefficient (M_1), J-Measure (M_6), Gini Index (M_7), Collective Strength (M_{16}), Jaccard (M_{17}), Information Gain (M_{19}), Interestingness Weighting Dependency (M_{24}), Normalized Mutual Information (M_{26}), Loevinger (M_{29}), Sorensen-Dice (M_{36}), Anderberg (M_{37}), and Ochiai II (M_{40}) could localize 53%, 52%, 49%, 53%, 48%, 51%, 50%, 50%, 51%, 47%, 47%, and 51% of the bugs respectively. Measures that perform similar to Tarantula are shown in Figure 10. They localize 45-46% of the bugs when 10% of program elements are inspected. The association measures included in Figures 13 and 14 perform worse than Tarantula and Ochiai.

For localizing bugs in Java programs, Tarantula performs better than Ochiai. Notice that when 10% of program elements are inspected, Tarantula and Ochiai could localize 26% and 20% of

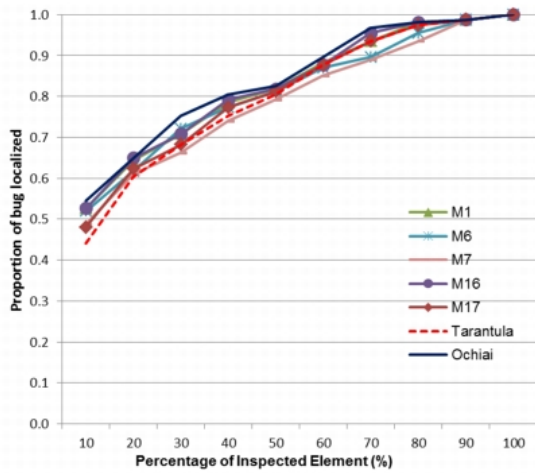


Figure 11. Comparing M_1 , M_6 , M_7 , M_{16} , M_{17} with Ochiai and Tarantula for C programs

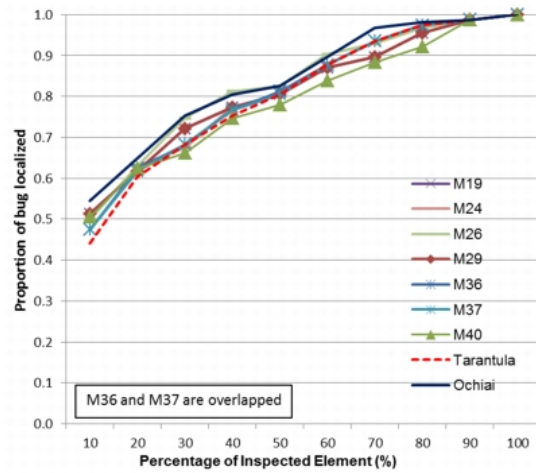


Figure 12. Comparing M_{19} , M_{24} , M_{26} , M_{29} , M_{36} , M_{37} , M_{40} with Ochiai and Tarantula for C programs

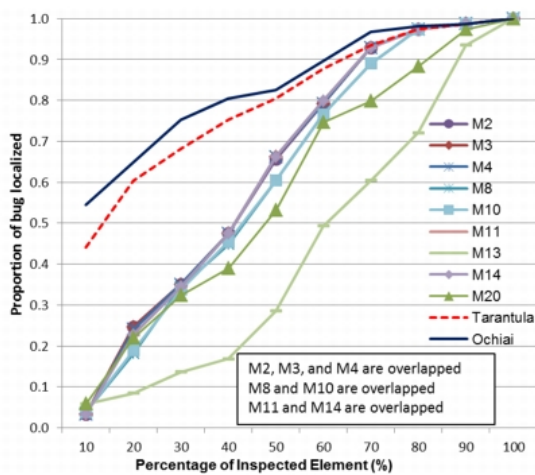


Figure 13. Comparing $M_2 - M_4$, M_8 , M_{10} , M_{11} , M_{13} , M_{14} , M_{20} with Ochiai and Tarantula for C programs

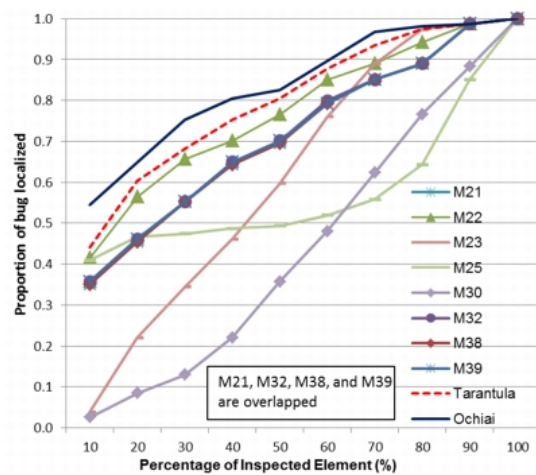


Figure 14. Comparing $M_{21} - M_{23}$, M_{25} , M_{30} , M_{32} , M_{38} , M_{39} with Ochiai and Tarantula for C programs

the bugs. Figure 15 shows measures that have similar performance as Tarantula. They could localize 26% of the bugs within 10% of inspected program elements. The measures are Confidence (M_9), Interest (M_{12}), Added Value (M_{15}), Sebag (M_{31}), Odd Multiplier (M_{33}), and Example and Counterexample Rate (M_{34}).

Figures 16, 17, and 18 show measures that have performance between Tarantula and Ochiai when 10% of program elements are inspected. Figures 19 and 20 show measures that perform

worse than Ochiai and Tarantula when localizing bugs in Java programs. Based on the results, we notice that there are a number of measures that could perform as good as Tarantula and Ochiai when using different programming language.

The differences in accuracies for localizing bug in C and Java programs possibly imply that there could be specific characteristic of spectra produced by different program language that could advantage or disadvantage spectrum-based fault localization. For example in Java, due to object oriented model, some program elements are executed more often than others e.g., class construction elements. Thus, their suspiciousness scores are the same.

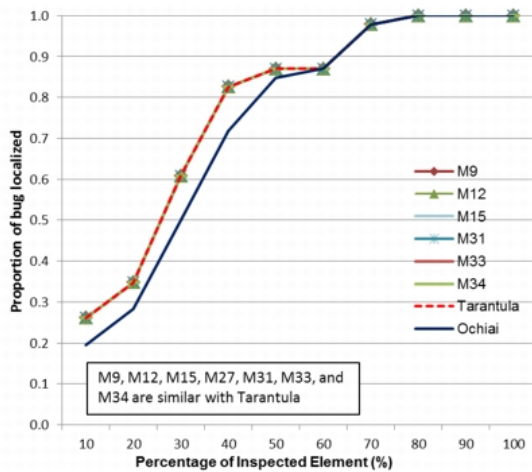


Figure 15. Comparing M_9 , M_{12} , M_{15} , M_{27} , M_{31} , M_{33} , M_{34} with Ochiai and Tarantula for Java programs

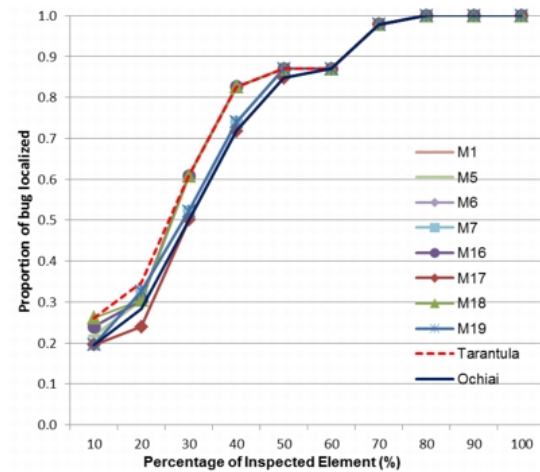


Figure 16. Comparing M_1 , M_5 , M_6 , M_7 , M_{16} , M_{17} , M_{18} , M_{19} with Ochiai and Tarantula for Java programs

5.5. Effectiveness for Various Kinds of Bugs.

We divide the bugs into several groups; our categorization is based on that by Kim et al. [29]. Kim et al. categorize different fixes that are made to various software systems. We categorize bugs based on how the bugs get fixed. The following paragraphs describe our bug categories and present the effectiveness of the various association measures, Tarantula, and Ochiai on the different bug categories.

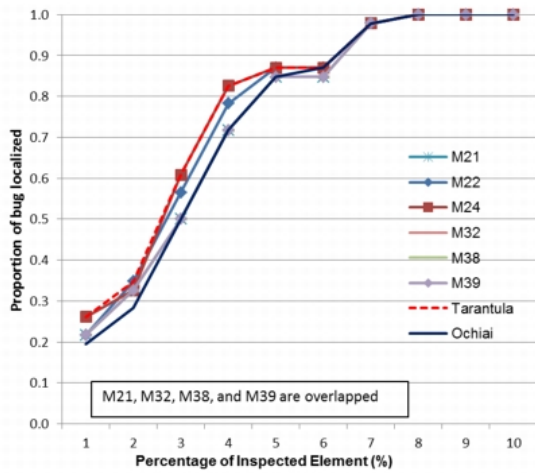


Figure 17. Comparing M_{21} , M_{22} , M_{24} , M_{32} , M_{38} , M_{39} with Ochiai and Tarantula for Java programs

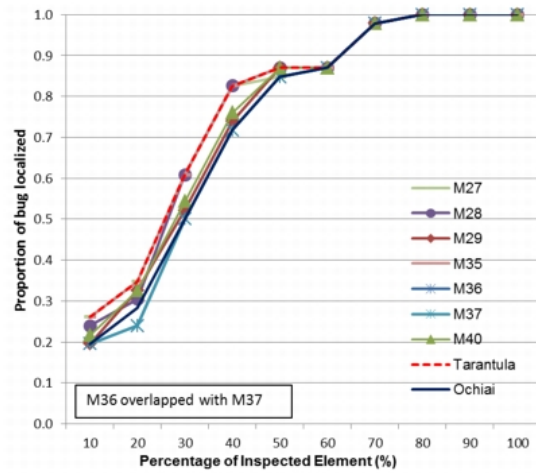


Figure 18. Comparing M_{27} - M_{29} , M_{35} - M_{37} , M_{40} with Ochiai and Tarantula for Java programs

