

# SURE: A Visualized Failure Indexing Approach using Program Memory Spectrum

YI SONG\* and XIHAO ZHANG\*, School of Computer Science, Wuhan University, China  
 XIAOYUAN XIE†, School of Computer Science, Wuhan University, China  
 SONGQIANG CHEN, The Hong Kong University of Science and Technology, China  
 QUANMING LIU, School of Computer Science, Wuhan University, China  
 RUIZHI GAO, Sonos Inc., USA

Failure indexing is a longstanding crux in software testing and debugging, the goal of which is to automatically divide failures (e.g., failed test cases) into distinct groups according to the culprit root causes, as such multiple faults residing in a faulty program can be isolated and thus be handled independently and simultaneously. The community of failure indexing has long been plagued by two challenges: 1) The effectiveness of division is still far from promising. Specifically, existing failure indexing techniques only employ a limited source of software run-time data, for example, code coverage, to be failure proximity and further divide them, which typically delivers unsatisfactory results. 2) The outcome can be hardly comprehensible. Specifically, a developer who receives the division result is just aware of how all failures are divided, without knowing why they should be divided the way they are. This leads to difficulties for developers to be convinced by the division result, which in turn affects the adoption of the results. To tackle these two problems, in this paper, we propose SURE, a viSualized failuRe indExing approach using the program memory spectrum. We first collect the run-time memory information (i.e., variables' names and values, as well as the depth of the stack frame) at several preset breakpoints during the execution of a failed test case, and transform the gathered memory information into a human-friendly image (called program memory spectrum, PMS). Then, any pair of PMS images that serve as proxies for two failures is fed to a trained Siamese convolutional neural network, to predict the likelihood of them being triggered by the same fault. Last, a clustering algorithm is adopted to divide all failures based on the mentioned likelihood. In the experiments, we use 30% of the simulated faults to train the neural network, and use 70% of the simulated faults as well as real-world faults to test. Results demonstrate the effectiveness of SURE: It achieves 101.20% and 41.38% improvements in faults number estimation, as well as 105.20% and 35.53% improvements in clustering, compared with the state-of-the-art technique in this field, in simulated and real-world environments, respectively. Moreover, we carry out a human study to quantitatively evaluate the comprehensibility of PMS, revealing that this novel type of representation can help developers better comprehend failure indexing results.

CCS Concepts: • **Software and its engineering** → **Software testing and debugging**.

\*Both authors contributed equally to this research.

†Corresponding author.

Authors' addresses: Yi Song, [yisong@whu.edu.cn](mailto:yisong@whu.edu.cn); Xihao Zhang, [zhangxihao@whu.edu.cn](mailto:zhangxihao@whu.edu.cn), School of Computer Science, Wuhan University, No. 299 Bayi Road, Wuchang District, Wuhan, Hubei, China, 430072; Xiaoyuan Xie, School of Computer Science, Wuhan University, No. 299 Bayi Road, Wuchang District, Wuhan, Hubei, China, [xxie@whu.edu.cn](mailto:xxie@whu.edu.cn); Songqiang Chen, The Hong Kong University of Science and Technology, Clear Water Bay, Hong Kong, China, [i9s.chen@connect.ust.hk](mailto:i9s.chen@connect.ust.hk); Quanming Liu, School of Computer Science, Wuhan University, No. 299 Bayi Road, Wuchang District, Wuhan, Hubei, China, [liuquanming@whu.edu.cn](mailto:liuquanming@whu.edu.cn); Ruizhi Gao, Sonos Inc., 2 Ave de Lafayette, Boston, MA, USA, [youtianzui.nju@gmail.com](mailto:youtianzui.nju@gmail.com).

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

© 2018 Association for Computing Machinery.

0004-5411/2018/8-ART111 \$15.00

<https://doi.org/XXXXXXX.XXXXXXX>

Additional Key Words and Phrases: failure proximity, program memory, siamese learning, failure indexing, parallel debugging

#### ACM Reference Format:

Yi Song, Xihao Zhang, Xiaoyuan Xie, Songqiang Chen, Quanming Liu, and Ruizhi Gao. 2018. SURE: A Visualized Failure Indexing Approach using Program Memory Spectrum. *J. ACM* 37, 4, Article 111 (August 2018), 37 pages. <https://doi.org/XXXXXXX.XXXXXXX>

## 1 INTRODUCTION

Due to the increments in the scale and complexity of modern software systems, faulty programs typically contain more than one fault [20, 25, 38, 78]. The co-existence of multiple faults has been shown to be harmful for fault localization tasks, that is, the existence of two or more faults will cause the localization of one or more of these faults to become harder [12, 20]. A mainstream way to tackle the challenge of multi-fault localization is *parallel debugging*<sup>1</sup>, in which a multi-fault localization task is decomposed into several sub-single-fault localization tasks by dividing<sup>2</sup> all failures<sup>3</sup> into distinct groups according to the culprit faults [18, 26, 28, 61, 71]. Such a division process is generally referred to as *failure indexing*, the goal of which is to isolate multiple faults into independent environments through the reorganization of failed test cases. Failure indexing has been proven to be diversely useful. For example, it can effectively alleviate the negative impact of multi-fault co-existence on fault localization by isolating faults. And for another example, the derived sub-tasks can be handled by different developers simultaneously, thus improving the efficiency of debugging [102].

For its promise, there are numerous works focusing on the topic of failure indexing [26, 33, 69]. Despite the fact that failure indexing has been integrated into many approaches or tools, it is plagued by two longstanding challenges:

- **The effectiveness of failure indexing is far from promising.** As mentioned above, the mission of failure indexing is to isolate multiple faults into separate environments by dividing failed test cases. It is obvious that an effective failure indexing process should satisfy two points: 1) The number of the derived fault-focused groups should be the same as the number of faults (i.e., *correct faults number estimation*), and 2) Failures in the same group should be triggered by the same fault, and conversely, failures in different groups should have different root causes (i.e., *promising clustering*). Many pioneering studies have pointed out that these two goals are difficult to achieve, because it is not easy to design an effective *failure proximity* (comprising the fingerprinting function and the distance metric, we will introduce them in detail in Section 2), which is the core of failure indexing [55]. The most prevalent and advanced failure proximities to date are the code coverage (CC)-based one and the statistical debugging (SD)-based one [33, 53]. However, they both rely on coverage information, which can be easily trapped and thus delivering unsatisfactory results in practice (we will further describe this point in Section 2.2 and Section 3).
- **The outcome of failure indexing is hard to comprehend.** Given a collection of failed test cases, a failure indexing process will index them to their root causes and thus produce several fault-focused groups. However, developers who receive the result are only aware of how failures are divided, without knowing why they should be divided in the current manner. It is recognized that software engineering is a human-centric discipline [68], the results provided by automated debugging techniques still require being comprehended by the person who

<sup>1</sup>In parallel debugging, people generally consider faults that do not interfere with one another.

<sup>2</sup>In the current field of parallel debugging, *clustering* is typically utilized for such *division*. Thus we use these two terms interchangeably in this paper.

<sup>3</sup>In dynamic testing, *failures* are typically *failed test cases*. In this paper, we use these two terms interchangeably.

performs the following task [12, 91]. In fact, nearly thirty years ago, the comprehensibility of software debugging tasks has been considered by researchers. Specifically, to make developers better comprehend fault localization results, Agrawal et al. and Jones et al. used color to visually reveal the linkage between program statements and failed/passed runs. These works have been demonstrated to be effective in helping a user localize faults [2, 38]. Parnin and Orso also realized this problem and carried out a study, explicitly indicating that human factors played a critical role in software debugging tasks [59]. Moreover, Xie et al. suggested integrating comprehension assistance into automated debugging, because the results delivered by automated debugging techniques still require being comprehended by human developers who perform the following jobs [91]. In 2023, Souza et al. continue highlighting the importance of comprehension in software debugging tasks. They employed Jaguar [64], an SBFL tool that visualizes fault localization results by coloring the suspicious elements from red (high) to green (low), and investigated to which extent such visualization provided hints for human developers to comprehend faulty behaviors [3]. The aforementioned representative works have garnered more than 2,700 citations. That is to say, comprehensibility has been recognized as an important topic in the field of software debugging. Seeing that failure indexing is an important class of tasks in software debugging, we can claim that **comprehension matters in failure indexing**: only when human developers comprehend the machine-produced failure indexing result, they will be convinced by the result and would like to further apply this result. Therefore, providing the failure indexing result and providing why the result is obtained are both important. However, to the best of our knowledge, existing studies only consider the former without focusing on the comprehensibility of failure indexing results.

These two challenges are worth addressing because they are in line with two problems in practice, respectively: First, the effectiveness of failure indexing determines *how good failure indexing will be*, and second, the comprehensibility of the result determines *whether the failure indexing result will be adopted by human developers*. In this paper, we propose SURE, a viSualized failuRe indEXing approach based on the program memory spectrum for tackling the aforementioned challenges:

- **To further improve the effectiveness of failure indexing**, we propose a novel type of failure proximity, namely, the Program Memory (PM)-based failure proximity. The shortcoming of using coverage to represent failures can be partly described by the PIE (Propagation, Infection, Execution) model [77], a classical theory in software testing and debugging. Specifically, the PIE model has demonstrated that a failure can be detected only if the fault<sup>4</sup> infects the program's internal state. However, with the coverage information only, it is very hard to explore the faulty internal state in depth during the program execution, because **being covered is a necessary but not sufficient condition for triggering a failure**. Thus, coverage cannot extract the signature of failures in deep insight. Based on this intuition, we conjecture that program internal dataflows (program memory in our context) can be a finer-grained failure representation when using coverage gives rise to unsatisfactory effectiveness, because program internal dataflows can effectively embody program internal states<sup>5</sup>. For the fingerprinting function, the PM-based failure proximity mines and utilizes the run-time memory information (i.e., variables' names, variables' values, and the depth of stack frame, **deeper insight to explore programs' internal state than coverage**) collected during the execution of a failed test case to represent it. Specifically, we first employ Spectrum-Based Fault Localization (SBFL) techniques to determine several highly-suspicious program statements, and set these statements as breakpoints. And then, we collect the run-time memory

<sup>4</sup>In this paper, we follow the practice of general fault localization to consider the non-omission fault.

<sup>5</sup>This conjecture has been verified by the experiments in Section 6.

information (i.e., variables' names and values, as well as the depth of the stack frame) at the preset breakpoints during executing a failed test case. We further reorganize the collected run-time memory information into Program Memory Spectrum (PMS), which is in the form of an image, to serve as the proxy for a failed test case. And for the distance metric, we train a Siamese convolutional neural network to predict the distance between a pair of PMS images (i.e., a pair of failed test cases), which reflects the likelihood that they are triggered by the same fault.

- **To make failure indexing results comprehensible to human developers**, we borrow the idea of Agrawal et al.'s work and Jones et al.'s work to make the failure indexing result visualized. In their work, they used color to visually map program statements in test suite executions, and help developers pinpoint faulty statements intuitively [2, 38]. In this paper, we design a novel algorithm to reorganize the collected run-time memory information (which is not easy for human developers to recognize) into PMS images (which are in a human-friendly form). Specifically, for a set of run-time memory information collected during running a failed test case, we first convert variables' names and values into the numeric form according to the Ascii code of characters, and then regard variables' names, variables' values, and the depth of the stack frame as the three channels of RGB (we will further introduce RGB in Section 4.3), to represent this failure in the form of PMS. As such, each failed test case is represented as a visualized image, so human developers can be easily aware of why all failures are divided the way they are when they receive the failure indexing result.

In the experiments, we obtain four C projects, Flex, Grep, Gzip, and Sed from the Software Infrastructure Repository (SIR) [21], and employ a mutation tool to inject simulated faults into the projects to generate faulty versions that contain one, two, three, four or five bugs (also referred to as 1-bug, 2-bug, 3-bug, 4-bug, and 5-bug faulty versions, respectively). Besides, we also use five Java projects, Chart, Closure, Lang, Math, and Time from Defects4J [39], as our benchmarks, and gather 1 ~ 5-bug real-world faulty versions according to the test case transplantation strategy proposed by An et al. [5]. In total, we generate 1,000 C simulated faulty versions and gather 100 Java real-world faulty versions. We use only 30% of the simulated faulty versions to train, and use 70% of the simulated faulty versions as well as real-world faulty versions to test. Experimental results significantly demonstrate the promise of the proposed approach: SURE can achieve 101.20% and 41.38% improvements in faults number estimation, as well as 105.20% and 35.53% improvements in clustering, compared with the state-of-the-art technique in the field, in simulated and real-world environments, respectively.

Moreover, to quantitatively measure to which extent the proposed visualized representation of failed test cases (i.e., PMS images) can be comprehensible to human developers, we conduct a human study involving 15 graduate students in Computer Science. The results show that compared with the most advanced and prevalent failure representation to date (i.e., the ranking-based and coverage-based ones), the representation of PMS images is more comprehensible and cost-effective: The comprehensibility of PMS images is 105.26% higher than the two baselines, while the time cost is only 15.35% and 14.87% of the two baselines, respectively.

The main contributions of this paper are as follows:

- (1) **A novel failure indexing approach.** We propose a visualized failure indexing approach, SURE, which utilizes the program memory spectrum to represent failures and uses Siamese convolutional neural networks to measure the distance between failures.
- (2) **A visualized representation of failed test cases.** The program memory spectrum is a visualized representation of failed test cases, which is in a human-friendly form (i.e., images) and thus can make failure indexing outcomes more comprehensible for human developers.

- (3) **A comprehensive evaluation.** We use both simulated and real-world faults in the experiments for more robust evaluation. It is worth mentioning that the training phase is performed on a small set of simulated faults while the test is conducted on a larger set of simulated faults and real-world faults. We also carry out a human study to quantitatively demonstrate the comprehensibility of the proposed PMS.

The remainder of this paper is organized as follows: Section 2 introduces the background knowledge and emphasizes the motivation of this paper. Section 3 provides a motivating example. Section 4 describes the technique details of SURE, and gives a running example to facilitate comprehension. Section 5 lists the research questions, datasets, evaluation metrics, etc. Section 6 analyzes the experimental results. Section 7 discusses some interesting topics. Section 8 is about the threats to validity. Section 9 reports related works. Conclusions and directions for future work are proposed in Section 10.

## 2 BACKGROUND

### 2.1 Failure Indexing

Parallel debugging is a well-recognized strategy for multi-fault localization, one of the most tedious and time-consuming problems in software testing and debugging [17, 42, 79, 85, 101]. The core of parallel debugging is to correctly divide all failures into distinct groups according to their root causes, i.e., failure indexing, thus multiple faults can be localized in isolated environments. It is obvious that failure indexing will directly determine the effectiveness of parallel debugging: the more promising failure indexing is, the better parallel debugging [36, 51, 87].

A collection of research has pointed out that failure indexing is indeed a difficult task. Failure indexing typically involves three components: A fingerprinting function that represents failures in a mathematical and structured form, a distance metric that calculates the distance between the proxies for failures, and a clustering algorithm that divides all failures into several groups based on the calculated distance information. Previous studies have made it clear that there is no clustering technique that is universally applicable in uncovering the variety of structures present in multidimensional data sets [34]. Thus, the most essential factors of failure indexing are the fingerprinting function and the distance metric (these two components are called *failure proximity*).

### 2.2 Fingerprinting Function

The directly available information regarding failed test cases to be indexed is twofold: the input data and the testing result (i.e., *failed*), which is too limited to effectively differentiate them. To tackle this problem, researchers typically design fingerprinting functions to mine dynamic information at run-time, so as to further extract the signature of failed test cases and use such data to represent failed test cases [16, 53, 62].

Among off-the-shelf fingerprinting functions, the code coverage-based strategy is the most prevalent currently [19, 31, 33], while the statistical debugging-based strategy is the most sophisticated and advanced [26, 36, 50, 69]. Here we use a faulty program containing  $l$  executable statements as an example. The code coverage (CC)-based fingerprinting function extracts the signature of failures from the execution path of failed test cases. It will represent failed test case  $f$  as a numerical vector of length  $l$ . There are two variants of the CC-based fingerprinting function:  $Cov_{hit}$  and  $Cov_{count}$ . As for the former, the  $i^{\text{th}}$  element of the numerical vector that represents  $f$  will be set to 1 if  $f$  covers the  $i^{\text{th}}$  statement  $s_i$ , and 0 otherwise. And as for the latter, the  $i^{\text{th}}$  element of the numerical vector that represents  $f$  will be set to the execution frequency of  $s_i$  if  $f$  covers  $s_i$ , and 0 otherwise. The statistical debugging (SD)-based fingerprinting function extracts the signature of failures from the suspiciousness ranking list of program statements. Specifically, a failed test case  $f$  is first combined

with passed test cases  $T_S$ , then the test suite  $f \cup T_S$  is executed on the faulty program and the coverage is collected simultaneously. Later, a spectrum-based fault localization (SBFL) technique is employed to calculate the suspiciousness value of each program statement based on coverage information [57, 90, 97], and produce a ranking list in which  $l$  program statements are descendingly ordered by their suspiciousness. As such,  $f$  can be represented as this ranking list of length  $l$ .

Despite the promise of the CC and the SD-based strategies, they share a common bottleneck: only relying on program coverage. Specifically, the intuition of the CC-based strategy is that failures triggered by the same fault should also have the same code coverage, and vice versa [98]. Thus, it creates numeric vectors that reflect the coverage information during the execution to represent failed test cases. The intuition of the SD-based strategy is that failures triggered by the same fault should target the same fault location, and vice versa [53, 55]. Thus, it incorporates SBFL techniques to further analyze raw coverage information, to convert it to a suspiciousness ranking list that suggests the fault location and thus can represent failed test cases. Obviously, both CC and SD purely rely on coverage information, which cannot work well if the coverage of the failures triggered by different faults is the same (we will exemplify such a situation in Section 3). That is to say, **even the most prevalent and advanced fingerprinting functions to date are still not able to deliver satisfactory results.**

### 2.3 Distance Metric

In failure proximity, the design of the distance metric is not an independent task, because it is tightly related to the form of the fingerprinting function. For example, the CC-based strategy represents a failed test case as a numerical vector and thus typically uses the Euclidean distance metric [76], because the Euclidean distance metric is a suitable and cheap way to handle the similarity measurement between such numerical vectors [33]. For another example, the SD-based strategy represents a failed test case as a ranking list and thus typically uses the Kendall tau distance metric [43], because the Kendall tau distance metric can properly measure how similar two ranking lists are by counting the number of pairwise disagreements between them [26]. From these, we can see that the distance metric is inseparable from the fingerprinting function: **only by designing a well-tailored distance metric can the signature of failures extracted by the fingerprinting function be fully exploited.**

## 3 MOTIVATING EXAMPLE

We use a motivating example in Table 1 to reveal the shortcoming of the CC-based and SD-based failure proximities mentioned above. This toy program containing 17 statements ( $s_1, s_2, \dots, s_{17}$ ) is used to identify and replace certain words in the input string, and then output the modified string and the log message. Specifically, if an input string contains “wordNone” or “wordNtwo”, these two words will be replaced with “\*1\*” and “\*2\*”, respectively. The log message records the operation of the program, for example, “wordNone recognized”, “wordNtwo recognized”, “both pattern recognized”, and “pass”. In this program, statements  $s_8$  and  $s_{16}$  each contain a fault.

Given a test suite containing 12 test cases:  $t_1 = \text{“speak wordNone”}$ ,  $t_2 = \text{“wordNone”}$ ,  $t_3 = \text{“word-Nonecontained”}$ ,  $t_4 = \text{“wwwwordNoneeee”}$ ,  $t_5 = \text{“has wordNtwo”}$ ,  $t_6 = \text{“wordNtwo”}$ ,  $t_7 = \text{“”}$ ,  $t_8 = \text{“mid*1*le”}$ ,  $t_9 = \text{“*1*2*”}$ ,  $t_{10} = \text{“a normal sentence”}$ ,  $t_{11} = \text{“wordnonewordNtw”}$ , and  $t_{12} = \text{“word-None and wordNtwo”}$ . The coverage information of them is given in Table 1, where a dot in the cell  $[t_i, s_j]$  indicates that  $t_i$  covers  $s_j$  ( $i = 1, 2, \dots, 12$ , and  $j = 1, 2, \dots, 17$ ). Among these test cases,  $t_1 \sim t_6$  are failed test cases due to the inconsistency between the actual output and the expected output ( $t_1 \sim t_6$  can be referred to as  $f_1 \sim f_6$ , respectively). More concretely,  $f_1 \sim f_4$  are triggered by *Fault*<sub>1</sub>, while  $f_5$  and  $f_6$  are triggered by *Fault*<sub>2</sub>. The remaining six test cases, i.e.,  $t_7 \sim t_{12}$ , pass the test (they can be referred to as  $T_S$ ).

Table 1. A motivating example

S	Program	$t_1$	$t_2$	$t_3$	$t_4$	$t_5$	$t_6$	$t_7$	$t_8$	$t_9$	$t_{10}$	$t_{11}$	$t_{12}$
$s_1$	public static String process(String s){	•	•	•	•	•	•	•	•	•	•	•	•
$s_2$	if(s.contains("**1**")    s.contains("**2**")){	•	•	•	•	•	•	•	•	•	•	•	•
$s_3$	return "";}								•	•			
$s_4$	int sign = 0;	•	•	•	•	•	•	•			•	•	•
$s_5$	int sum_1 = 0;	•	•	•	•	•	•	•			•	•	•
$s_6$	sum_1 = s.contains("wordNone") ? 1 : 0;	•	•	•	•	•	•	•			•	•	•
$s_7$	sign += sum_1;	•	•	•	•	•	•	•			•	•	•
$s_8$	s = s.replaceAll("wordNone", "?1?");	•	•	•	•	•	•	•			•	•	•
	//Fault1: "?1?" should be "**1**"												
$s_9$	int sum_2 = 0;	•	•	•	•	•	•	•			•	•	•
$s_{10}$	sum_2 = s.contains("wordNtwo") ? 2 : 0;	•	•	•	•	•	•	•			•	•	•
$s_{11}$	sign += sum_2;	•	•	•	•	•	•	•			•	•	•
$s_{12}$	s = s.replaceAll("wordNtwo", "**2**");	•	•	•	•	•	•	•			•	•	•
$s_{13}$	if(sign == 3){	•	•	•	•	•	•	•			•	•	•
$s_{14}$	return "both pattern recognized";}												•
	String msg = sign == 1 ?												
$s_{15}$	"wordNone recognized" : "pass";	•	•	•	•	•	•	•				•	•
	msg = sign > 2 ?												
$s_{16}$	"wordNtwo recognized" : msg;	•	•	•	•	•	•	•				•	•
	//Fault2: ">2" should be "==2"												
$s_{17}$	return s + "/" + msg;}	•	•	•	•	•	•	•				•	•
	<b>Result</b>	<b>F</b>	<b>F</b>	<b>F</b>	<b>F</b>	<b>F</b>	<b>F</b>	<b>S</b>	<b>S</b>	<b>S</b>	<b>S</b>	<b>S</b>	<b>S</b>

Table 2. The representation for  $f_2 \sim f_6$ 

Failure Proximity	Representation
The CC-based	[1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1]
The SD-based	[1, 1, 3, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 15, 15, 15]

As mentioned in Section 2.2, the CC-based failure proximity represents a failure as a code coverage vector, i.e., the corresponding column in Table 1. For example,  $f_1$  covers all program statements except  $s_3$  and  $s_{14}$ , thus, as for  $Cov_{hit}$ ,  $f_1$  can be represented as a binary vector of length 17: [1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1]. By observing Table 1, we can find that the CC-based representation for the other five failed test cases ( $f_2 \sim f_6$ ) is the same as that for  $f_1$ , as shown in the second row of Table 2. And as for  $Cov_{count}$ , the number "1" in the binary vector will be replaced with the execution frequency, which does not help to differentiate the six failures either. Thus, neither  $Cov_{hit}$  nor  $Cov_{count}$  can effectively perform failure indexing in this example, i.e., the scenario where failures caused by different faults have the same execution coverage.

Also recall the description in Section 2.2, the SD-based failure proximity represents a failed test case as a suspiciousness ranking list of program statements based on this failed test case and passed test cases  $T_S$ . For example, to represent  $f_1$ , we can combine  $f_1$  with  $T_S$  to form a new test suite,  $f_1 \cup T_S$ , run this test suite on the program, and collect the coverage information. The coverage will be delivered to an SBFL technique to calculate risk values of being faulty for each program statement, and then produce a suspiciousness ranking list by descendingly ordering statements by their risk

values. Here we employ a recognized SBFL technique, GP03 [95], whose expression is given in Formula 1, to complete this process<sup>6</sup>.

$$\text{Suspiciousness}_{GP03} = \sqrt{|N_{CF}^2 - \sqrt{N_{CS}}|}. \quad (1)$$

The generated suspiciousness ranking list is  $[1, 1, 3, 4, 4, 4, 4, 4, 4, 4, 4, 4, 15, 15, 15]$ <sup>7</sup>, which will serve as the proxy for  $f_1$ . However, the SD-based representation for the other five failed test cases ( $f_2 \sim f_6$ ) is the same as that for  $f_1$ , as shown in the third row of Table 2. That is, despite the incorporation of SBFL techniques, the SD-based failure proximity can still not differentiate the six failed test cases, because the SD-based strategy essentially relies on program coverage as well, while the six failed test cases triggered by different faults here exhibit the same execution coverage.