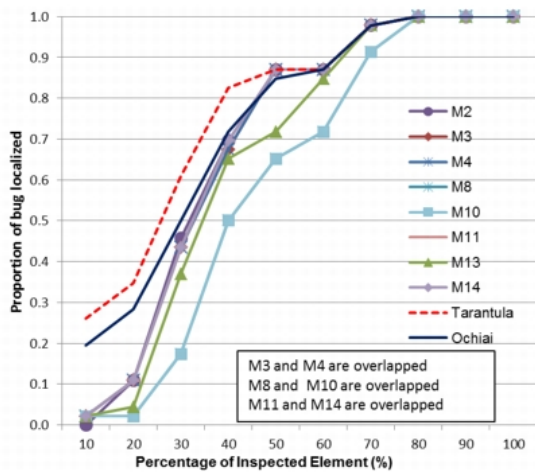


Figure 19. Comparing M_2 - M_4 , M_8 , M_{10} , M_{11} - M_{14} with Ochiai and Tarantula for Java programs

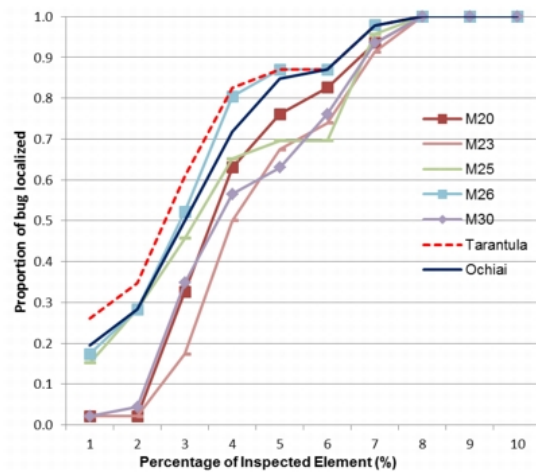


Figure 20. Comparing M_{20} , M_{23} , M_{25} , M_{26} , M_{30} with Ochiai and Tarantula for Java programs

5.5.1. Bug Categories. Kim et al. [29] describes a number of bug categories based on the way bugs are fixed. In this paper, we refer to their categories to analyze the bug in our subject programs. We categorize bugs in our subject programs into 8 categories. Table XV shows the categories and number of buggy versions for each type of bug. We add two categories that are not included by Kim et al. [29], i.e., change of return expression (CH-RET) and others (OTH).

Bug category	Total version
Addition/removal of conditional statement (CH-CS)	28
Addition/removal of non-conditional statement (CH-NCS)	10
Change in method calls (CH-MC)	26
Change of assignment expression (AS-CE)	60
Change of if condition expression (IF-CC)	44
Change of loop predicate (LP-CC)	11
Change of return expression (CH-RET)	18
Others (OTH)	3

Table XV. Bug categories

We categorize a bug into addition or removal of conditional statement category (CH-CS) when a bug could be fixed by adding a conditional check statement (e.g. if statement) or removing an inappropriate check statement. This type of bugs occurs when there is a missing precondition or postcondition check of some variables, or an extraneous conditional check. Bugs that could be fixed by adding or removing statement which is not a conditional check statement are put into addition/removal of non-conditional statement (CH-NCS) category. An example of this type of bug is missing or extraneous assignment statements.

Change in method calls (CH-MC) category includes bugs that can be fixed by adding or removing a method call in a program, or by changing the parameter values of a method call. Bugs that could be fixed by changing the right hand side of an assignment expression are put into change of assignment expression (AS-CE) category. When a bug could be fixed by modifying a conditional expression within an if statement, we put the bug into change of if condition expression (IF-CC) category. Similarly, when a bug could be fixed by modifying a conditional expression in a looping statement, then we put this bug into change of loop predicate (LP-CC) category. When a bug could be fixed by changing the value or the expression in a return statement, then we put the bug into change of return expression (CH-RET) category.

We create a category named other (OTH) to include other bugs that are not covered by the above categories, e.g., a bug is fixed by changing an if statement to a for loop statement, etc. We ignore bugs in category (OTH).

5.5.2. *Effectiveness.* We first evaluate the effectiveness of the various measures on each category by the first accuracy criterion (i.e., percentage of code inspected). Based on this criterion, we create several partial orders based on statistically significantly better relationships among the measures. We highlight measures that are at the top of the partial orders as shown in Table XVI.

Based on these partial orders, Tarantula is at the top of the partial orders of four bug categories. Ochiai, Information Gain (M_{19}), Normalized Mutual Information (M_{26}), Two-way Support Variation (M_{29}) are at the top of the partial orders of six bug categories. Klosgen is at the top of the partial orders of all bug categories.

Next, we evaluate the effectiveness of the measures on each bug category by the second accuracy criterion (i.e., proportion of bugs localized). We compute the percentage of bugs localized when up to 10% of program elements are inspected, as shown in Table XVII and XVIII.

For each category shown in Table XVII and XVIII, different measures have different effectiveness to localize bugs. The star (*) marks the measures that could localize most bugs in each category. J-Measure (M_6), Gini index (M_7), Klosgen (M_{18}), Information Gain (M_{19}), Two-way Support (M_{28}), and Two-way Support Variation (M_{29}) could localize the most number of bugs in addition or removal of conditional statement (CH-CS) category. They could localize 50% of the bugs in this category. ϕ -Coefficient (M_1), Added Value (M_{15}), Collective Strength (M_{16}), Klosgen (M_{18}), Interestingness Weighting Dependency (M_{24}), Two-way Support (M_{28}), and Ochiai could better localize bugs in addition or removal of non-conditional statement (CH-NCS) category than other measures. They could localize 70% of the bugs in this category.

For localizing change in method calls (CH-MC) bugs, at most 35% of the bugs in this category could be localized by the measures. There are number of measures that could localize 35% of the bugs in this category, they are Confidence (M_9), Interest (M_{12}), Added Value (M_{15}), Collective Strength (M_{16}), Klosgen (M_{18}), Accuracy (M_{21}), Interestingness Weighting Dependency (M_{24}), One-way Support (M_{27}), Two-way Support (M_{28}), Sebag (M_{31}), Least Contradiction (M_{32}), Odd Multiplier (M_{33}), Example and Counterexample Rate (M_{34}), Zhang (M_{35}), Simple-Matching (M_{38}), Rogers and Tanimoto (M_{39}), and Tarantula.

Bug category	Top Measures
Addition/removal of conditional statement (CH-CS)	<u>Klosgen</u> (M_{18}), <u>Ochiai</u> , <u>Information Gain</u> (M_{19}), <u>Normalized Mutual Information</u> (M_{26}), <u>Two-way Support Variation</u> (M_{29}), Tarantula, Coefficient (M_1), Kappa (M_5), J-Measure (M_6), Confidence (M_9), Added Value (M_{15}), Collective Strength (M_{16}), Interestingness Weighting Dependency (M_{24}), One-way Support (M_{27}), Two-way Support (M_{28}), Sebag (M_{31}), Odd Multiplier (M_{33}), Example and Counterexample Rate (M_{34}), Zhang (M_{35}), Sorensen-Dice (M_{36}), Anderberg (M_{37}), and Ochiai II (M_{40})
Addition/removal of non-conditional statement (CH-NCS)	<u>Klosgen</u> (M_{18}), <u>Ochiai</u> , <u>Information Gain</u> (M_{19}), <u>Normalized Mutual Information</u> (M_{26}), <u>Two-way Support Variation</u> (M_{29}), Coefficient (M_1), Kappa (M_5), Gini Index (M_7), J-Measure (M_6), Added Value (M_{15}), Collective Strength (M_{16}), Jaccard (M_{17}), Interestingness Weighting Dependency (M_{24}), Two-way Support (M_{28}), Anderberg (M_{37}), and Ochiai II (M_{40})
Change in method calls (CH-MC)	<u>Klosgen</u> (M_{18}), <u>Ochiai</u> , <u>Information Gain</u> (M_{19}), <u>Normalized Mutual Information</u> (M_{26}), <u>Two-way Support Variation</u> (M_{29}), Tarantula, Coefficient (M_1), Confidence (M_9), Interest (M_{12}), Added Value (M_{15}), Collective Strength (M_{16}), Interestingness Weighting Dependency (M_{24}), One-way Support (M_{27}), Two-way Support (M_{28}), Sebag (M_{31}), Odd Multiplier (M_{33}), Example and Counterexample Rate (M_{34}), Zhang (M_{35}), and Ochiai II (M_{40})
Change of assignment expression (AS-CE)	<u>Klosgen</u> (M_{18}), <u>Ochiai</u> , and <u>Normalized Mutual Information</u> (M_{26})
Change of if condition expression (IF-CC)	<u>Klosgen</u> (M_{18}), <u>Ochiai</u> , <u>Information Gain</u> (M_{19}), <u>Normalized Mutual Information</u> (M_{26}), <u>Two-way Support Variation</u> (M_{29}), and J-Measure (M_6)
Change of loop predicate (LP-CC)	<u>Klosgen</u> (M_{18}), <u>Ochiai</u> , <u>Information Gain</u> (M_{19}), <u>Normalized Mutual Information</u> (M_{26}), <u>Two-way Support Variation</u> (M_{29}), Tarantula, Coefficient (M_1), Kappa (M_5), Collective Strength (M_{16}), J-Measure (M_6), Gini Index (M_7), Confidence (M_9), Added Value (M_{15}), Jaccard (M_{17}), Interestingness Weighting Dependency (M_{24}), Goodman Kruskal (M_{25}), One-way Support (M_{27}), Two-way Support (M_{28}), Sebag (M_{31}), Example and Counterexample Rate (M_{34}), Zhang (M_{35}), Sorensen-Dice (M_{36}), Anderberg (M_{37}), and Ochiai II (M_{40})
Change of return expression (CH-RET)	<u>Klosgen</u> (M_{18}), <u>Information Gain</u> (M_{19}), <u>Normalized Mutual Information</u> (M_{26}), <u>Two-way Support Variation</u> (M_{29}), Tarantula, Coefficient (M_1), J-Measure (M_6), Gini Index (M_7), Confidence (M_9), Interest (M_{12}), Added Value (M_{15}), Leverage (M_{21}), Interestingness Weighting Dependency (M_{24}), One-way Support (M_{27}), Two-way Support (M_{28}), Sebag (M_{31}), Odd Multiplier (M_{33}), Example and Counterexample Rate (M_{34}), Zhang (M_{35}), and Ochiai II (M_{40})

Table XVI. Measures that are at the top of the partial orders for each bug categories. The underline marks measures that are at the top of the partial orders of more than five bug categories.

Measures	CH-CS	CH-NCS	CH-MC	AS-CE	IF-CC	LP-CC	CH-RET
ϕ -Coefficient(M_1)	46%	*70%	31%	47%	41%	*73%	50%
Odds Ratio(M_2)	0%	0%	0%	3%	5%	9%	0%
Yule's Q(M_3)	0%	0%	4%	3%	5%	9%	0%
Yule's Y(M_4)	0%	0%	4%	3%	5%	9%	0%
Kappa(M_5)	43%	60%	31%	43%	30%	64%	50%
J-Measure(M_6)	*50%	60%	27%	45%	*48%	45%	50%
Gini Index(M_7)	*50%	60%	27%	45%	36%	45%	50%
Support(M_8)	0%	0%	4%	5%	5%	9%	0%
Confidence(M_9)	43%	60%	*35%	40%	32%	64%	56%
Laplace(M_{10})	0%	0%	4%	5%	5%	9%	0%
Conviction(M_{11})	0%	0%	4%	3%	5%	9%	0%
Interest(M_{12})	43%	60%	*35%	40%	30%	64%	56%
Pietatsky-Shapiro(M_{13})	7%	0%	4%	8%	5%	0%	0%
Certainty Factor(M_{14})	0%	0%	4%	3%	5%	9%	0%
Added Value(M_{15})	50%	*70%	*35%	*48%	*48%	*73%	56%
Collective Strength(M_{16})	46%	*70%	*35%	*48%	39%	*73%	50%
Jaccard(M_{17})	43%	60%	31%	43%	32%	*73%	50%
Kloggen(M_{18})	*50%	*70%	*35%	*48%	*48%	*73%	*61%
Information Gain(M_{19})	*50%	60%	27%	45%	45%	45%	50%
Coverage(M_{20})	4%	0%	8%	5%	7%	9%	0%
Accuracy(M_{21})	43%	40%	*35%	37%	18%	18%	39%
Leverage(M_{22})	43%	40%	31%	35%	27%	64%	56%
Relative Risk(M_{23})	0%	0%	4%	5%	5%	9%	0%
Int. Weighting Dependency(M_{24})	43%	*70%	*35%	43%	39%	*73%	56%
GoodMan and Kruskal(M_{25})	39%	30%	23%	42%	25%	45%	50%
Normalized Mutual Info.(M_{26})	46%	60%	27%	43%	43%	45%	50%
One-Way Support(M_{27})	43%	60%	*35%	40%	30%	64%	56%
Two-Way Support(M_{28})	*50%	*70%	*35%	47%	43%	*73%	50%
Two-Way Support Variation(M_{29})	*50%	60%	27%	45%	45%	45%	50%
Loevinger(M_{30})	0%	0%	4%	5%	0%	9%	0%
Sebag(M_{31})	43%	60%	*35%	40%	30%	64%	56%
Least Contradiction(M_{32})	43%	40%	*35%	37%	20%	18%	39%
Odd Multiplier(M_{33})	43%	60%	*35%	40%	30%	64%	56%
Example and Counter.(M_{34})	43%	60%	*35%	40%	30%	64%	56%
Zhang(M_{35})	43%	60%	*35%	40%	30%	64%	56%

Table XVII. Effectiveness of bug localization for (M_1) to (M_{35}) for each bug category when up to 10% program elements are inspected [Part I]. The star (*) marks the measures that could localize most of the bugs for each category.

Measures	CH-CS	CH-NCS	CH-MC	AS-CE	IF-CC	LP-CC	CH-RET
Sorensen-Dice (M_{36})	43%	60%	31%	43%	30%	*73%	50%
Anderberg (M_{37})	43%	60%	31%	43%	30%	*73%	50%
Simple-Matching (M_{38})	43%	40%	*35%	37%	18%	18%	39%
Rogers and Tanimoto (M_{39})	43%	40%	*35%	37%	20%	18%	39%
Ochiai II (M_{40})	43%	40%	31%	45%	39%	*73%	50%
Tarantula	43%	60%	*35%	40%	27%	64%	56%
Ochiai	46%	*70%	31%	45%	*48%	*73%	50%

Table XVIII. Effectiveness of bug localization for (M_{36}) to (M_{40}), Tarantula, and Ochiai for each bug category when up to 10% program elements are inspected [Part II]. The star (*) marks the measures that could localize most of the bugs for each category.

Added Value (M_{15}), Collective Strength (M_{16}), Klosgen (M_{18}) could localize 48% of the bugs in change of assignment expression (AS-CE) category. J-Measure (M_6), Added Value (M_{15}), Klosgen (M_{18}), and Ochiai could localize 48% of the bugs in change of if condition expression (IF-CC) category.

ϕ -Coefficient (M_1), Added Value (M_{15}), Collective Strength (M_{16}), Jaccard (M_{17}), Klosgen (M_{18}), Interestingness Weighting Dependency (M_{24}), Two-way Support (M_{27}), Sorensen-Dice (M_{36}), Anderberg (M_{37}), Ochiai II (M_{40}), and Ochiai could localize 73% of the bugs in change of loop predicate (LP-CC) category. For bugs in change of return expression (CH-RET) category, Klosgen (M_{18}) could localize these bugs better than others, i.e., 61% of the bugs could be localized.

By only inspecting up to 10% of the program elements, bugs in change in method calls (CH-MC) category are not easy to be localized. The best measures could only localize 35% of these bugs. On the other hand, bugs in addition or removal of non-conditional statement (CH-NCS) and change of loop predicate (LP-CC) categories could be better localized by the measures (i.e., up to 70% and 73% respectively). It is interesting to note that Klosgen (M_{18}) is the measure that could localize the most number of bugs in all categories, followed by Added Value (M_{15}) that could localize the most number of bugs in five categories. Two-way Support (M_{28}) and Collective Strength (M_{16}) could localize the most number of bugs in four bug categories. Ochiai and Tarantula could localize the most number of bugs in 3 and 1 categories respectively.

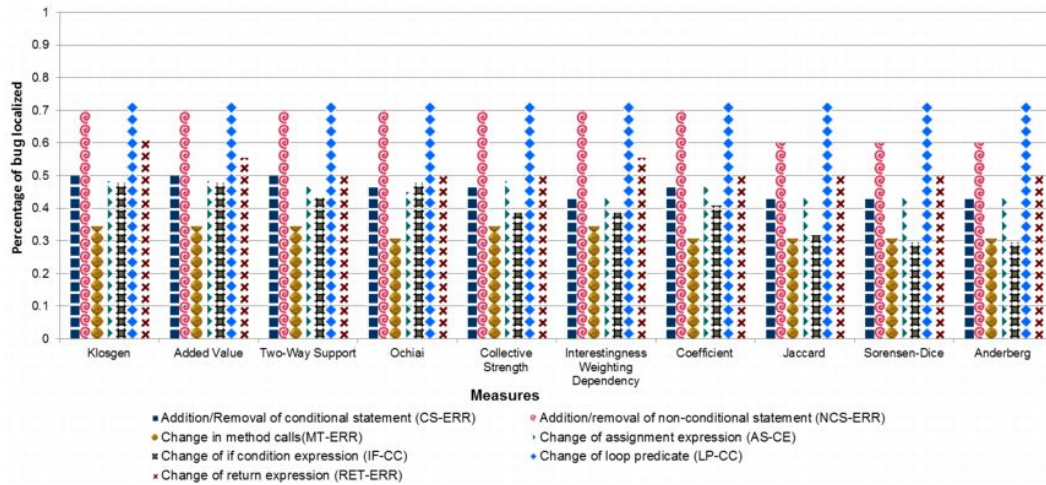


Figure 21. Effectiveness of measures to localize bugs by inspecting up to 10% of code [Part I]

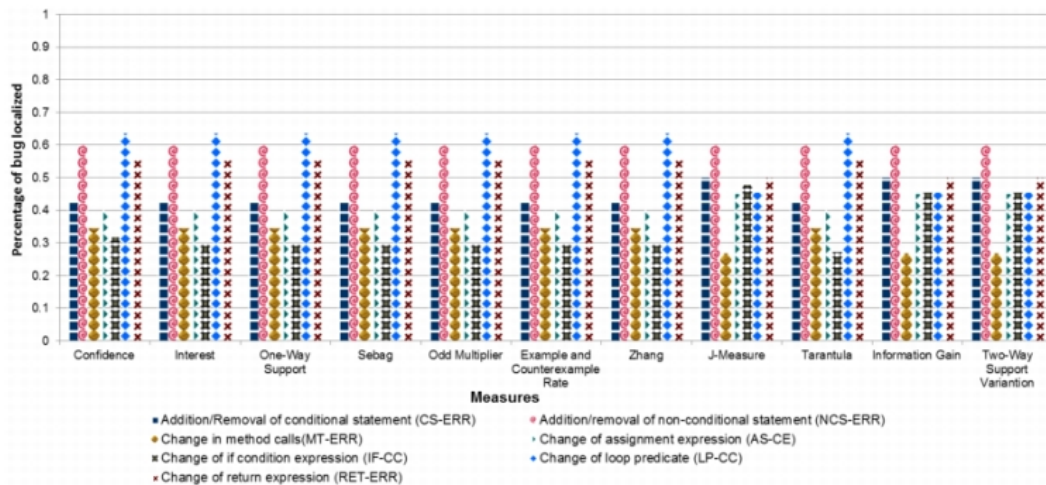


Figure 22. Effectiveness of measures to localize bugs by inspecting up to 10% of code [Part II]

Figures 21, 22, 23 and 24 show the effectiveness of the measures in localizing bugs for each category. As shown in Figure 24, Odds Ratio (M_2), Yule's Q (M_3), Yule's Y (M_4), Support (M_8), Laplace (M_{10}), Conviction (M_{11}), Pietatsky-Shapiro (M_{13}), Certainty Factor (M_{14}), Relative Risk (M_{23}), Loevinger (M_{30}), and Coverage (M_{20}) could only localize a small number of bugs for each bug category. Figures 21, 22, and 23 show measures that could localize more bugs for most of the categories. A number of measures could localize at least 50% of the bugs in addition or removal of non-conditional statement (CH-NCS), change of loop predicate (LP-CC), and change of return statement (CH-RET) categories.