That is to say, when failures caused by different faults have identical coverage, even the most widespread and state-of-the-art failure proximities are still not enough to deliver promising outcomes. This situation, is, unfortunately, very common in practice, as a prestigious work in the field of failure indexing pointed out previously:

*“A significant portion of execution profiles would be the same even if these failures are due to different faults” [55].*

Thus, developing a better failure proximity is of great significance.

Moreover, neither the CC-based nor the SD-based failure proximity represents failures in a human-friendly form, it is hard for human developers to comprehend the failure indexing result based on such forms of failure representation. Specifically, whether a numerical vector in the CC-based strategy or a ranking list in the SD-based strategy, its length is equal to the number of executable statements of the faulty program. Unlike our motivating example, in real-world debugging, a faulty program can easily have tens of thousands of statements. As a consequence of which, the proxy for a failure will be a numerical vector or a ranking list of great length. It is obvious that in such a scenario, given failure indexing results, human developers are hard to comprehend why these numerical vectors or ranking lists should be clustered into the same (or different) group(s).

This motivating example confirms the challenges we point out in Section 1, namely, 1) The effectiveness of failure indexing is far from promising, and 2) The outcome of failure indexing is hard to comprehend for human developers. In this paper, we propose SURE, a novel failure indexing approach based on the program memory spectrum, which is more effective than the state-of-the-art failure indexing technique to date. Program memory spectrum is in the form of visualized images, human developers can easily comprehend the failure indexing results given such human-friendly failure representation.

## 4 APPROACH

The program memory-based failure proximity utilizes the run-time memory information to represent failures, and measures the distance between failures based on the characteristics of memory

<sup>6</sup>SBFL techniques typically involve four spectrum notations:  $N_{CF}$  and  $N_{CS}$  represent the number of test cases that execute the statement and return the testing result of failed or passed, respectively,  $N_{UF}$  and  $N_{US}$  represent the number of test cases that do not execute it and return the testing result of failed or passed, respectively.

<sup>7</sup>In light of the experience of previous works [26, 33, 92], if several statements with the same suspiciousness form a tie, the rankings of all statements in the tie will be set to the beginning position of this tie.



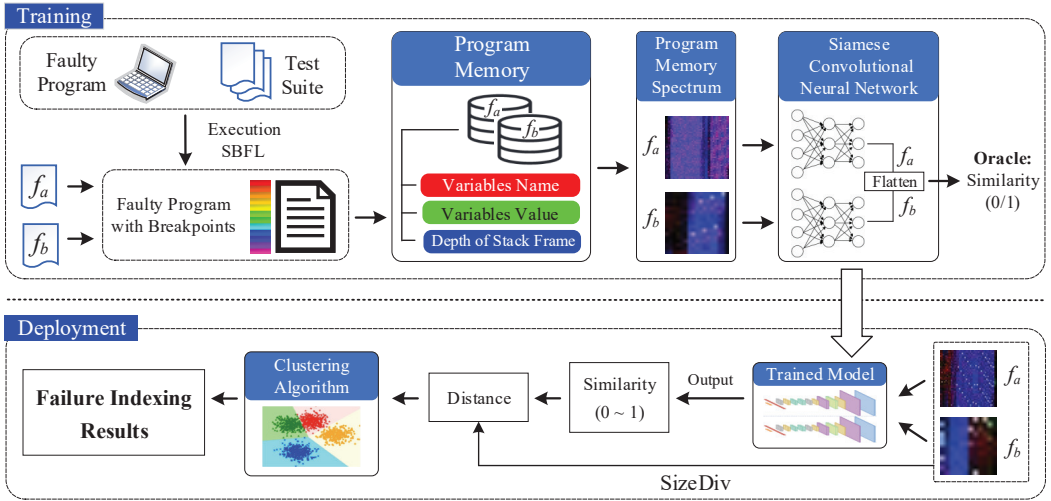


Fig. 1. Overview of SURE

information. Based on this description, we propose our approach, SURE. The overview of SURE is depicted in Figure 1. For a faulty program under test and a given test suite containing multiple failed test cases, SURE performs the following steps to achieve failure indexing:

- (1) **Breakpoint determination.** we first employ an SBFL technique to obtain the suspiciousness of all program statements, and then set the Top- $x\%$  riskiest statements as the breakpoints (the determination of the value of  $x$  will be investigated in the first research question of Section 5.1).
- (2) **Collection of run-time program memory.** For a failed test case, we run it on the faulty program, and collect the run-time memory information at the preset breakpoints. Specifically, the run-time memory information we collect includes variables' names, variables' values, and the depth of the stack frame.
- (3) **Generation of program memory spectrum.** The collected memory information will be converted to program memory spectrum (PMS), which is in the form of visualized images in a human-friendly way.
- (4) **Model training.** For any pair of PMS images (i.e., any pair of failed test cases), we will assign a label to it: If these two failed test cases are triggered by the same fault, the label is 1, indicating that the similarity between them should be 1. Otherwise, if these two failed test cases have different root causes, the label is 0. Then, we feed such labeled pairs of PMS images to a Siamese convolutional neural network to train.
- (5) **Prediction.** Once the model is trained, it can be used to handle the prediction of the similarity between two unseen PMS images.
- (6) **Distance calculation.** For the value of the similarity between a pair of failed test cases, we first subtract it from 1, and then multiply this difference by a parameter, *SizeDiv*, as the final distance. *SizeDiv* reflects the divergence between the sizes of two PMS images.
- (7) **Clustering.** The clustering algorithm will receive the distances between any pair of failed test cases, and based on which produces several clusters of failed test cases.

We detailedly introduce the above seven steps of SURE in Section 4.1 ~ Section 4.7, respectively. Moreover, we give a running example to facilitate readers' comprehension in Section 4.8.

#### 4.1 Breakpoint determination

Given a faulty program  $P$  containing  $l$  statements, a test suite  $T$  comprising failed test cases  $T_F$  and passed test cases  $T_S$  ( $T_F \cup T_S = T$  and  $T_F \cap T_S = \emptyset$ ). We employ an SBFL technique to calculate the suspiciousness of  $l$  program statements of  $P$  based on  $T$ , and rank all statements in descending order of suspiciousness. Then, we determine the Top- $x\%$  riskiest statements as the breakpoints (the value of  $x$  will be investigated in the first research question of Section 5.1). As such, we can get  $q$  breakpoints,  $bp_1, bp_2, \dots, bp_j, \dots, bp_q$ , where the value of  $q$  is  $\lfloor l \times x\% \rfloor$ , and  $bp_j$  is the breakpoint ranked  $j^{\text{th}}$  in suspiciousness.

The intuition behind this strategy is that the statements with higher risk values are more likely to be faulty, and run-time memory information gathered at these positions could have a stronger capability to reveal faults, thus can contribute more to representing failures.

As a reminder, SBFL is not required to determine the exact position of faults in this process. We simply need to identify a collection of highly-risky program statements, which are sufficient to specify breakpoints that can monitor the execution of abnormal programs.

#### 4.2 Collection of run-time program memory

For the  $i^{\text{th}}$  failed test case in  $T_F$ ,  $f_i \in T_F$ , executing it on  $P$  and gathering the run-time memory information at each breakpoint. Specifically, the collected memory information of  $f_i$  is:

$$MI_i = \{bp'_1 : V_1^i, bp'_2 : V_2^i, \dots, bp'_j : V_j^i, \dots, bp'_q : V_q^i\}, \quad (2)$$

where  $V_j^i$  is the memory information collected at  $bp'_j$  during executing  $f_i$ . Notice the difference between  $bp'_j$  here and  $bp_j$  in Section 4.1:  $bp'_j$  is the  $j^{\text{th}}$  breakpoint to be executed during the actual running of  $f_i$ . This arrangement integrates control flows information into our method, because program control flows are also helpful for failure representation to an extent. Supposing that there are  $m_{ij}$  variables queried at  $bp'_j$  during the running of  $f_i$ ,  $V_j^i$  comprises three lists with length  $m_{ij}$ :  $V_{name_j}^i$ ,  $V_{value_j}^i$ , and  $V_{depth_j}^i$ :

$V_{name_j}^i$  contains the names of the  $m_{ij}$  variables queried at  $bp'_j$  during running  $f_i$ :

$$V_{name_j}^i = [name_1, name_2, \dots, name_{m_{ij}}], \quad (3)$$

$V_{value_j}^i$  contains the values of these  $m_{ij}$  variables<sup>8</sup>:

$$V_{value_j}^i = [value_1, value_2, \dots, value_{m_{ij}}], \quad (4)$$

$V_{depth_j}^i$  contains the depth of the stack frame of these  $m_{ij}$  variables:

$$V_{depth_j}^i = [depth_1, depth_2, \dots, depth_{m_{ij}}]. \quad (5)$$

If we are from the perspective of a failed test case, that is, considering all breakpoints together, we can integrate the memory information collected at all breakpoints into a hunk:

$$V_{name}^i = V_{name_1}^i \oplus V_{name_2}^i \oplus \dots \oplus V_{name_j}^i \oplus \dots \oplus V_{name_q}^i, \quad (6)$$

$$V_{value}^i = V_{value_1}^i \oplus V_{value_2}^i \oplus \dots \oplus V_{value_j}^i \oplus \dots \oplus V_{value_q}^i, \quad (7)$$

$$V_{depth}^i = V_{depth_1}^i \oplus V_{depth_2}^i \oplus \dots \oplus V_{depth_j}^i \oplus \dots \oplus V_{depth_q}^i \quad (8)$$

<sup>8</sup>In our experiments, we find that regarding a variable's value as a string is beneficial for distance measurement. For some special types of variables, further action will be taken. For example, for a pointer, we will further index its value by address.

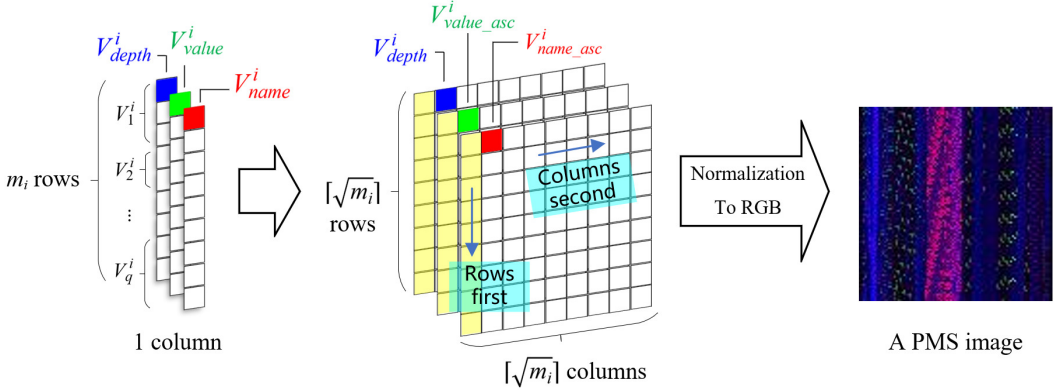


Fig. 2. Workflow of PMS images generation

where  $\oplus$  denotes the append operation between two lists. We use  $m_i$  to denote the length of  $V_{name}^i$ ,  $V_{value}^i$ , and  $V_{depth}^i$ , which is the number of variables queried at all breakpoints during the running of  $f_i$ :

$$m_i = |V_{name}^i| = |V_{value}^i| = |V_{depth}^i| = m_{i1} + m_{i2} + \dots + m_{ij} + \dots + m_{iq} \quad (9)$$

During the execution of a failed test case, when the program stops at a breakpoint, only the memory data at the position before the execution of that breakpoint can be collected. To ensure that memory information is gathered regarding the predetermined positions produced by SBFL, when faulty programs stop at each of the breakpoints, we continue executing a further step and then carry out the collection operation. Besides, if a program statement is executed for more than once, we collect variables' values when the execution is completed, since the latest value reflects the entire accumulation of the execution. As a reminder, for the sake of cost saving, if a program statement contains function calls, we concentrate on the original location of the preset breakpoint rather than iteratively navigating to the callee to query memory information (Experiments in Section 6 demonstrate that this strategy is good enough for failure indexing).

### 4.3 Generation of program memory spectrum

For each  $f_i \in T_F$ , we have collected the run-time memory information regarding it, i.e.,  $MI_i$ . Now we introduce based on  $MI_i$ , how to generate the corresponding program memory spectrum  $PMS_i$ .

The workflow of the generation of PMS images is illustrated in Figure 2. Specifically,  $MI_i$  is a matrix with the shape of  $(1 \times m_i \times 3)$ , as shown in **the leftmost sub-figure of Figure 2**, where  $m_i$  is the number of variables queried at all breakpoints during the running of  $f_i$ , and "3" indicates the three dimensions: variables' names, variables' values, and the depth of the stack frame. The arrangement of this  $(1 \times m_i \times 3)$  matrix is in the order of execution of breakpoints. Specifically,  $V_1^i$ , i.e., the run-time memory information collected at the first executed breakpoint during running  $f_i$ , is at the top. Similarly,  $V_q^i$ , i.e., the run-time memory information collected at the last executed breakpoint (there are  $q$  breakpoints in total) during running  $f_i$ , is at the bottom.

In practical development, the value of  $m_i$  is generally very large, thus the shape of this  $(1 \times m_i \times 3)$  matrix is quite unbalanced. We reshape this matrix to a  $(\lfloor \sqrt{m_i} \rfloor \times \lfloor \sqrt{m_i} \rfloor \times 3)$  matrix. The principle of this reshaping process is "row first, column second". Specifically, we take the first  $\lfloor \sqrt{m_i} \rfloor$  elements from top to bottom in the  $(1 \times m_i \times 3)$  matrix as the first column of the  $(\lfloor \sqrt{m_i} \rfloor \times$

$\lceil\sqrt{m_i}\rceil \times 3$  matrix, and then take the  $(\lceil\sqrt{m_i}\rceil + 1)^{\text{th}}$  to  $2\lceil\sqrt{m_i}\rceil^{\text{th}}$  elements from top to bottom in the  $(1 \times m_i \times 3)$  matrix as the second column of the  $(\lceil\sqrt{m_i}\rceil \times \lceil\sqrt{m_i}\rceil \times 3)$  matrix, and so on. Notice that if  $m_i$  is not a perfect square number, the last column of the  $(\lceil\sqrt{m_i}\rceil \times \lceil\sqrt{m_i}\rceil \times 3)$  matrix cannot be filled. In such a scenario, the few remaining elements will be set to the number zero. After the reshaping process is completed, a  $(\lceil\sqrt{m_i}\rceil \times \lceil\sqrt{m_i}\rceil \times 3)$  matrix can be obtained, as shown in **the middle sub-figure of Figure 2**.

The reshaped matrix can be regarded as three  $(\lceil\sqrt{m_i}\rceil \times \lceil\sqrt{m_i}\rceil \times 1)$  sub-matrices:  $[:, :, 0]$ ,  $[:, :, 1]$ , and  $[:, :, 2]$ , in which the first sub-matrix is actually  $V_{name}^i$ , and the second as well as the third sub-matrices are  $V_{value}^i$  and  $V_{depth}^i$ , respectively. Among them, the elements of  $V_{name}^i$  and  $V_{value}^i$  are in the form of characters, and the elements of  $V_{depth}^i$  are integers. To make them have a uniform format, we employ an existing method in the reference [49] to convert elements of  $V_{name}^i$  and  $V_{value}^i$  from the character form to the numeric form. Specifically, for a variable's name or a variable's value in the string form (denoted as  $str$ ), we use the following formula to convert it to an integer (denoted as  $str_{asc}$ ):

$$str_{asc} = \sum_{i=0}^{|str|} (i + 1) * Ascii(str[i]), \quad (10)$$

where  $Ascii()$  gets the Ascii code given a character (the strategy of converting strings to values by employing ascii codes has been adopted by many previous works [56, 58, 75]). The multiplication by  $(i + 1)$  is to handle the case where two strings have the same batch of characters but in different order. For example, a string "Ee01" can be converted to an integer 611, which is calculated by  $1 * Ascii('E') + 2 * Ascii('e') + 3 * Ascii('0') + 4 * Ascii('1')$ .

Now, each element of the  $(\lceil\sqrt{m_i}\rceil \times \lceil\sqrt{m_i}\rceil \times 3)$  matrix is already an integer, this three-channel matrix is ready to be transformed into a colored image. Here we employ RGB (Red, Green, and Blue), a broadly-used additive color model, to complete this transformation process, because of its prevalence, availability, and simplicity. In RGB, three primary colors each have values ranging from 0 to 255, they can be combined together in various ways to reproduce a broad array of colors in the visible spectrum [23, 30]. To make the  $(\lceil\sqrt{m_i}\rceil \times \lceil\sqrt{m_i}\rceil \times 3)$  matrix adapt to the RGB model, we normalize the elements in each of its sub-matrices, namely,  $[:, :, 0]$ ,  $[:, :, 1]$ , and  $[:, :, 2]$  separately, according to the mentioned value range of RGB. The normalization can also alleviate the imbalance among the three dimensions, because the three sub-matrices reflect different resources of memory information thus have different scales. After normalization, the  $(\lceil\sqrt{m_i}\rceil \times \lceil\sqrt{m_i}\rceil \times 3)$  matrix can be directly transformed into a colored image, as shown in **the rightmost sub-figure of Figure 2**.

The motivation for transforming run-time memory information into PMS images is to facilitate the following task of distance calculation. Specifically, memory information is raw data collected for failure representation, its scale could be typically large, and its form is not easy to recognize or handle by machines either. **Images are a proper candidate form to rephrase raw memory data, because the form of images can effectively integrate various sources of information of memory data into a single block, such a re-representation for raw data can provide well-organized and highly-structured objects for the downstream deep learning technique that performs distance calculation.** Moreover, the form of visualized images can also boost human developers' comprehension of failure indexing results, because images can be smoothly recognized by human beings at a glance.

#### 4.4 Model training

The training phase of SURE is illustrated in the upper part of Figure 1. Specifically, for a pair of failed test cases, we first collect the corresponding two sets of memory information on the faulty program

with breakpoints, and then transform them into two PMS images, respectively, as mentioned in Section 4.1 ~ 4.3.

We use the idea of Siamese learning to construct a similarity prediction model for two failed test cases (in the form of PMS images), because Siamese learning is a classical supervised learning strategy in similarity learning tasks [11, 13, 15]. The goal of a Siamese learning-based model is to learn a similarity metric between pairs of input samples, it typically consists of two identical sub-networks for feature extraction that share the same structure and weights, as well as fully connected layers that are responsible for taking the learned representations from the mentioned feature extraction networks, thus producing a final output that reflects the similarity between two input samples.

The architecture of our Siamese learning model is designed as follows: two identical feature extraction networks (will be specified in the second research question of Section 5.1), and two fully connected layers (with a ReLU activation function in between). The loss function we use is the binary cross entropy loss (BCELoss), which is commonly used in deep learning especially for binary classification tasks, as shown in Formula 11:

$$L = -\frac{1}{N} \sum_{i=1}^N [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)], \quad (11)$$

where  $N$  is the number of the data samples,  $y_i$  is the true binary label for the  $i^{\text{th}}$  sample, and  $p_i$  is the predicted probability that the  $i^{\text{th}}$  sample is positive. As a reminder, the true binary label of a pair of PMS images is whether “1” (i.e., these two failed test cases are triggered by the same fault, thus the similarity between them should be 1) or “0” (i.e., these two failed test cases have different root causes, thus the similarity between them should be 0). We denote pairs of PMS images with the label “1” as *positive samples*, while denote pairs of PMS images with the label “0” as *negative samples*.

In particular, we use only 30% of the simulated faults to train the model, and use 70% of the simulated faults as well as real-world faults to test the model (please refer to Section 5.3 for the details of the faults used in the experiments).

#### 4.5 Prediction

Once the similarity prediction model is trained, it can be used to predict the similarity between two unseen PMS images (i.e., two failed test cases). This similarity reflects the likelihood of these two failures being triggered by the same fault. Before sending two images to the model, we first preprocess them into a uniform size. This is because handling two images with different sizes is hard for Siamese-based models, in light of the structure of the two feature extraction networks must be identical.

#### 4.6 Distance calculation

The output of our Siamese learning model is the similarity between two failures. Before sending the value of the produced similarity to the downstream clustering algorithm, we further process it by Formula 12, to make up for the loss caused by harmonizing the sizes of images in the previous step.

$$Distance = (1 - Similarity) \times SizeDiv. \quad (12)$$

We first explain the first half of this formula, i.e.,  $1 - Similarity$ . The similarity produced by the Siamese model measures how similar two PMS images are. To adapt it for clustering algorithms, we subtract the value of the similarity between two failures from 1, to reflect how dissimilar they are. Next, we explain the multiplication of the factor *SizeDiv*. As mentioned in Section 4.3, a PMS image

representing the failed test case  $f_i$  is a square whose side length is  $\lceil \sqrt{m_i} \rceil$ , where  $m_i$  is the number of variables queried at all breakpoints during the execution of  $f_i$ . It is obvious that different  $f_i$  may correspond to different  $m_i$ , thus resulting in different side lengths of PMS images. Considering that Siamese models are hard to handle two images with different sizes, in the prediction phase, we preprocess two input images into the same size before feeding them to the model. But as mentioned previously, the side length of PMS images embodies the number of queried variables during running a failed test case, which can reveal the characteristics of a failure to some extent, and thus can also contribute to the distance measurement between two failed test cases. To alleviate this issue, we design *SizeDiv*, a factor that is able to quantitatively reflect the divergence between the sizes of two images. Specifically, *SizeDiv* is calculated by dividing the side length of the larger image by the side length of the smaller image.

#### 4.7 Clustering

In Section 2.1, we have pointed out that there is no clustering technique that is universally applicable in uncovering the variety of structures present in multidimensional data sets [34]. That is, the clustering algorithm is not the core of failure indexing techniques. Therefore, we employ the clustering component of MSeer [26], the state-of-the-art failure indexing technique to date, to complete the clustering phase of SURE. The clustering algorithm is twofold, namely, the faults number estimation and the clustering. The former aims to predict the number of clusters (i.e., the number of faults) given the distance information between data samples (i.e., failed test cases), and the latter is responsible for clustering data samples (i.e., failed test cases) into different groups. Next, we give a concise description of these two steps, for more details please refer to the reference [26].

**For the faults number estimation.** It is well-recognized that one of the trickiest challenges in clustering lies in the estimation of the number of clusters [24, 45, 73]. Putting it into the context of failure indexing, we can claim that predicting the number of faults given a number of failures is very important. The adopted clustering algorithm presented a novel mountain method-based approach inspired by the previous works [14, 94], to perform the faults number estimation and the assignment of initial medoids to these clusters at the same time. Specifically, the adopted clustering algorithm first calculates a potential value for each failed test case according to the density of its surroundings, such a potential value is used to measure the possibility of a data point being set as a medoid. And then, 1) The failure with the highest potential value will be chosen as the first medoid. 2) The potential values of all failed test cases will then be updated in accordance with their distance from the newest medoid. 3) Repeating the above two steps iteratively, until the maximum potential value falls within a certain threshold.

**For the clustering.** Once the number of clusters and the initial medoids are determined, all failures are ready to be clustered. The adopted clustering algorithm utilizes K-medoids, a widely-used clustering strategy, to complete this process. The K-medoids strategy sets actual (not virtual) data points as medoid and thus can be more applicable to SURE, because the mean of memory information is difficult to define. Moreover, the K-medoids strategy has been shown to be very robust to the existence of noise or outliers [41].

#### 4.8 Running example

In the motivating example in Section 3, we use a toy program to show that the most prevalent and advanced approach to date can still not good enough to produce satisfactory failure indexing results. That toy program successfully reveals the bottleneck of existing approaches with only 17 program statements. But such a small-scale program is not suitable to illustrate the workflow of SURE, because a small number of program statements will correspondingly result in a small number of breakpoints, thus further resulting in a small number of queried variables. In such cases,

the generated PMS images will be very small, which does not affect the running of SURE but can hinder readers' comprehension as a running example. Therefore, here we use a complete running example to illustrate the workflow of SURE, where CC and SD do not work.

We employ Grep, a classical tool aiming to print lines in specific files that contain a match of the given patterns. We obtain Grep (v2.4) which contains 13,274 lines from the Software Infrastructure Repository (SIR) [67], and randomly mutate three positions of the original clean program to inject three faults into the program, to obtain a 3-bug faulty version. The three mutants are shown in Listing 1.