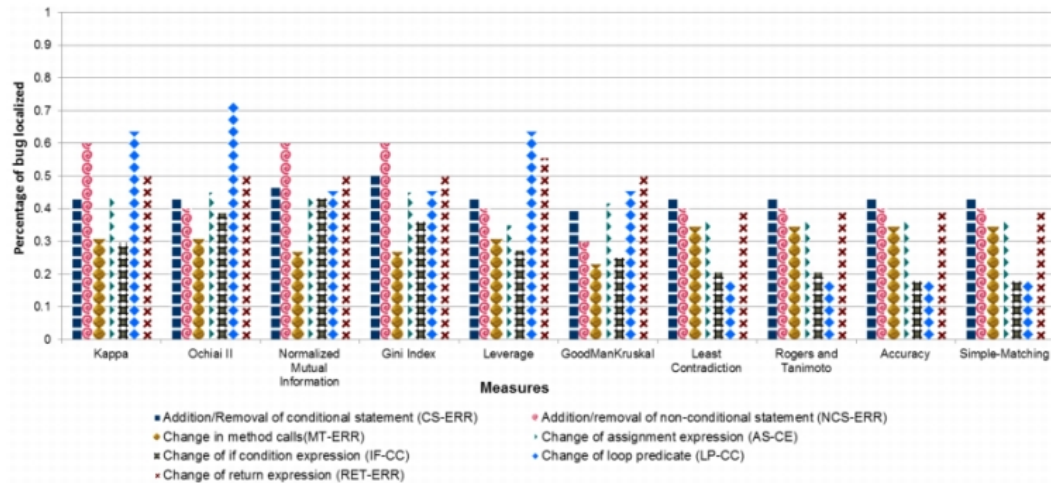


Figure 23. Effectiveness of measures to localize bugs by inspecting up to 10% of code [Part III]

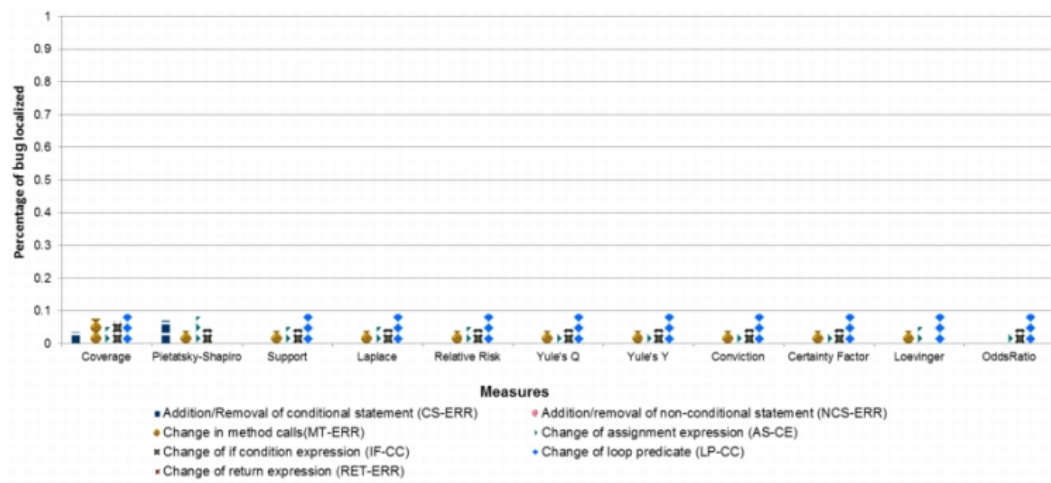


Figure 24. Effectiveness of measures to localize bugs by inspecting up to 10% of code [Part IV]

5.6. Effectiveness on Multiple-bug Versions

We evaluate the effectiveness of association measures, Tarantula, and Ochiai to localize multiple bugs in programs. We refer these programs as multiple-bug versions. A multiple-bug version contains a number of bugs where each bug only involves one line in the program and different bugs affect different lines [15, 17].

We generate 173 multiple-bug versions of C programs as shown in Table XIX. As print_token and schedule2 datasets only have 4 and 7 bugs that involve one line, we randomly insert two bugs for every version, while for other datasets, we randomly insert five bugs for every multiple-bug version. Also, we ensure that each bug has been inserted at least in one of the versions.

Dataset	Num. of Bug in a Version	Num. of Single Bug in Dataset	Language	Num. of Buggy Version
print.token	2	4	C	10
schedule2	2	5	C	10
print.token2	5	9	C	10
replace	5	25	C	32
schedule	5	7	C	9
tcas	5	30	C	41
tot.info	5	19	C	23
space	5	19	C	38

Table XIX. Multiple-bug datasets

We generate multiple-bug versions for each dataset as many as the number of single bug versions in the dataset. For example, there are 38 single bug versions for Space, so we randomly generate 38 multiple-bug versions for Space, each of which contains 5 bugs. For each print.tokens and schedule2 dataset, we generate 10 multiple-bug versions. Thus, we have 20 multiple-bug versions that contain two bugs (minimum number of multiple bugs) and 153 versions that contain five bugs.

We evaluate the effectiveness of 40 association measures, Tarantula, and Ochiai to localize multiple bugs in programs. Tables XX, XXI, and XXII show the overall mean and standard deviation of percentage of code inspected of the measures to localize bugs in all multiple-bug versions, versions containing five bugs, and versions containing two bugs respectively.

Generally, the effectiveness of all measures to localize multiple bugs are not as good as localizing single bugs. The overall ranges of mean of percentage of code inspected of the measures to localize all multiple-bug versions, five bugs, and two bugs are between 45% to 71%, 43% to 71%, and 34% to 74% respectively. Measures that could localize five bugs in the programs with the smallest percentage of code inspected (43%) are Odds ratio (M_2), Yule's Q (M_3), and Yule's Y (M_4). For localizing five bugs in the programs, the measures are Accuracy (M_{21}), Least Contradiction (M_{32}), Simple-Matching (M_{38}), and Rogers and Tanimoto (M_{39}).

We notice that performance of the measures in localizing five and two bugs in the programs are different. Tarantula performs slightly better than Ochiai when localizing five bugs (i.e., 50.23% for Tarantula, 50.59% for Ochiai), in contrast Ochiai performs better than Tarantula when localizing

Association Measures	Mean	StdDev	Association Measures	Mean	StdDev
Odds ratio (M_2)	44.51%	28.67%	Example Rate (M_{34})	49.57%	33.25%
Yule's Q (M_3)	44.55%	28.81%	Klogsen (M_{18})	49.93%	32.19%
Yule's Y (M_4)	44.69%	28.68%	Kappa (M_5)	50.21%	32.28%
Conviction (M_{11})	45.78%	27.80%	Laplace (M_{10})	50.31%	25.79%
Certainty Factor (M_{14})	45.80%	27.79%	Two Way Support (M_{28})	50.36%	32.01%
Relative Risk (M_{23})	46.06%	28.81%	Rogers and Tanimoto (M_{39})	50.43%	33.74%
Coverage (M_{20})	46.20%	30.40%	Accuracy (M_{21})	50.46%	33.77%
Normalized Mutual Information (M_{26})	47.50%	32.25%	Simple-Matching (M_{38})	50.46%	33.77%
Added Value (M_{15})	48.18%	33.44%	Least Contradiction (M_{32})	50.53%	33.90%
ϕ -Coefficient (M_1)	48.83%	32.20%	Leverage (M_{22})	50.80%	32.48%
One Way Support (M_{27})	49.23%	33.01%	J-Measure (M_6)	51.28%	32.18%
Interestingness Weighting Dependency (M_{24})	49.33%	33.09%	Ochiai II (M_{40})	51.29%	32.22%
Jaccard (M_{17})	49.35%	31.31%	Sorensen-Dice (M_{36})	51.42%	31.19%
Odd Multiplier (M_{33})	49.35%	32.97%	Anderberg (M_{37})	51.45%	31.17%
Zhang (M_{35})	49.36%	32.95%	Information Gain (M_{19})	51.58%	31.97%
Sebag (M_{31})	49.38%	32.93%	Two Way Support Variantion (M_{29})	51.58%	31.97%
Tarantula	49.39%	32.98%	Support (M_8)	51.70%	26.67%
Collective Strength (M_{16})	49.41%	32.07%	Gini Index (M_7)	52.17%	32.18%
Interest (M_{12})	49.41%	32.91%	Loevinger (M_{30})	52.95%	31.65%
Ochiai	49.42%	30.65%	Piatetsky-Shapiro's (M_{13})	62.02%	30.89%
Confidence (M_9)	49.43%	32.92%	GoodMan Kruskal (M_{25})	71.41%	35.57%

Table XX. Overall mean and standard deviation (in parentheses) of accuracy values (smaller the better) of all multiple-bug versions

two bugs in the programs (i.e., 42.95% for Tarantula, 40.50% for Ochiai). There are 8 measures that perform better than Tarantula and Ochiai in localizing two bugs in the programs, while 17 measures perform better for localizing five bugs. Among these measures, Added Value (M_{15}) and Normalized Mutual Information (M_{26}) consistently perform better than Tarantula and Ochiai for localizing programs that contain both two and five bugs. Interestingly these measures also have good performance in localizing single bug. The overall mean of percentage of code inspected for Added Value (M_{15}) and Normalized Mutual Information (M_{26}) are 23.30% and 23.22% respectively, which are only slightly smaller than the smallest mean of percentage of code inspected to localized all single bug programs (22.42%). In addition, Added Value (M_{15}) and Normalized Mutual Information (M_{26}) are in the top of partial order of 5 and 6 bug categories out of 7 categories.