---

```

3164  if (leftoversf)
3165  {
3166      copyset(leftovers, labels[ngroups]);
3167      copyset(intersect, labels[j]);
3168 -   MALLOC(grps[ngroups].elems, position, d->nleaves);
3168 +   MALLOC(grps[ngroups].elems, position, d->nleaves*-1);
3169      copy(&grps[j], &grps[ngroups]);
3170      ++ngroups;
3171  }
...
4681  for (ep = text + size - 11 * len;;)
4682  {
4683      while (tp <= ep)
4684      {
4685 -         d = d1[U(tp[-1])], tp += d;
4685 +         d = d1[U (! tp[-1])], tp += d;
4686         d = d1[U(tp[-1])], tp += d;
4687         if (d == 0)
4688             goto found;
...
7480  compile_stack.stack = TALLOC (INIT_COMPILE_STACK_SIZE, compile_stack_elt_t);
7481  if (compile_stack.stack == NULL)
7482      return REG_ESPACE;
7484 -  compile_stack.size = INIT_COMPILE_STACK_SIZE;
7484 +  compile_stack.size = 0;
7485  compile_stack.avail = 0;

```

---

Listing 1. Mutants in the running example

We compile the mutated program and run the accompanying test suite, observing a series of failures and also getting a number of passed test cases. Then, we conduct failure indexing following the aforementioned steps: For the breakpoint determination (Section 4.1), we employ an SBFL technique (e.g., DStar [82]) to determine the suspiciousness of program statements and setting breakpoints. For the collection of run-time program memory (Section 4.2), we run failed test cases and collect the run-time memory information *MI* in accordance with Formula 2. *MI* comprises variables' names, variables' values, and the depth of the stack frame, whose content and structure are defined as Formula 3 ~ Formula 8, and *MI*'s scale is determined by Formula 9. For the generation of PMS images (Section 4.3), we first convert variables' names and values to the numeric form using Formula 10, and then transform each set of memory information into a PMS image.

It is hard and not necessary to completely illustrate this running example because there are too many failed test cases. Here we randomly select seven of them, namely,  $t_{43}$ ,  $t_{178}$ ,  $t_{180}$ ,  $t_{182}$ ,  $t_{224}$ ,  $t_{252}$ , and  $t_{279}$ , to show their PMS images and describe the following steps<sup>9</sup>. Among them,  $t_{43}$  and  $t_{279}$  are triggered by the mutant located in Line 3168,  $t_{178}$  and  $t_{180}$  are triggered by the mutant located in Line 4685, and  $t_{182}$ ,  $t_{224}$  as well as  $t_{252}$  are triggered by the mutant located in Line 7484. The

<sup>9</sup>The remaining failures omitted here can also support the conclusion of this example.

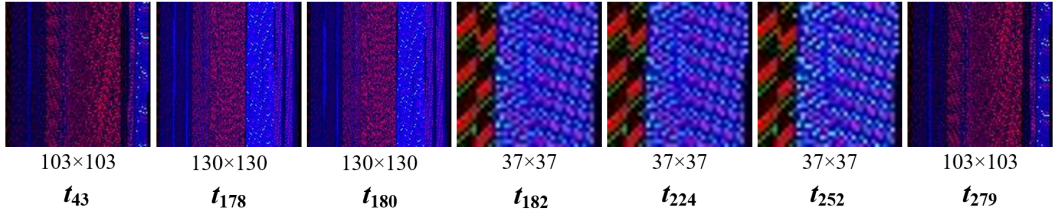


Fig. 3. The PMS images of the failed test cases

Table 3. Distance information between failures

	$t_{43}$	$t_{178}$	$t_{180}$	$t_{182}$	$t_{224}$	$t_{252}$	$t_{279}$
$t_{43}$	0.00	1.26	1.26	2.78	2.78	2.78	0.01
$t_{178}$	1.26	0.00	0.02	3.51	3.51	3.51	1.26
$t_{180}$	1.26	0.02	0.00	3.51	3.51	3.51	1.26
$t_{182}$	2.78	3.51	3.51	0.00	0.01	0.02	2.78
$t_{224}$	2.78	3.51	3.51	0.01	0.00	0.02	2.78
$t_{252}$	2.78	3.51	3.51	0.02	0.02	0.00	2.78
$t_{279}$	0.01	1.26	1.26	2.78	2.78	2.78	0.00

seven PMS images as well as their size, are given in Figure 3. Then, we employ the trained model (Section 4.4) to predict the similarity between pairs of failures (Section 4.5), and get the distance information (Section 4.6), as given in Table 3. Finally, all failures are clustered into different groups using the two phases of the clustering algorithm: the faults number estimation and the clustering (Section 4.7). The outcome of the clustering algorithm is:  $\{t_{43}, t_{279}\}$ ,  $\{t_{178}, t_{180}\}$ , and  $\{t_{182}, t_{224}, t_{252}\}$ , which is consistent with the characteristic of the distance information given in Table 3. Obviously, this clustering outcome achieves promising failure indexing since the number of clusters (i.e., the number of faults, three) is correctly estimated, and all failures are grouped correctly according to the root cause.

Moreover, we can observe that the characteristics of the PMS images in Figure 3 are highly identical to the failure indexing result. Specifically, the PMS images of  $t_{43}$  and  $t_{279}$  are similar and have the same size. Likewise, we can observe the same phenomenon on  $\{t_{178}, t_{180}\}$  and  $\{t_{182}, t_{224}, t_{252}\}$ . With the support of the PMS images, developers can be easily convinced by the failure indexing outcome and thus adopt it.

## 5 EXPERIMENTAL SETUP

### 5.1 Research Question

- **RQ1: Does the value of the breakpoint determination threshold impact SURE's effectiveness?**

As mentioned in Section 4.1, SURE will take the Top- $x\%$  riskiest program statements as breakpoints. We investigate how the value of  $x$  impacts the effectiveness of SURE. Specifically, we compare the effectiveness of SURE when such a threshold is set from 0% to 100% with 10% increments, i.e., 10%, 20%, 30%, ..., 100%.

- **RQ2: Does the selection of the feature extraction network impact SURE's effectiveness?**



For a pair of failures, SURE uses a Siamese learning-based model to predict its distance. Because the input samples in our context are in the form of PMS images, the feature extraction network is correspondingly convolutional. In this RQ, we investigate how the feature extraction network impacts the effectiveness of SURE. Specifically, we compare the effectiveness of SURE when the network configures VGG-16 [66], AlexNet [46], and ResNet-18 [29].

• **RQ3: How does SURE perform compared with the most prevalent and advanced failure indexing technique?**

In Section 2, we have introduced that the CC-based strategy has been extensively used by the failure indexing community due to its simplicity, and the SD-based strategy has been recognized as the state-of-the-art solution. Thus, we compare SURE with these two types of techniques for more convincing evaluation.

For the CC-based strategy, we select *Cov<sub>hit</sub>* as our baseline, because it is the most common technique in the CC class. Specifically, it represents a failed test case as a binary numeric vector with a length equal to the number of program statements, the  $i^{\text{th}}$  element of the vector will be set to 1 if the  $i^{\text{th}}$  statement is covered by this failed test case, and 0 otherwise. Moreover, seeing that there have emerged some works concerning the impact of the execution frequency of program statements on debugging [65, 74, 80], we also adopt another technique in the CC class, *Cov<sub>count</sub>*, as our baseline. Specifically, *Cov<sub>count</sub>* also represents a failed test case as a numeric vector with a length equal to the number of program statements, while the  $i^{\text{th}}$  element of the vector will be set to the actual execution frequency if the  $i^{\text{th}}$  statement is covered (rather than a constant, 1), and 0 otherwise.

For the SD-based strategy, we select *MSeer* as our baseline, because it is the state-of-the-art technique not only in the SD class but also in the current field of failure indexing. Specifically, it first runs a failed test case along with passed test cases against the faulty program, and inputs the gathered coverage information to Crosstab [84], an SBFL technique, to calculate the suspiciousness of being faulty for each program statement. And then, ranking all statements in descending order according to their suspiciousness values. Such a ranking list will serve as the proxy for this failed test case. Moreover, in our previous empirical study, we comprehensively investigated the factors that could impact the effectiveness of MSeer, finding that the selection of SBFL formulas matters [69]. According to our result, if Crosstab is replaced by GP19 (another SBFL formula evolved by genetic programming in the reference [95]), MSeer will do much better. Thus, we manually upgrade MSeer to *MSeer<sub>GP19</sub>*, which employs GP19 to calculate suspiciousness, and include it as our baseline<sup>10</sup>.

To summarize, we compare SURE with four techniques, i.e., *Cov<sub>hit</sub>* and *Cov<sub>count</sub>* (the CC-based strategy), as well as *MSeer* and *MSeer<sub>GP19</sub>* (the SD-based strategy).

• **RQ4: To what extent can SURE help human developers comprehend failure indexing results?**

As a visualized failure indexing technique, SURE represents failed test cases as human-friendly PMS images. In this way, developers who receive the failure indexing result can not only know how all failures are grouped, but also be aware of why they are grouped in the current manner. In this RQ, we quantitatively investigate to what extent developers can comprehend the failure indexing result with the support of PMS images. Similar to RQ3, we compare the comprehensibility of SURE with that of CC-based and SD-based techniques. Specifically, we conduct a comparison of developers' performance on comprehending failure indexing results among given PMS images (i.e., the fingerprinting function of SURE), code coverage vectors (i.e., the fingerprinting function of CC-based strategies), and suspiciousness ranking lists (i.e., the fingerprinting function of SD-based

<sup>10</sup>As a reminder, *MSeer<sub>GP19</sub>* is not an existing published failure indexing technique, it is manually created by us to further evaluate the competitiveness of SURE.

strategies). This human study involves 9 tasks and 15 experienced graduate students in Computer Science from Wuhan University.

## 5.2 Parameter Setting

SURE needs to first determine the suspiciousness value of program statements, and thus sets breakpoints at those highly-risky positions. In this phase, we use DStar [82], the state-of-the-art SBFL technique that has been utilized broadly. Considering the preference for DStar in many other studies (such as the references [8, 60, 81]), we set the value of  $*$  in DStar to 2, the most thoroughly-explored value, in our experiments. Such a choice is not hard-coded but can be configurable, any other fault localization techniques working at the granularity of statements can be adapted to this phase.

In the training phase, we set the value of batch size as 16. The initial value of the learning rate is defined as  $1e-4$ , and it will be multiplied by 0.96 after each epoch.

## 5.3 Dataset

For comprehensive and robust evaluation, in the experiments, we use the benchmark involving both simulated faults and real-world faults. For simulated benchmarks, we manually seed different types of faults into clean programs, to obtain simulated faulty programs, in light of the fact that previous research has confirmed that mutation-based faults can provide credible results for experiments in software testing and debugging [6, 7, 22, 40, 52, 63]. The clean programs are obtained from the Software Infrastructure Repository (SIR) [21]. For real-world benchmarks, we search for the Defects4J projects [39], one of the most popular datasets in the current field of software testing and debugging, according to the test case transplantation strategy proposed by An et al. [5]. All faulty programs contain one, two, three, four, or five bugs (also referred to as 1-bug, 2-bug, 3-bug, 4-bug, and 5-bug faulty versions, respectively), because such numbers of faults are most-widely investigated in previous studies [26, 48, 88, 100].

**5.3.1 Simulated faulty programs.** SIR is a classical repository in the community of software testing and debugging, which has been employed in numerous pioneering studies [32, 44, 83, 96, 101]. From SIR we obtain Flex, Grep, Gzip, and Sed, four projects that have been extensively adopted in earlier works [10, 27, 70], as the clean programs to mutate (i.e., seed faults). The concise information about them can be found in Table 4. We utilize an open-source tool with “13” fork and “23” star on GitHub to perform mutation, which defines 67 types of points that can be mutated, and provides several mutation operators for each [9]. The mutation operators we leverage can be categorized into the following two classes:

- **Assignment Fault** [35]: Editing the value of constants in the statement, or replacing the operators such as addition, subtraction, multiplication, division, etc. with each other.
- **Predicate Fault** [93]: Reversing the *if-else* predicate, or deleting the *else* statement, or modifying the decision condition, and so on.

To create an  $r$ -bug faulty version ( $r = 2, 3, 4, 5$ ), the faults from  $r$  individual 1-bug faulty versions are injected into the same program. Such a strategy has been adopted by the majority of the published studies in the field of multi-fault debugging [1, 31, 33, 47, 99]. In total, we obtain 1,000 SIR faulty versions containing 1~5 faults.