Association Measures	Mean	StdDev	Association Measures	Mean	StdDev
Odds ratio (M_2)	42.88%	29.45%	Example Rate (M_{34})	50.43%	34.19%
Yule's Q (M_3)	42.90%	29.58%	Ochiai	50.59%	30.68%
Yule's Y (M_4)	43.06%	29.45%	Collective Strength (M_{16})	50.62%	32.31%
Relative Risk (M_{23})	44.42%	29.66%	Klosgen (M_{18})	51.16%	32.66%
Conviction (M_{11})	44.54%	28.71%	Kappa (M_5)	51.18%	33.05%
Certainty Factor (M_{14})	44.56%	28.71%	Two Way Support (M_{28})	51.40%	32.62%
Coverage (M_{20})	44.66%	31.14%	Ochiai II (M_{40})	51.65%	33.10%
Normalized Mutual Information (M_{26})	48.54%	32.63%	J-Measure (M_6)	52.17%	32.69%
Laplace (M_{10})	49.20%	26.60%	Anderberg (M_{37})	52.22%	31.59%
Added Value (M_{15})	49.29%	34.25%	Sorensen-Dice (M_{36})	52.24%	31.57%
ϕ -Coefficient (M_1)	49.65%	32.68%	Leverage (M_{22})	52.31%	33.83%
Jaccard (M_{17})	49.90%	31.76%	Information Gain (M_{19})	52.50%	32.46%
One Way Support (M_{27})	50.05%	33.93%	Two Way Support Variantion (M_{29})	52.50%	32.46%
Odd Multiplier (M_{33})	50.18%	33.88%	Rogers and Tanimoto (M_{39})	52.58%	34.50%
Zhang (M_{35})	50.20%	33.86%	Accuracy (M_{21})	52.61%	34.54%
Interestingness Weighting Dependency (M_{24})	50.22%	33.94%	Simple-Matching (M_{38})	52.61%	34.54%
Sebag (M_{31})	50.22%	33.84%	Least Contradiction (M_{32})	52.69%	34.67%
Tarantula	50.23%	33.89%	Gini Index (M_7)	53.03%	32.89%
Interest (M_{12})	50.25%	33.82%	Loevinger (M_{30})	53.51%	32.14%
Confidence (M_9)	50.28%	33.83%	Piatetsky-Shapiro's (M_{13})	64.78%	29.65%
Support (M_8)	50.41%	27.30%	GoodMan Kruskal (M_{25})	71.12%	36.17%

Table XXI. Overall mean and standard deviation (in parentheses) of accuracy values (smaller the better) of versions containing five bugs

We also perform statistical tests for each pair of measures including Tarantula and Ochiai using Wilcoxon signed rank test [69] at 0.05 statistical significance threshold to see if some measures are statistically significantly better than others in localizing multiple-bug versions. Table XXIII shows the measures that are on the top of the partial order (no other measures that statistically significantly perform better than the measure) for localizing both all multiple-bug versions, versions containing five bugs, and versions containing two measures. We notice that Odds ratio (M_2), Yule's Q (M_3), Added Value (M_{15}), and Coverage (M_{20}) are in the top of the partial order of all multiple-bug versions and at least in the top order of one of partial order of two or five bugs versions. They are statistically significantly better than Tarantula and Ochiai for all multiple-bug versions.

We also plot the curve showing the proportion of code that are investigated (x-axis) vs. the proportion of bugs localized (y-axis) for all multiple-bug versions. We split the large graphs into

Association Measures	Mean	StdDev	Association Measures	Mean	StdDev
Accuracy (M_{21})	34.00%	21.43%	Zhang (M_{35})	42.95%	24.62%
Least Contradiction (M_{32})	34.00%	21.43%	Tarantula	42.95%	24.62%
Simple-Matching (M_{38})	34.00%	21.43%	J-Measure (M_6)	44.45%	27.69%
Rogers and Tanimoto (M_{39})	34.00%	21.43%	Information Gain (M_{19})	44.55%	27.62%
Leverage (M_{22})	39.30%	15.39%	Two Way Support Variantion (M_{29})	44.55%	27.62%
Normalized Mutual Information (M_{26})	39.55%	28.69%	Jaccard (M_{17})	45.15%	28.04%
Added Value (M_{15})	39.70%	25.55%	Sorensen-Dice (M_{36})	45.15%	28.04%
Collective Strength (M_{16})	40.15%	29.21%	Anderberg (M_{37})	45.50%	27.69%
Ochiai	40.50%	29.68%	Gini Index (M_7)	45.60%	25.87%
Klosgen (M_{18})	40.55%	27.20%	Ochiai II (M_{40})	48.55%	25.01%
Piatetsky-Shapiro's (M_{13})	40.85%	32.75%	Loevinger (M_{30})	48.65%	27.97%
Two Way Support (M_{28})	42.45%	26.25%	Conviction (M_{11})	55.30%	17.14%
Interestingness Weighting Dependency (M_{24})	42.55%	25.33%	Certainty Factor (M_{14})	55.30%	17.14%
ϕ -Coefficient (M_1)	42.60%	28.15%	Odds ratio (M_2)	57.00%	17.70%
Kappa (M_5)	42.80%	25.05%	Yule's Q (M_3)	57.20%	17.98%
Confidence (M_9)	42.95%	24.62%	Yule's Y (M_4)	57.20%	17.98%
Interest (M_{12})	42.95%	24.62%	Coverage (M_{20})	58.00%	21.11%
One Way Support (M_{27})	42.95%	24.62%	Relative Risk (M_{23})	58.60%	16.90%
Sebag (M_{31})	42.95%	24.62%	Laplace (M_{10})	58.80%	16.61%
Odd Multiplier (M_{33})	42.95%	24.62%	Support (M_8)	61.55%	18.97%
Example Rate (M_{34})	42.95%	24.62%	GoodMan Kruskal (M_{25})	73.65%	31.30%

Table XXII. Overall mean and standard deviation (in parentheses) of accuracy values (smaller the better) of versions containing two bugs

Num. of Bugs in a Version	Top Measures
Two and five bugs	<u>Odds Ratio</u> (M_2), <u>Yule's Q</u> (M_3), <u>Added Value</u> (M_{15}), and <u>Coverage</u> (M_{20})
Five bugs	<u>Odds Ratio</u> (M_2), <u>Yule's Q</u> (M_3), and <u>Coverage</u> (M_{20})
Two bugs	<u>Added Value</u> (M_{15}), Kappa (M_5), Piatetsky-Shapiro's (M_{13}), Collective Strength (M_{16}), Klosgen (M_{18}), Accuracy (M_{21}), Leverage (M_{22}), Normalized Mutual Information (M_{26}), Loevinger (M_{30}), Least Contradiction (M_{32}), Simple-Matching (M_{38}), and Rogers and Tanimoto (M_{39})

Table XXIII. Measures that are at the top of the partial orders for different number of bugs within a buggy version. The underline marks measures that are at the top of the partial orders of two types of multiple-bug versions.

several smaller graphs so that measures that have similar accuracies would be grouped together, as shown in Figures 25, 26, 27, 28, 29, and 30. For each graph, we compare a number of association measures with Tarantula and Ochiai. Measures in Figure 25 shows measures that could localize more buggy versions as compare to Tarantula and Ochiai when less than 10% of code is inspected.

Measures that have similar performance with Tarantula are shown in Figure 26. Figures 27 and 28 show measures that could localize similar number of buggy version as compare to Tarantula and Ochiai when less than 10% of code is inspected. Measures that perform worse than Tarantula and Ochiai are shown in Figures 29 and 30. In this paper, we omit showing the curves for versions that contain five bugs because the curves are the similar with curves of all multiple-bug versions.

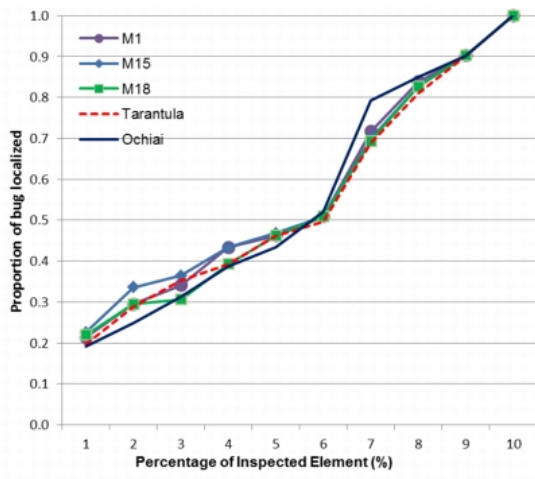


Figure 25. Comparing M_1 , M_{15} , M_{18} with Ochiai and Tarantula for all multiple-bugs versions

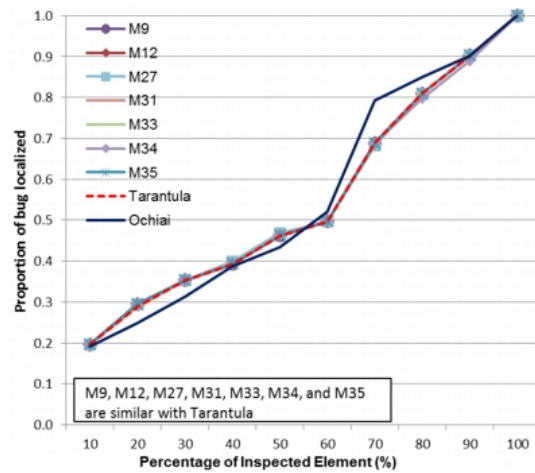


Figure 26. Comparing M_9 , M_{12} , M_{27} , M_{31} , M_{33} - M_{35} with Ochiai and Tarantula for all multiple-bug versions

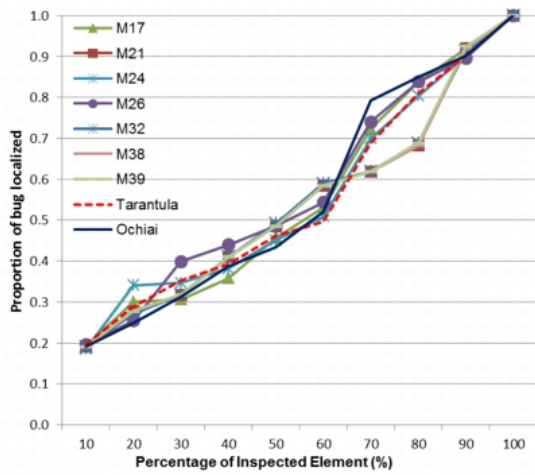


Figure 27. Comparing M_{17} , M_{21} , M_{24} , M_{26} , M_{32} , M_{38} , M_{39} with Ochiai and Tarantula for all multiple-bug versions

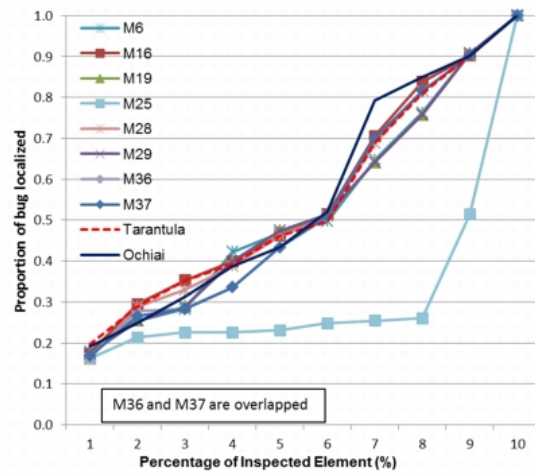


Figure 28. Comparing M_6 , M_{16} , M_{19} , M_{25} , M_{28} , M_{29} , M_{36} , M_{37} with Ochiai and Tarantula for all multiple-bug versions