**5.3.2 Real-world faulty programs.** Defects4J is one of the most popular benchmarks in the current field of software testing and debugging, due to its realism and ease of use [39]. Defects4J is typically for single-fault scenarios, namely, no matter how many bugs are contained in a faulty program, the provided test suite is only sufficient to reveal one of them. Such a characteristic hinders its

Table 4. Benchmarks

Language	Project	Version	kLOC	Functionality
C	Flex	2.5.3	14.5	Parser generator
	Grep	2.4	13.5	Text matcher
	Gzip	1.2.2	7.3	File archiver
	Sed	3.02	10.2	Stream editor
	Chart	2.0.0	96.3	Chart library
	Closure	2.0.0	90.2	Closure compiler
Java	Lang	2.0.0	22.1	Apache commons-lang
	Math	2.0.0	85.5	Apache commons-math
	Time	2.0.0	28.4	Date and time library

use in failure indexing studies. To adapt Defects4J to multi-fault scenarios, An et al. presented a test case transplantation strategy [5]. Specifically, the majority of Defects4J faulty versions are indexed chronologically according to the revision date, a lower ID indicates a more recent version. Therefore, the fault in a newer version is also likely to be contained in an older version. For example, the fault of the faulty version Math-5b is found to exist in the faulty version Math-6b as well, it is not revealed in Math-6b simply due to the absence of the fault-revealing test case. If this test case is transplanted to Math-6b, the enhanced test suite is able to reveal both of the two faults, thus a 2-bug faulty version can be obtained [5]. Following this strategy, we search for a collection of multi-fault Defects4J programs, from Chart, Closure, Lang, Math, and Time, as shown in Table 4. Because Defects4J multi-fault versions are obtained by searching in real-life environments, their number is not very large. In total, we obtain 100 Defects4J faulty versions containing 1~5 faults.

#### 5.4 Metric

The mission of a failure indexing technique is to identify the mutual relationship among failures by clustering, i.e., determine which failures are triggered by the same (or different) fault(s). Correspondingly, the outcome of failure indexing is several clusters of failed test cases. For evaluation, we need to quantitatively measure to which extent the delivered clusters correctly reflect the true relationship among failures. There are two typical types of metrics in evaluating clustering outcomes, namely, external metrics [86] and internal metrics [72]. The former compares clustering results with the oracle (true linkages between failures and faults), while the latter examines inherent properties of the delivered clusters, such as compactness and separation, etc., typically when the oracle is inaccessible. In our controlled experiments, the oracle clusters can be obtained easily in advance. Specifically, for an  $r$ -bug SIR faulty version, it is generated by combining  $r$  individual single-bug faulty versions. We can run the test suite against these  $r$  faulty versions separately, thus becoming aware of the culprit fault of each failure. And for an  $r$ -bug Defects4J faulty version, it is generated by transplanting the failed test cases triggered by a fault to the faulty version that contains another fault. This transplantation process itself explicitly indicates the relationship between failures and underlying faults. Thus, we use external metrics to measure the effectiveness of failure indexing techniques. Concretely speaking, we employ four external metrics, the Fowlkes and Mallows Index (FMI), the Jaccard Coefficient (JC), the Precision Rate (PR), and the Recall Rate (RR), in our experiments, because they are all classical metrics for clustering and have been adopted by many published studies [86, 89, 103]. Among them, FMI and JC are pair of test cases-based, while PR and RR are single test case-based.

Table 5. Four scenarios in the pair of test cases-based metrics

Notation	Results of failure indexing	
	In the generated cluster	In the oracle cluster
SS	Same	Same
SD	Same	Difference
DS	Difference	Same
DD	Difference	Difference

Table 6. Four scenarios in the single test case-based metrics

Notation	Results of failure indexing	
	In the generated cluster	In the oracle cluster
TP	Positive	Positive
FP	Positive	Negative
TN	Negative	Negative
FN	Negative	Positive

**5.4.1 Pair of test cases-based metrics.** The pair of test cases-based metric refers to comparing the indexing consistency of each pair of failed test cases in the generated cluster with the oracle cluster. Four possible scenarios in the comparison are given in Table 5. Specifically, supposing that there are  $n$  failed test cases that need to be clustered, they can form  $C_n^2$  pairs by combining any two ones. For a pair, if its two member failures are determined to be triggered by the *Same* fault in the generated cluster, and these two are also divided into the *Same* group in the oracle cluster, this pair falls into the category of “SS”. Similarly, a pair is also possible to fall into the categories of “SD”, “DS”, or “DD”. We use  $X_{SS}$ ,  $X_{SD}$ ,  $X_{DS}$ , and  $X_{DD}$  to denote the numbers of the pairs fallen into “SS”, “SD”, “DS”, and “DD”, respectively. Obviously, the sum of  $X_{SS}$ ,  $X_{SD}$ ,  $X_{DS}$ , and  $X_{DD}$  will be  $C_n^2$ , because any pair of failures will belong to one of the four scenarios. FMI and JC can integrate these four notations into one metric, to comprehensively reflect the similarity between the generated cluster and the oracle cluster, as shown in Formula 13 and Formula 14, respectively.

$$FMI = \sqrt{\frac{X_{SS}}{X_{SS} + X_{SD}} \times \frac{X_{SS}}{X_{SS} + X_{DS}}}, \quad (13)$$

$$JC = \frac{X_{SS}}{X_{SS} + X_{SD} + X_{DS}}. \quad (14)$$

It can be proved that the intervals of FMI and JC are both  $[0, 1]$ , and that the larger the value in this range, the more effective clustering is.

**5.4.2 Single test case-based metrics.** The single test case-based metric refers to comparing the classification result of each failed test case in the generated cluster with the oracle cluster. Four possible scenarios in the comparison are given in Table 6. Specifically, for a failed test case, if it is predicted to be triggered by a fault in the generated cluster (i.e., *Positive*), and its root cause is indeed this fault (i.e., *True* prediction result), this failure falls into the category of “TP”. Likewise, if a failed test case is predicted not to be triggered by a fault in the generated cluster (i.e., *Negative*), while its root cause is actually this fault (i.e., *False* prediction result), this failure falls into the category of “FN”. Similarly, a failed test case is also possible to fall into the categories of “FP” and

“TN”. We use  $X_{TP}$ ,  $X_{FN}$ ,  $X_{FP}$ , and  $X_{TN}$  to denote the numbers of the failed test cases fallen into “TP”, “FN”, “FP”, and “TN”, respectively. Obviously, the sum of “TP”, “FN”, “FP”, and “TN” will be the number of failures, because any failure will belong to one of the four scenarios. PR and RR can integrate these four notations into one metric, to comprehensively reflect the similarity between the generated cluster and the oracle cluster, as shown in Formula 15 and Formula 16, respectively.

$$PR = \frac{X_{TP}}{X_{TP} + X_{FP}}, \quad (15)$$

$$RR = \frac{X_{TP}}{X_{TP} + X_{FN}}. \quad (16)$$

It can be proved that the intervals of PR and RR are both  $[0, 1]$ , and that the larger the value in this range, the more effective clustering is.

Notice that FMI, JC, PR, and RR evaluate the effectiveness of clustering of one faulty version. In our experiments, there are a series of faulty versions to make our evaluation abundant. Thus, for a failure indexing technique  $T$ , we determine its effectiveness by adapting the FMI, JC, PR, and RR metrics on one faulty version to that on multiple faulty versions:

$$S\_M_M^T = \sum_i^{V_{equal}^T} M_i, \quad (17)$$

where “S\_M” is the abbreviation for “Sum\_Metrics”.  $V_{equal}^T$  is the number of faulty versions whose number of faults is correctly predicted by using  $T$ .  $M_i$  is the metric value ( $M$  takes FMI, JC, PR, or RR) on the  $i^{\text{th}}$  “ $k == r$ ” faulty version<sup>11</sup>. As a reminder, an  $r$ -bug faulty version contains  $r$  faults, which is hard to know in advance. A failure indexing technique should first correctly predict the number of faults  $r$ , and then cluster all failures into  $r$  groups. For an  $r$ -bug faulty version, if the predicted number of faults  $k$  is not equal to  $r$ , we do not evaluate the clustering effectiveness on this faulty version, because prior studies have pointed out that the performance of a failure indexing technique can be mainly determined by those “ $k == r$ ” faulty versions [26, 69], that is, the contribution of “ $k != r$ ” ones is marginal. And also, it is indeed difficult to compare  $k$  delivered clusters with  $r$  oracle clusters. Notice that “ $k == r$ ” is just an ideal scenario (not necessary) for SURE. Even if  $k != r$ , SURE can also work. We will further explain this point in Section 8.

In summary, there are five evaluation metrics used in the experiments, which correspond to the two goals of failure indexing mentioned previously. Specifically,  $V_{equal}^T$  reflects the goal of *correct faults number estimation*, and  $S\_M_{FMI}^T$ ,  $S\_M_{JC}^T$ ,  $S\_M_{PR}^T$ , and  $S\_M_{RR}^T$  reflect the goal of *promising clustering*.

## 5.5 Environment

We collect program coverage and run-time memory information on Ubuntu 16.04.1 LTS with GCC 5.4.0 and JDB 1.8. The Siamese neural network model is trained and deployed on 4 GPUs of GeForce RTX 2080 Ti. The clustering process runs on a server equipped with 96 Intel Xeon(R) Gold 5218 CPU cores with 2.30GHz and 160 GB of memory.

<sup>11</sup>For an  $r$ -bug faulty version, if the predicted number of faults  $k$  is equal to  $r$ , we label this faulty version as “ $k == r$ ”.

Table 7. Performance of SURE on SIR

		10%	20%	30%	40%	50%	60%	70%	80%	90%	100%
VGG-16	$V_{equal}^T$	161	147	122	143	153	164	142	145	151	155
	$S\_M_{FMI}^T$	<b>126.78</b>	114.32	96.11	114.36	121.28	131.10	112.59	112.26	118.91	122.75
	$S\_M_{JC}^T$	<b>106.39</b>	94.79	80.36	96.97	102.25	111.44	94.95	93.18	99.82	103.33
	$S\_M_{PR}^T$	<b>120.24</b>	106.67	91.60	111.40	115.48	127.75	108.06	110.88	119.31	123.99
	$S\_M_{RR}^T$	<b>80.67</b>	71.20	59.52	72.80	75.92	84.55	69.20	63.20	72.66	73.17
AlexNet	$V_{equal}^T$	167	143	139	169	187	160	160	176	163	150
	$S\_M_{FMI}^T$	<b>131.15</b>	113.31	109.04	135.30	148.08	125.76	126.53	137.60	127.59	118.99
	$S\_M_{JC}^T$	<b>109.70</b>	95.43	91.11	115.11	124.54	105.42	106.46	115.04	106.70	100.56
	$S\_M_{PR}^T$	<b>123.10</b>	109.06	102.31	131.53	145.66	124.69	127.87	136.72	127.4	120.04
	$S\_M_{RR}^T$	<b>80.44</b>	70.11	69.80	86.33	89.91	76.78	80.89	82.49	79.07	74.08
ResNet-18	$V_{equal}^T$	118	153	151	157	150	161	157	147	172	137
	$S\_M_{FMI}^T$	<b>92.71</b>	121.25	119.68	125.56	115.85	126.42	122.81	115.38	136.39	109.13
	$S\_M_{JC}^T$	<b>77.48</b>	102.14	100.77	106.47	95.82	106.07	102.51	96.42	114.98	92.45
	$S\_M_{PR}^T$	<b>82.63</b>	115.51	115.36	122.86	114.62	121.60	121.79	116.69	135.98	109.16
	$S\_M_{RR}^T$	<b>59.37</b>	79.27	76.55	78.96	68.59	77.27	73.57	68.52	84.51	67.48

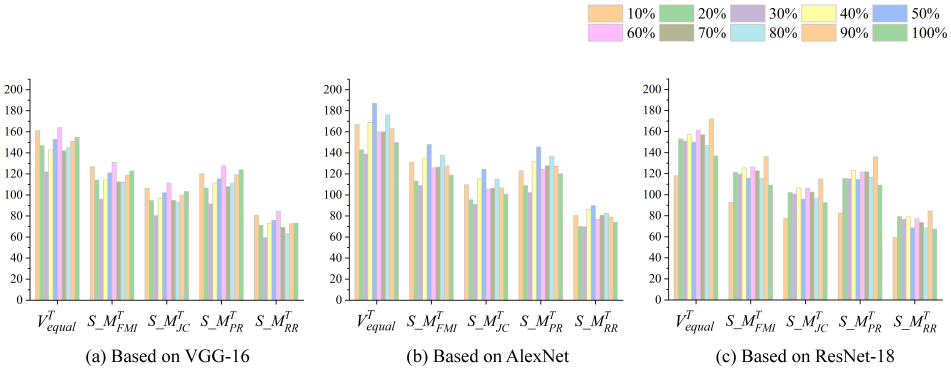


Fig. 4. The impact of the breakpoint determination threshold on SURE on SIR

## 6 RESULT AND ANALYSIS

### 6.1 RQ1: The value of the breakpoint determination threshold.

As the first step, SURE needs to determine the Top- $x$ % riskiest program statements as breakpoints. In this RQ, we investigate the effectiveness of SURE when the breakpoint determination threshold is set from 0% to 100% with 10% increments. Notice that in the training and deployment phases of SURE, feature extraction networks are necessary. Though RQ2 will investigate the selection of such networks, we must adopt a feature extraction network in this RQ to make the model run. To avoid bias, we will investigate this research question based on three feature extraction networks, i.e., VGG-16, AlexNet, and ResNet-18, separately. The results are given in Table 7 (in simulated environments) and Table 8 (in real-world environments).

**6.1.1 In simulated environments.** Let us first focus on the impact of the breakpoint determination threshold on the effectiveness of SURE, based on VGG-16 (i.e., the row “VGG-16” of Table 7). In this table, each cell denotes the value of a specific metric when using a variant of SURE. For example, the cell (“VGG-16”, “10%”,  $V_{equal}^T$ ) is 161, indicating that when SURE takes the Top-10% riskiest

Table 8. Performance of SURE on Defects4J

		10%	20%	30%	40%	50%	60%	70%	80%	90%	100%
VGG-16	$V_{equal}^T$	35	37	31	30	31	37	36	28	32	31
	$S_{M_{FMI}}^T$	<b>34.91</b>	36.84	30.95	29.90	30.91	36.85	35.86	27.86	31.85	30.94
	$S_{M_{JC}}^T$	<b>34.83</b>	36.69	30.90	29.81	30.84	36.72	35.73	27.74	31.73	30.88
	$S_{M_{PR}}^T$	<b>33.03</b>	34.7	30.00	28.72	29.52	35.02	33.68	25.99	29.81	30.04
	$S_{M_{RR}}^T$	<b>33.05</b>	34.41	29.92	28.58	29.37	34.67	33.51	25.85	29.48	29.95
AlexNet	$V_{equal}^T$	41	40	32	34	31	38	34	32	33	32
	$S_{M_{FMI}}^T$	<b>40.81</b>	39.90	31.92	33.96	30.87	37.80	33.90	31.97	32.91	31.89
	$S_{M_{JC}}^T$	<b>40.63</b>	39.80	31.83	33.92	30.75	37.63	33.82	31.94	32.84	31.79
	$S_{M_{PR}}^T$	<b>37.64</b>	37.66	30.35	33.25	28.83	34.83	32.77	31.25	30.91	30.21
	$S_{M_{RR}}^T$	<b>37.44</b>	37.58	30.35	33.17	28.58	34.58	32.61	31.25	31.01	30.10
ResNet-18	$V_{equal}^T$	37	31	33	36	34	33	33	29	33	33
	$S_{M_{FMI}}^T$	<b>36.85</b>	30.90	32.91	35.92	33.88	32.89	32.92	28.88	32.90	32.93
	$S_{M_{JC}}^T$	<b>36.71</b>	30.82	32.83	35.84	33.78	32.81	32.85	28.79	32.80	32.88
	$S_{M_{PR}}^T$	<b>34.36</b>	29.41	31.07	34.04	32.03	30.83	31.75	27.04	31.16	32.04
	$S_{M_{RR}}^T$	<b>34.20</b>	29.44	30.93	34.01	31.92	30.80	31.51	26.98	31.08	31.93

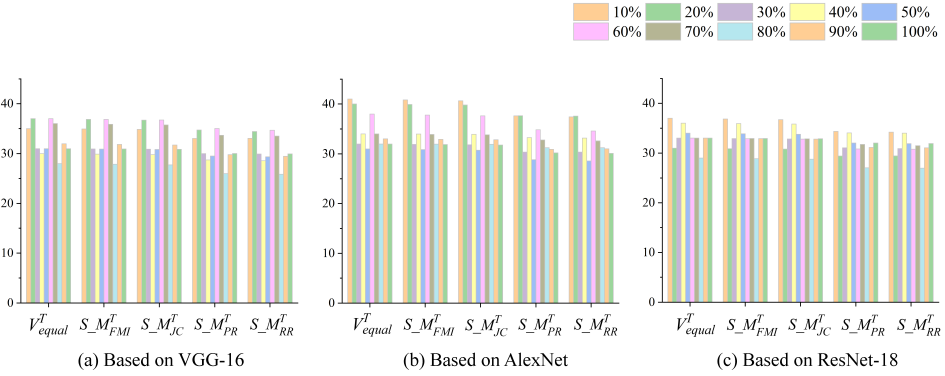


Fig. 5. The impact of the breakpoint determination threshold on SURE on Defects4J

statements as breakpoints and adopts VGG-16 as the feature extraction network (SURE in such a configuration can be referred to as  $SURE_{VGG-16}^{10\%}$ , similarly below), it can correctly predict the number of faults on 161 faulty versions. Similarly, when  $T$  takes  $SURE_{VGG-16}^{10\%}$ , the values of  $S_{M_{FMI}}^T$ ,  $S_{M_{JC}}^T$ ,  $S_{M_{PR}}^T$ , and  $S_{M_{RR}}^T$  are 126.78, 106.39, 120.24, and 80.67, respectively. We can find that when we use VGG-16 as the feature extraction network, it does not matter much if the breakpoint determination threshold takes 10%, 20%, ..., or 100%, because all the five metrics seem to be stable regardless of how many breakpoints are set. We illustrate this trend (i.e., the row “VGG-16” of Table 7) in Figure 4(a).

A similar trend can be found when AlexNet and ResNet-18 serve as the feature extraction network. Specifically, if we adopt AlexNet or ResNet-18 as the feature extraction network (i.e., the row “AlexNet” or the row “ResNet-18” of Table 7), there is no explicit trend toward change in the five metrics with the increase of  $x$  either, as shown in Figure 4(b) and Figure 4(c), respectively.

**6.1.2 In real-world environments.** Likewise, let us first use VGG-16 as the feature extraction network, to investigate the impact of the breakpoint determination threshold on the effectiveness of SURE. The results are reported in the row “VGG-16” of Table 8. Similar to the trend observed in simulated environments, we can find that the effectiveness of SURE is likely to be stable, regardless of how many breakpoints are set according to the suspiciousness of program statements. For example, when the threshold is set from 10% to 100% with 10% increments, the value of  $V_{equal}^T$  is 35, 37, 31, 30, 31, 37, 36, 28, 32, and 31. And the values of  $S_{FMI}^T$ ,  $S_{JC}^T$ ,  $S_{PR}^T$ , and  $S_{RR}^T$  seem to be not largely influenced by the breakpoint determination threshold either. We depict this trend in Figure 5(a). When the feature extraction model is replaced by AlexNet and ResNet-18, this conclusion can still be maintained, as shown in the row “AlexNet” and the row “ResNet-18” of Table 8, as well as Figure 5(b) and Figure 5(c), respectively.

The result of RQ1 can be explained to mean that the memory information collected at the Top-10% riskiest statements is already sufficient to represent failures. This result confirms our intuition of the breakpoint determination phase (Section 4.1), i.e., statements with higher risk values are more likely to be faulty, and run-time memory information gathered at these positions could have a stronger capability to reveal faults, and thus can contribute more to representing failures. Moreover, the result of RQ1 also reveals the necessity of the breakpoint determination component of SURE. Specifically, more memory information cannot contribute more to failure representation, thus we can first use SBFL to determine a few breakpoints, and then collect memory information within a small range. This is practically important because the overhead of running SBFL is less than that of collecting run-time memory information. To put it another way, we preliminarily use a lower-cost technique to avoid the following procedure with relatively higher costs, while such a procedure could not bring gains for our approach. Therefore, **in the following RQs, the investigation will be based on the breakpoint determination threshold of Top-10%.**

## 6.2 RQ2: The adopted feature extraction network.

In addition to the breakpoint determination threshold, another configurable hyperparameter of SURE is the feature extraction network. In this RQ, we select three deep convolutional neural networks (CNNs), VGG-16, AlexNet, and ResNet-18, to investigate how the effectiveness of SURE will be influenced when adopting different networks. The three investigated networks are all classical and commonly-used. Specifically, VGG-16 was introduced by the Visual Geometry Group (VGG) from the University of Oxford in 2014 [66], it is part of the VGG network series and is named “16” because it consists of 16 layers. VGG-16 is known for its simple yet effective architecture, which allows it to learn rich and complex features from images. Also, VGG-16’s parameter count is relatively large thus it requires more computational resources. AlexNet is a pioneering deep CNN introduced by Krizhevsky et al. in 2012 [46], which marks a significant breakthrough in computer vision. AlexNet popularized the use of ReLU activation functions, which helps alleviate the vanishing gradient problem and accelerate training. ResNet-18 is a variant of the residual neural network (ResNet) architecture. It was introduced as part of the ResNet family by He et al. in 2015 [29], which uses residual blocks to address the vanishing gradient problem that often occurs in deep neural networks. In this RQ, we analyze how the effectiveness of SURE will be influenced when adopting VGG-16, AlexNet, and ResNet, based on the run-time memory information collected at the Top-10% riskiest statement.

**6.2.1 In simulated environments.** We first analyze the results in simulated environments, as shown in the column “10%” of Table 7 (This column is in bold). We can find that there are 167 “ $k = r$ ” faulty versions when using AlexNet as the feature extraction network, while 161 and 118 ones when using VGG-16 and ResNet-18, respectively. As for the metrics  $S_{FMI}^T$ ,  $S_{JC}^T$ ,  $S_{PR}^T$ , and  $S_{RR}^T$ ,



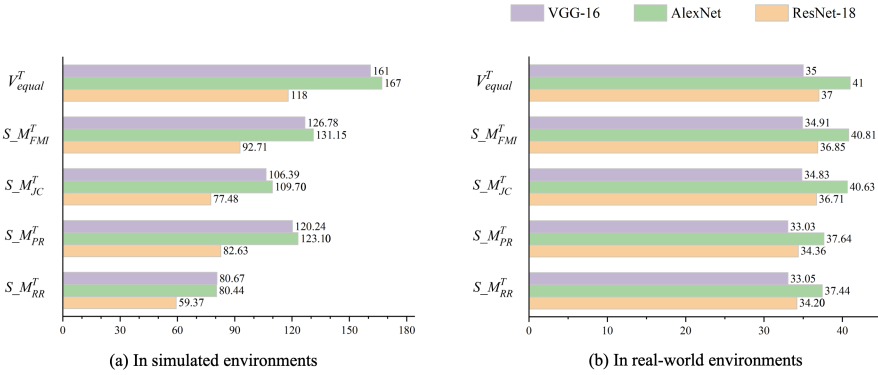


Fig. 6. The impact of the feature extraction network on SURE

AlexNet also dominates VGG-16 and ResNet-18 (with only one exception where  $S_{M_{RR}}^{SURE^{10\%}_{VGG-16}}$ : 80.67 is slightly higher than  $S_{M_{RR}}^{SURE^{10\%}_{AlexNet}}$ : 80.44). We illustrate these results in Figure 6(a) as well for clear understanding. We can find that in simulated environments, AlexNet can better extract features for PMS images thus deliver more promising failure indexing outcomes.

**6.2.2 In real-world environments.** Now we navigate to real-world environments, the results are given in the column “10%” of Table 8 (This column is in bold). We can find that there are 41 “ $k = r$ ” faulty versions when using AlexNet as the feature extraction network, while 35 and 37 ones when using VGG-16 and ResNet-18, respectively. As for the metrics  $S_{M_{FMI}}^T$ ,  $S_{M_{JC}}^T$ ,  $S_{M_{PR}}^T$ , and  $S_{M_{RR}}^T$ , AlexNet also dominates VGG-16 and ResNet-18. We illustrate these results in Figure 6(b) as well for clear understanding. Similar to the conclusion in simulated environments, we can find that AlexNet is more suitable for extracting features for PMS images.