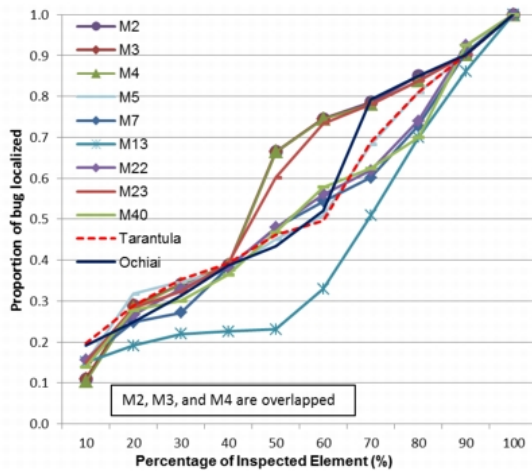


Figure 29. Comparing M_2 - M_5 , M_7 , M_{13} , M_{22} , M_{23} , M_{40} with Ochiai and Tarantula for all multi-bug programs

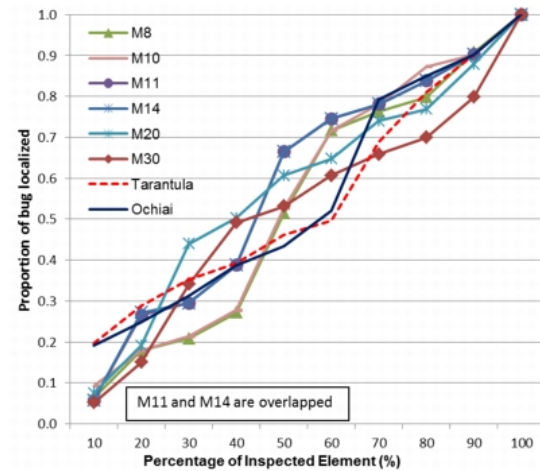


Figure 30. Comparing M_8 , M_{10} , M_{11} , M_{14} , M_{20} , M_{30} with Ochiai and Tarantula for all multi-bug programs

We also plot the curve showing the proportion of code that are investigated (x-axis) vs. the proportion of bugs localized (y-axis) for versions that contain two bugs, as shown in Figures 31, 32, 33, 34, 35, and 36. Measures shown in Figures 31, 32, and 33 shows measures that could localize more buggy versions as compare to Tarantula and Ochiai when less than 10% of code is inspected. Measures that have similar performance with Tarantula are shown in Figure 34. Figure 35 show measures that could localize similar number of buggy version as compare to Tarantula and Ochiai when less than 10% of code is inspected. Measures that perform worse than Tarantula and Ochiai are shown in Figure 36. We summarize the findings of this section in the answer for RQ6 in the next Section.

5.7. Discussion

In this section, we summarize the answers to the research questions mentioned in Section 1.

RQ1. We are interested to find if off-the-shelf association measures are powerful enough to locate bugs. Based on the mean accuracy values of the measures, it could be noted that the 40 association measures could help to find all the bugs when an average of 22% - 56% of the program elements are inspected. Fifty percent of the association measures are able to help find bugs by inspecting

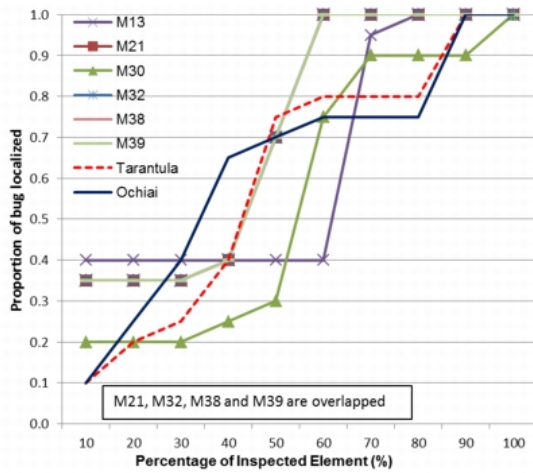


Figure 31. Comparing M_{13} , M_{21} , M_{30} , M_{32} , M_{38} , M_{39} with Ochiai and Tarantula for versions containing two bugs

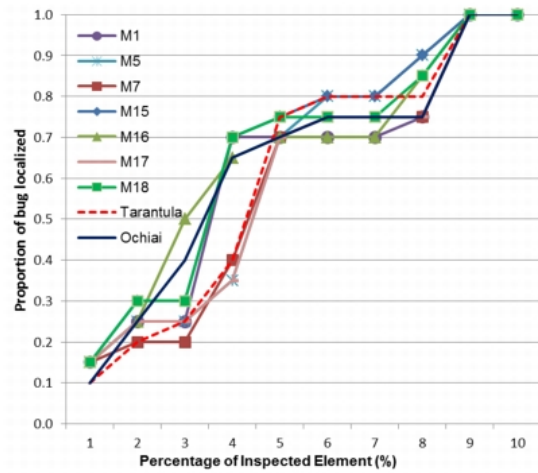


Figure 32. Comparing M_1 , M_5 , M_7 , M_{15} , M_{16} , M_{17} , M_{18} with Ochiai and Tarantula for versions containing two bugs

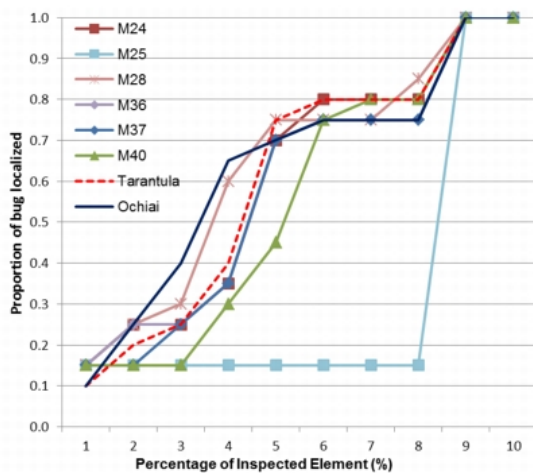


Figure 33. Comparing M_{24} , M_{25} , M_{28} , M_{36} , M_{37} , M_{40} with Ochiai and Tarantula for versions containing two bugs

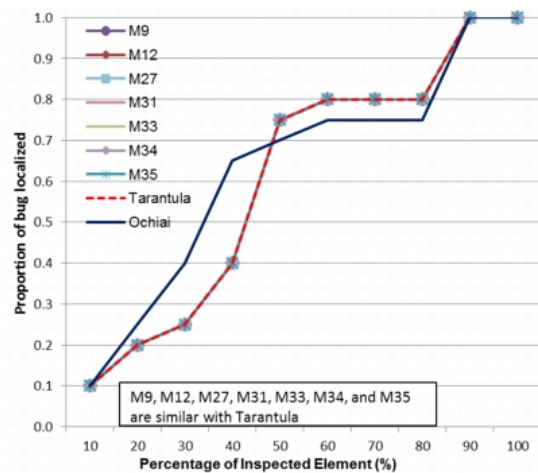


Figure 34. Comparing M_9 , M_{12} , M_{27} , M_{31} , M_{33} - M_{35} with Ochiai and Tarantula for versions containing two bugs

an average of 22-25% of elements, while Tarantula and Ochiai require debuggers to inspect approximately 23% and 25% of program elements respectively.

RQ2. Next, we are interested to find which association measures are better than others in localizing single bug programs. The answer to this research question is the partial order shown in Figure 2. At the top of the partial order there are 2 off-the-shelf association measures

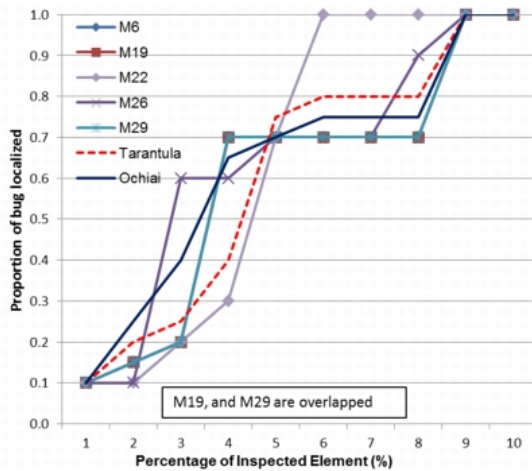


Figure 35. Comparing M_6 , M_{19} , M_{22} , M_{26} , and M_{29} with Ochiai and Tarantula for versions containing two bugs

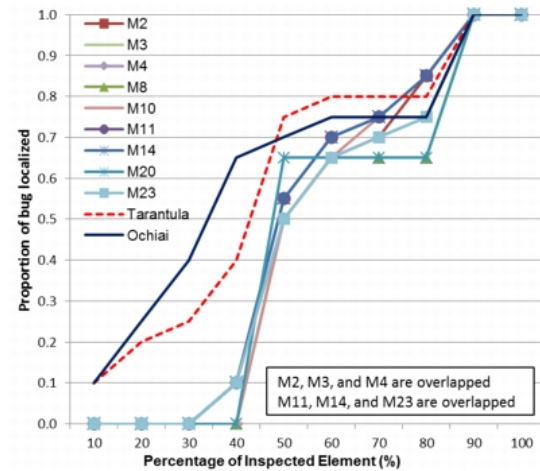


Figure 36. Comparing $M_2 - M_4$, M_8 , M_{10} , M_{11} , M_{14} , M_{20} , M_{23} with Ochiai and Tarantula for versions containing two bugs

namely: Klossgen (M_{18}), and Normalized Mutual Information (M_{26}) that perform comparably to Ochiai. They are statistically significantly better than the other off-the-shelf association measures and Tarantula. There are 12 other measures that perform statistically significantly better than Tarantula. These are: ϕ -coefficient (M_1), Added Value (M_{15}), Collective Strength (M_{16}), J-Measure (M_6), Information Gain (M_{19}), Two-way Support (M_{28}), Two-way Support Variation (M_{29}), Interestingness Weighting Dependency (M_{24}), Gini Index (M_7), Example and Counterexample Rate (M_{34}), Kappa (M_5), and Ochiai II (M_{40}).

RQ3. Finally, we would like to know the relative accuracy of the association measures versus those of well-known suspiciousness measures for fault localization. By applying statistical significance tests under 0.05 significance threshold, Klossgen (M_{18}) and Normalized Mutual Information (M_{26}) are comparable to Ochiai and are statistically significantly better than Tarantula. Based on the proportion of bugs localized, Klossgen (M_{18}) and Added Value (M_{15}) could localize more bugs than Tarantula and Ochiai by inspecting up to 10% of the program elements. They could localize 49% and 50% of the bugs respectively, while Tarantula and Ochiai could localize 40% and 47% of the bugs respectively.

RQ4. We find that most measures perform better for the C programs than the Java programs that we analyze. To evaluate the measures in terms of percentage of code inspected, we compute two partial orders of the 40 association measures, Tarantula, and Ochiai, by performing statistical significance tests. For the C programs, Ochiai and Klosgen (M_{18}) are the measures that are at the top of the partial orders (i.e., no measures are statistically significantly better than them). For the Java programs, a number of measures are at the top including Ochiai, Klosgen (M_{18}), Confidence (M_9), Interest (M_{12}), Added Value (M_{15}), Sebag (M_{31}), Example and Counterexample Rate (M_{34}), Zhang (M_{35}), Tarantula, Interestingness Weighting Dependency (M_{24}), and Normalized Mutual Information (M_{26}). In terms of proportion of bugs localized when up to 10% of code is inspected, for the C programs, Klosgen (M_{18}) and Added Value (M_{15}) outperform Ochiai and Tarantula. Tarantula and Ochiai could localize 44% and 55% of the bugs respectively, while these measures could localize 56% of the bugs. For the Java programs, Tarantula, Confidence (M_9), Interest (M_{12}), Added Value (M_{15}), Klosgen (M_{18}), Relative Risk (M_{23}), Loevinger (M_{30}), Normalized Mutual Information (M_{26}), Odd Multiplier (M_{33}), Example and Counterexample Rate (M_{34}), and Least Contradiction (M_{32}) could localize the most number of bugs (i.e., 26% of the bugs could be localized). Ochiai, on the other hand, could only localize 20% of the bugs.

RQ5. In terms of percentage of code inspected, again we compute several partial orders (one per bug category) by performing statistical significance tests. Tarantula is at the top of four partial orders. Ochiai, Information Gain (M_{19}), Normalized Mutual Information (M_{24}), and Two-way Support Variation (M_{29}) are at the top of six partial orders. Klosgen (M_{18}) is at the top of seven partial orders. In terms of proportion of bugs localized, we notice that Klosgen (M_{18}) could localize the most number of bugs for all categories as compared to other measures when up to 10% of the program elements are inspected. The categories of bugs that could be better localized by most of the measures are addition or removal of non-conditional statements (CH-CS) and change of loop predicates (LP-CC).

RQ6. When localizing multiple bugs in programs, all measures does not perform as good as localizing single bugs in the programs. The percentage of code inspected required to localize all the bugs, versions contain five bugs, and versions contain two bugs are 45% to 71%, 43% to 71%, and 34% to 74% respectively. Measures that are at the top of the partial order for localizing all multiple-bug versions are Odds ratio (M_2), Yule's Q (M_3), Added Value (M_{15}), and Coverage (M_{20}). The percentage of code inspected of the measures in localizing versions that contain different number of bugs inserted are different (i.e, two and five bugs). However, we notice that Added Value (M_{15}) and Normalized Mutual Information (M_{26}) consistently outperform Tarantula and Ochiai when localizing versions that contain two and five bugs. In addition, these measures also have good accuracies in localizing single bug programs. The overall mean of percentage of code inspected for Added Value (M_{15}) and Normalized Mutual Information (M_{26}) are 23.30% and 23.22% respectively, which are only slightly smaller then the smallest mean of percentage of code inspected to localized all single bug programs (22.42%). Also, Added Value (M_{15}) and Normalized Mutual Information (M_{26}) are in the top of partial order of 5 and 6 bug categories out of 7 categories.

Based on our results, we find that there is no single best measure for all cases. For localizing C and Java programs containing single bug, Klosgen (M_{18}) and Ochiai always outperform other measures. They are comparable. When we evaluate the effectiveness of the measures to localize different bug categories, there is no other measures that could outperform Klosgen (M_{18}) for all bug categories. Ochiai could outperform in 5 out of 7 bug categories. However, when localizing multiple bugs, there are several measures that outperform these measures. Added Value (M_{18}), Odds Ratio (M_2), Yule's Q (M_3), and Coverage (M_{20}) are the best measures. We notice that Added Value also has good accuracy in localizing single bug, even though it is not as good as Ochiai and Klosgen, but its accuracy is only slightly lower than Ochiai and Klosgen. It can be interesting future work to explore ways to compose various measures together so that the combined *meta-measure* may perform better for all cases than every individual measure.

5.8. Threats to Validity

Threats to construct validity refers to the suitability of our evaluation criteria. In this work, we use two criteria: percentage of program elements inspected to find all bugs, and proportion of bugs localized when at most a given percentage of program elements are inspected. We believe these two evaluation criteria reasonably measure the effectiveness of a fault localization approach. They have also been used before in prior studies on fault localization [4, 16].

Threats to internal validity include bias and human errors. The accuracy of a measure in localizing bugs is influenced by the granularity level considered during program instrumentation and trace generation (statement, basic block, or method levels). Different granularity levels may produce different accuracies since there would be different total numbers of elements which would affect the percentages of inspected elements. We choose to use block-hit spectra in our evaluation since it has a suitable balance between instrumentation costs and bug-reveal powers, and the focus of our study is to compare the effectiveness of different association measures on the *same* spectra. We hypothesize that the relative performance of different association measures on other spectra may remain the same as that on block-hit spectra, but it remains interesting future work for us to verify. Also, we manually instrument the C programs; we might miss instrumenting some blocks or add extraneous instrumentation code. For Java program, we automatically instrument the programs where there could be possible implementation errors. We manually assign the bugs into categories; there might be some errors in our assignments. In order to minimize such errors, we carefully checked the instrumented programs and assigned bug category labels.

Threats to external validity refers to the generalizability of our findings. We have tried to reduce this threat by considering a number of programs of various sizes written in two popular programming languages: C and Java. In the future, we would like to reduce this threat further by analyzing more programs written in various programming languages.

6. CONCLUSION

In this work, we investigate the effectiveness of a comprehensive number of association measures for fault localization. These measures gauge the strength of association between two variables expressible as a dichotomy matrix. We consider and compare 40 association measures with two well-known fault localization measures, namely Tarantula and Ochiai.

In terms of average percentage of code inspected, we find that Klosgen outperforms Ochiai and Tarantula. Many measures, including Normalized Mutual Information, ϕ -coefficient, Added Value, Collective Strength, J-Measure, Information Gain, Two-way Support, Two-way Support Variation, Interestingness Weighting Dependency, Gini Index, Example and Counterexample Rate, Kappa, and Ochiai II outperform Tarantula. The percentages of code inspected for different buggy program versions are different; such percentage values for each measure form a distribution, and we could employ statistical tests to compare the accuracy of different measures. We find that three measures, Klosgen, Normalized Mutual Information, and Ochiai can be statistically significantly better than other measures. In terms of proportion of bugs found when up to 10% of the code is inspected, Klosgen outperforms Ochiai. Also, Klosgen, Ochiai, Collective Strength, Added Value, ϕ -coefficient, Normalized Mutual Information, Two-way Support, Interestingness Weighting Dependency, J-Measure, Information Gain, Two-way Support Variation, Example and Counterexample Rate, Kappa, Confidence, Odd Multiplier, Sebag, Interest, Zhang, One-way Support, Gini Index, Jaccard, Sorensen-Dice, Anderberg, and Ochiai II outperform Tarantula. Thus, we can conclude that association measures are also promising to be used for fault localization.

We find that most measures perform better for the C programs than the Java programs that we analyze. For the C programs, in terms of proportions of bugs localized, Added Value and Normalized Mutual Information outperform Ochiai and Tarantula. For the Java programs, in terms of proportions of bugs localized, Tarantula, Confidence, Interest, Added Value, Klosgen, Relative Risk, Loevinger, Normalized Mutual Information, Odd Multiplier, Example and Counterexample Rate, and Least Contradiction measures, outperform the other measures and Ochiai.

We have also grouped the bugs into 7 categories and analyze the effectiveness of the measures to localize each category of bugs. The categories of bugs that could be better localized by most of the measures are addition or removal of non-conditional statements (CH-CS) and change of loop predicates (LP-CC). Klosgen is among the best measures for all bug categories.

The effectiveness of all measures in localizing multiple bugs in programs is not as good as localizing single bug in programs. The smallest percentage of code inspected required to localize the bugs is 45%. The measures that outperform other measures in localizing multiple-bug versions are Odds Ratio, Yule's Q, Added Value, and Coverage. We notice that Added Value is consistently outperform Ochiai and Tarantula for localizing both versions that contain two and five bugs.

In the future, we would like to integrate promising association measures for fault localization to popular IDEs and debugging tools such as Eclipse, Visual Studio.Net, etc.

Dataset. Our dataset and tool are made publicly available at:
<http://www.mysmu.edu/phdis2009/lucia.2009/jsme/Dataset.htm>.

REFERENCES

1. Tassey G. The economic impacts of inadequate infrastructure for software testing. *National Institute of Standards and Technology. Planning Report 02-3.2002* 2002; .
2. Beizer B. *Software Testing Techniques*. 2nd edn., International Thomson Computer Press: Boston, 1990.
3. Liblit B, Naik M, Zheng AX, Aiken A, Jordan MI. Scalable statistical bug isolation. *Proc. ACM SIGPLAN 2005 Int. Conf. Programming Language Design and Implementation (PLDI'05)*, 2005.
4. Jones J, Harrold M. Empirical evaluation of the tarantula automatic fault-localization technique. *Proc. of International Conference on Automated Software Engineering*, 2005.
5. Renieris M, Reiss S. Fault localization with nearest neighbor queries. *Proc. of Int. Conf. on Automated Software Engineering*, 2003; 141–154.
6. Manevich R, Sridharan M, Adams S, Das M, Yang Z. PSE: Explaining program failures via postmortem static analysis. *FSE*, 2004. URL citeseer.ist.psu.edu/manevich04pse.html.
7. Ko AJ, Myers BA. Debugging reinvented: asking and answering why and why not questions about program behavior. *International Conference on Software Engineering*, 2008.
8. Gupta N, He H, Zhang X, Gupta R. Locating faulty code using failure-inducing chops. *ASE*, 2005; 263–272.

9. Mayer W, Stumptner M. Model-Based Debugging – State of the Art And Future Challenges. *Electronic Notes in Theoretical Computer Science (ENTCS)* 2007; **174**(4).
10. Zeller A. *Why Programs Fail: A Guide to Systematic Debugging*. 2nd edn., Morgan Kaufmann, 2009.
11. Qi D, Roychoudhury A, Liang Z, Vaswani K. Darwin: An approach for debugging evolving programs. *ESEC / SIGSOFT FSE '09*, ACM: Amsterdam, The Netherlands, 2009.
12. CLiu, Yan X, Fei L, Han J, Midkiff S. Sober: Statistical model-based bug localization. *Proc. of Joint Meeting of European Software Engineering Conference and ACM SIGSOFT Symposium on the Foundations of Software Engineering*, Lisbon, Portugal, 2005.
13. Abreu R. Spectrum-based fault localization in embedded software. PhD Thesis, Delft University of Technology 2009.
14. Chilimbi TM, Liblit B, Mehra K, Nori AV, Vaswani K. Holmes: Effective statistical debugging via efficient path profiling. *2009 IEEE 31st International Conference on Software Engineering*, IEEE Computer Society, 2009; 34–44, doi:10.1109/ICSE.2009.5070506.
15. Jones J, Harrold M, Stasko J. Visualization of test information to assist fault localization. *Proc. of International Conference on Software Engineering*, Orlando, Florida, 2002; 467–477.
16. Abreu R, Zoetewij P, van Gemund AJC. On the Accuracy of Spectrum-based Fault Localization. *Mutation Testing: Academic and Industrial Conference Practice and Research Techniques (TAICPART-MUTATION)*, 2007.
17. Abreu R, Zoetewij P, van Gemund AJ. Spectrum-Based Multiple Fault Localization. *IEEE/ACM International Conference on Automated Software Engineering*, Auckland, New Zealand, 2009.
18. Agresti A. *An Introduction to Categorical Data Analysis*. John Wiley & Sons, 1996.
19. Yule GU. On the association of attributes in statistics. *Philosophical Transactions of the Royal Society* 1900; **A194**:257–319.
20. Yule GU. On the methods of measuring association between two attributes. *Journal of the Royal Statistical Society* 1912; **75**:579–642.
21. Abreu R, Zoetewij P, Golsteijn R, van Gemund AJC. A practical evaluation on spectrum-based fault localization. *The Journal of System and Software* 2009; **82**:1780–1792.
22. Siemens, Harrold M, Rothermel G. *Aristotle Analysis System – Siemens Programs, HR Variants*. <http://www.cc.gatech.edu/aristotle/Tools/subjects/>.
23. Do H, Elbaum SG, Rothermel G. Supporting controlled experimentation with testing techniques: An infrastructure and its potential impact. *Empirical Software Engineering: An International Journal* 2005; **10**(4):405–435.
24. Liblit B, Aiken A, Zheng AX, Jordan MI. Bug isolation via remote program sampling. *Proc. ACM SIGPLAN 2003 Conf. Programming Language Design and Implementation (PLDI 2003)*, San Diego, CA, 2003; 141–154.
25. Santelices R, Jones J, Yu Y, Harrold M. Lightweight fault-localization using multiple coverage types. *Proc. of ICSE*, 2009.
26. Artzi S, Dolby J, Tip F, Pistoia M. Directed test generation for effective fault localization. *ISSTA'10*, 2010.
27. Artzi S, Dolby J, Tip F, Pistoia M. Practical fault localization for dynamic web applications. *ICSE'10*, 2010.

28. Bandyopadhyay A, Ghosh S. On the effectiveness of the tarantula fault localization technique for different fault classes. *HASE'11*, 2011; 317–324.
29. Pan K, Kim S, JE JW. Toward understanding bug fix patterns. *Empirical Software Engineering* 2009; **14**:286–315.
30. BJiang, WKChan, THTse. On practical adequate test suites for integrated test case prioritization and fault localization. *Q SIC'11*, 2011; 21–30.
31. Zhang X, Gupta N, Gupta R. Locating faults through automated predicate switching. *International Conference on Software Engineering*, 2006.
32. Sterling CD, Olsson RA. Automated bug isolation via program chipping. *Software: Practice and Experience (SP&E)* August 2007; **37**(10):1061–1086. John Wiley & Sons, Inc.
33. Jeffrey D, Gupta N, Gupta R. Fault localization using value replacement. *International Symposium on Software Testing and Analysis*, 2008.
34. Feldman A, van Gemund A. A two-step hierarchical algorithm for model-based diagnosis. *Proceedings of the 21st National Conference on On Artificial Intelligence*, AAAI Press: Boston, Massachusetts, 2006; 827–833.
35. Mayer W, Stumptner M. Evaluating models for model-based debugging. *ASE*, 2008; 128–137.
36. Mayer W, Stumptner M. Abstract interpretation of programs for model-based debugging. *IJCAI*, 2007; 471–476.
37. Mayer W, Abreu R, Stumptner M, van Gemund A. Prioritising model-based debugging diagnostic reports. *Proceedings of the International Workshop on Principles of Diagnosis (DX)*, 2009.
38. Liblit B, Aiken A. Building a better backtrace: Techniques for postmortem program analysis. *Technical Report CSD-02-1203*, UC Berkeley 2002.
39. Tallam S, Tian C, Gupta R. Dynamic slicing of multithreaded programs for race detection. *Proc. of ICSM*, 2008.
40. Parnin C, Orso A. Are automated debugging techniques actually helping programmers? *ISSTA*, 2011; 199–209.
41. Altman RB, Klein TE. Challenges for biomedical informatics and pharmacogenomics. *Annu. Rev. Pharmacol. Toxicol.* 2002; **42**:113–133.
42. Healey JF. *Statistics: A Tool for Social Research*. 8th edn., Wadsworth Publishing, 2008. URL <http://www.amazon.com/Statistics-Research-Joseph-F-Healey/dp/0495096555>.
43. Agrawal R, Srikant R. Fast algorithms for mining association rules. *Proc. of Int. Conf. on Very Large Data Bases*, 1994.
44. Tan PN, Kumar V, Srivastava J. Selecting the right interestingness measure for association patterns. *Proc. 2002 ACM SIGKDD Int. Conf. Knowledge Discovery in Databases (KDD'02)*, Edmonton, Canada, 2002; 32–41.
45. Geng L, Hamilton H. Interestingness measures for data mining: A survey. *ACM Computing Surveys*, 2006.
46. Reps T, Ball T, Das M, Larus J. The use of program profiling for software maintenance with applications to the year 2000 problem. *ESEC/FSE*, 1997.
47. Harrold M, Rothermel G, Sayre K, Wu R, Yi L. An empirical investigation of the relationship between spectra differences and regression faults. *Software Testing, Verification and Reliability* 2000; **10**(3):171–194.
48. Cohen J. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement* 1960; **20**(1):37–46.

49. Smyth P, Goodman R. An information theoretic approach to rule induction from databases. *IEEE Trans. Knowledge and Data Eng.* 1992; **4**(4):301–316.
50. Gini C. Variability and mutability, contribution to the study of statistical distributions and relations. *Studi Economico-Giuridici della R. Università de Cagliari* 1912; **3**(part 2)(i-iii):3–159.
51. Clark P, Boswell R. Rule induction with cn2: Some recent improvements. *In Machine Learning - EWSL-91*, 1991; 151163.
52. Brin S, Motwani R, Ullman JD, Tsur S. Dynamic itemset counting and implication rules for market basket analysis. *Proceedings of 1997 ACM SIGMOD Int. Conf. Management of Data (SIGMOD'97)*, Tucson, AZ, 1997; 255–264.
53. Piatesky-Shapiro G. Discovery, analysis, and presentation of strong rules. *Knowledge Discovery in Databases*, 1991.
54. Shortliffe E, Buchanan B. A model of inexact reasoning in medicine. *Mathematical Biosciences* 1975; **23**:351379.
55. Aggarwal C, Yu PS. A new framework for itemset generation. *Symposium on Principles of Database Systems (PODS)*, 1998.
56. Hand DJ, Mannila H, Smyth P. *Principles of Data Mining*. MIT Press, 2001.
57. Klosgen W. Explora: A multipattern and multistrategy discovery assistant. *Advances in Knowledge Discovery and Data Mining*, 1996.
58. Quinlan J. *C4.5: Programs for Machine Learning*. Morgan Kaufmann, 1993.
59. Cheng H, Lo D, Zhou Y, Wang X, Yan X. Identifying bug signatures using discriminative graph mining. *ISSTA*, 2009.
60. Ohsaki M, Kitaguchi S, Okamoto K, Yokoi H, Yamaguchi T. Evaluation of rule interestingness measures with a clinical dataset on hepatitis. *Proc. of the 8th European Conference on Principles of Data Mining and Knowledge Discovery (PKDD 2004)*, 2004; 362–373.
61. Goodman L, Kruskal W. Measures of association for cross classifications. *J. Amer. Statistic Association* 1954; **49**:732–764.
62. Dice LR. Measures of the amount of ecologic association between species. *Ecology*, 1945.
63. Sorensen T. A method of establishing groups of equal amplitude in plant sociology based on similarity of species content and its application to analyses of the vegetation on danish commons. *Vidensk Selsk Biol Skr*, 1948.
64. Anderberg M. *Clustering Analysis for Applications*. London, Academic Press, 1973.
65. Sokal RR, Michener C. A statistical method for evaluating systematic relationships. *Univ Kans Sci Bull*, 1958.
66. Rogers J, Tanimoto TT. A computer program for classifying plants. *Science*, 1960.
67. Ochiai A. Zoogeographic studies on the soleoid fishes found in japan and its neighbouring regions. *Bull Jnp Soc Sci Fish*, 1957.
68. Hutchins M, Foster H, Goradia T, Ostrand T. Experiments of the effectiveness of dataflow- and controlflow-based test adequacy criteria. *Proc. of ICSE*, 1994; 191–200.
69. Wilcoxon F. Individual comparisons by ranking methods. *Biometrics Bulletin available: <http://www.jstor.org/stable/3001968>* 1945; **1**:80–83.