A possible explanation for the promise of AlexNet in our experiments could be its relatively shallow architecture. Specifically, PMS images are typically in a simple form (not too many complex features involved in PMS images). It can be naturally conjectured that a shallow CNN network might be enough to extract their features. Among the selected three networks, AlexNet consists of 8 layers, including 5 convolutional layers and 3 fully connected layers, while VGG-16 and ResNet-18 both have much deeper architectures than AlexNet. Thus, it is intuitive that SURE is more suitable to configure AlexNet. Therefore, **in the following RQs, the investigation will be based on the feature extraction network of AlexNet.**

### 6.3 RQ3: The competitiveness of SURE.

In this RQ, we investigate to which extent SURE exceeds the existing most prevalent and advanced failure indexing techniques. Given the conclusions of RQ1 and RQ2, we configure 10% as the breakpoint determination threshold and AlexNet as the feature extraction network for SURE in this RQ. To put it another way, we will use  $SURE_{AlexNet}^{10\%}$  to compare with baseline techniques (for convenience, we use the term “SURE” to denote “ $SURE_{AlexNet}^{10\%}$ ” hereafter).

**6.3.1 In simulated environments.** The results in simulated environments are given in Table 9 and Figure 7(a). It can be observed that SURE outperforms all the baseline techniques substantially. Specifically, regarding the capability of faults number estimation, SURE can correctly predict the number of faults on 167 faulty versions on SIR, with 421.88%, 209.26%, 101.20%, and 30.47%

Table 9. Competitiveness of SURE in simulated environments

T	$V_{equal}^T$	$S\_M_{FMI}^T$	$S\_M_{JC}^T$	$S\_M_{PR}^T$	$S\_M_{RR}^T$
<i>Cov<sub>hit</sub></i>	32	25.31	21.34	23.18	14.92
<i>Cov<sub>count</sub></i>	54	43.33	36.75	42.37	31.30
<i>MSeer</i>	83	62.86	51.72	63.08	39.26
<i>MSeer<sub>GP19</sub></i>	128	99.98	84.64	104.09	73.38
<b>SURE</b>	<b>167</b>	<b>131.15</b>	<b>109.70</b>	<b>123.10</b>	<b>80.44</b>

Table 10. Competitiveness of SURE in real-world environments

T	$V_{equal}^T$	$S\_M_{FMI}^T$	$S\_M_{JC}^T$	$S\_M_{PR}^T$	$S\_M_{RR}^T$
<i>Cov<sub>hit</sub></i>	20	20.00	20.00	20.00	20.00
<i>Cov<sub>count</sub></i>	24	23.98	23.96	23.50	23.50
<i>MSeer</i>	29	28.99	28.98	28.75	28.75
<i>MSeer<sub>GP19</sub></i>	32	31.98	31.95	31.25	31.25
<b>SURE</b>	<b>41</b>	<b>40.81</b>	<b>40.63</b>	<b>37.64</b>	<b>37.44</b>

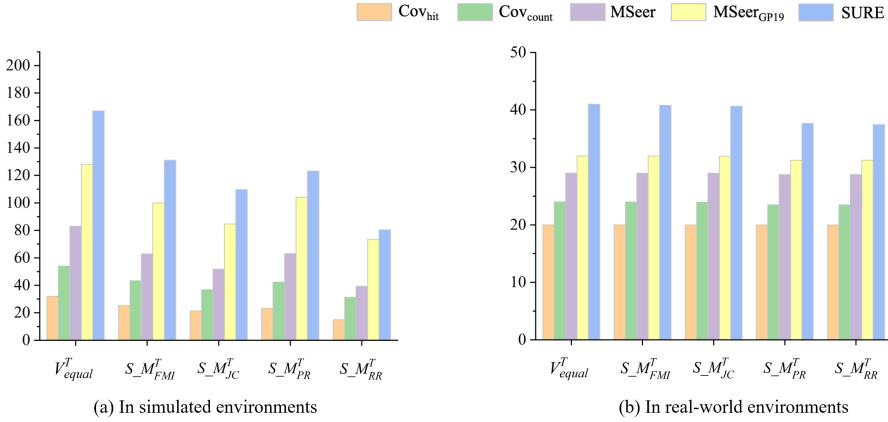


Fig. 7. Competitiveness of SURE

improvements compared with *Cov<sub>hit</sub>* (32), *Cov<sub>count</sub>* (54), *MSeer* (83), and *MSeer<sub>GP19</sub>* (128), respectively. *SURE* consistently exceeds the baselines on all four clustering metrics. For instance, if we focus on the comparison between *SURE* and *Cov<sub>hit</sub>*, improvements are 418.17%, 414.06%, 431.06%, and 439.14%, regarding  $S\_M_{FMI}^T$ ,  $S\_M_{JC}^T$ ,  $S\_M_{PR}^T$ , and  $S\_M_{RR}^T$ , respectively. Considering the four clustering metrics globally, the average improvement of *SURE* over *Cov<sub>hit</sub>* is 425.61%. Similarly, in the contexts of comparing *SURE* with *Cov<sub>count</sub>*, *MSeer*, and *MSeer<sub>GP19</sub>*, the average improvements can be calculated as 187.18% and 105.20%, and 22.17% respectively.

**6.3.2 In real-world environments.** The results in simulated environments are given in Table 10 and Figure 7(b). It can be observed that *SURE* outperforms all the baseline techniques substantially. Specifically, regarding the capability of faults number estimation, *SURE* can correctly predict the number of faults on 41 faulty versions on Defects4J, with 105.00%, 70.83%, 41.38%, and 28.13%

improvements compared with  $Cov_{hit}$  (20),  $Cov_{count}$  (24),  $MSeer$  (29), and  $MSeer_{GP19}$  (32), respectively. *SURE* consistently exceeds the baselines on all four clustering metrics. For instance, if we focus on the comparison between *SURE* and  $Cov_{hit}$ , improvements are 104.05%, 103.15%, 88.20%, and 87.20%, regarding  $S_{FMI}^T$ ,  $S_{JC}^T$ ,  $S_{PR}^T$ , and  $S_{RR}^T$ , respectively. Considering the four clustering metrics globally, the average improvement of *SURE* over  $Cov_{hit}$  is 95.65%. Similarly, in the contexts of comparing *SURE* with  $Cov_{count}$ ,  $MSeer$ , and  $MSeer_{GP19}$ , the average improvements can be calculated as 64.81%, 35.53%, and 23.76%, respectively.

As a reminder,  $MSeer$  is the state-of-the-art technique in the field of failure indexing.  $MSeer_{GP19}$  is not an existing published failure indexing technique, it is manually created by us to further evaluate the competitiveness of *SURE*.

#### 6.4 RQ4: The comprehensibility of *SURE*.

In Section 1, we point out that one of the two longstanding challenges in failure indexing is the lack of comprehensibility, that is, the result of failure indexing is hard to comprehend for human developers, which could prevent the result from being applied. The comprehensibility of a failure indexing technique essentially lies in the comprehensibility of the failure proximity, because it is the failure proximity that represents failures and hence provides human developers with the evidence of failure indexing. Here we describe the comprehensibility of failure indexing techniques as follows:

*A failure indexing technique is considered to have strong comprehensibility, if given the failure indexing result as well as the evidence (i.e., the proxies for failed test cases) for how this result was obtained, **developers can be immediately convinced by the result based on the evidence.***

But this description cannot directly guide the investigation of RQ4, because it is subjective for developers to determine whether they are convinced by the result or not, and it is difficult to quantitatively measure the extent to which they are convinced. To tackle this problem, we propose the following strategy to quantify the comprehensibility of a failure indexing technique:

*Providing developers with only the evidence (i.e., the proxies for failed test cases) of a failure indexing process, and letting them manually cluster failed test cases according to the evidence. **The closer their manual clustering outcomes are to the failure indexing results, the better they comprehend the evidence, and thus, this failure indexing technique has stronger comprehensibility.***

The intuition behind this strategy can be described in Figure 8. Specifically, if only the evidence of a failure indexing process is provided to developers, and the result of their manual clustering based on that evidence is similar or even identical to the result of the failure indexing technique, then when the failure indexing result as well as the evidence are provided together, they are able to comprehend the evidence.

As such, the clustering effectiveness made by human developers can serve as the comprehensibility of a failure indexing technique. Hence, the evaluation metrics introduced in Section 5.4, i.e.,  $V_{equal}^T$ ,  $S_{FMI}^T$ ,  $S_{JC}^T$ ,  $S_{PR}^T$ , and  $S_{RR}^T$ , can be directly employed in this RQ.

We recruit 15 participants from Wuhan University, all of whom are graduate students in Computer Science and have at least four years of programming experience. They are requested to manually finish 9 failure indexing tasks (each task corresponds to an  $r$ -bug faulty program, where  $r = 1, 2, 3, 4, \text{ or } 5$ ). Specifically, in each task, only a series of failed test cases are provided to participants.

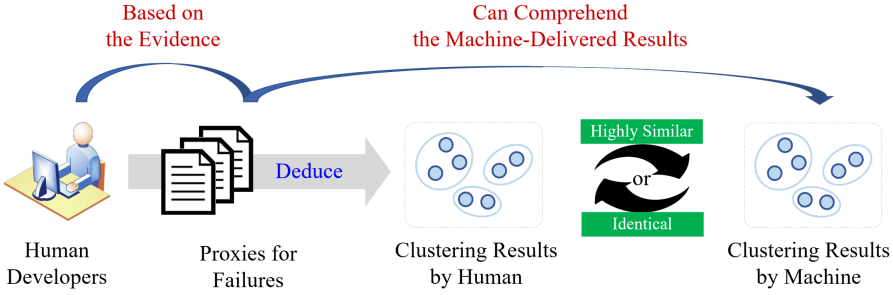


Fig. 8. Quantitatively measure the comprehensibility of failure indexing techniques

Table 11. Nine failure indexing tasks

Project	Language	# Failures	# Faults	Failure representation for developers		
				PMS	Ranking list*	Coverage <sup>†</sup>
Task-1	C	10	1	1, 4, 7, 10, 13	2, 5, 8, 11, 14	3, 6, 9, 12, 15
Task-2	C	29	2	1, 4, 7, 10, 13	2, 5, 8, 11, 14	3, 6, 9, 12, 15
Task-3	C	10	3	1, 4, 7, 10, 13	2, 5, 8, 11, 14	3, 6, 9, 12, 15
Task-4	C	19	4	3, 6, 9, 12, 15	1, 4, 7, 10, 13	2, 5, 8, 11, 14
Task-5	C	17	5	3, 6, 9, 12, 15	1, 4, 7, 10, 13	2, 5, 8, 11, 14
Task-6	Java	6	2	3, 6, 9, 12, 15	1, 4, 7, 10, 13	2, 5, 8, 11, 14
Task-7	Java	3	3	2, 5, 8, 11, 14	3, 6, 9, 12, 15	1, 4, 7, 10, 13
Task-8	Java	4	4	2, 5, 8, 11, 14	3, 6, 9, 12, 15	1, 4, 7, 10, 13
Task-9	Java	5	5	2, 5, 8, 11, 14	3, 6, 9, 12, 15	1, 4, 7, 10, 13

\* "Ranking list" denotes "Suspiciousness ranking lists by GP19". The same in Table 12.

<sup>†</sup> "Coverage" denotes "Coverage vectors of execution frequency". The same in Table 12.

These failures are triggered by how many underlying fault(s), as well as the mutual relationship among these failures, are hidden from participants. Participants will manually divide all failures through two steps: 1) Estimating these failures should be divided into how many clusters, i.e., there are how many underlying faults behind these failures. 2) Manually dividing failures according to the estimated number of clusters.

As for each task, its failed test cases are given in one of three forms of representation: program memory spectrum (PMS) images, suspiciousness ranking lists by GP19, and coverage vectors of execution frequency. Specifically, PMS images are the new failure representation proposed in this paper, suspiciousness ranking lists by GP19 are the failure representation of  $MSeer_{GP19}$ , and coverage vectors of execution frequency are the failure representation of  $Cov_{count}$ . We do not investigate suspiciousness ranking lists by Crosstab (the failure representation of  $MSeer$ ), and coverage vectors of binary indicators (the failure representation of  $Cov_{hit}$ ). This is because  $MSeer$  and  $MSeer_{GP19}$  both belong to the SD-based failure proximity, while  $MSeer_{GP19}$  has been demonstrated to be better than  $MSeer$  in Section 6.3. Also,  $Cov_{hit}$  and  $Cov_{count}$  both belong to the CC-based failure proximity, while  $Cov_{count}$  has been demonstrated to be better than  $Cov_{hit}$  in Section 6.3.

We concisely introduce the 9 tasks in Table 11, including the programming language, the number of observed failures in testing, as well as the number of underlying faults. **From the perspective of tasks**, each task will be carried out by all of the 15 participants, among them, five are based on

Table 12. Comprehensibility of three forms of failure representation

		Participant															Total
		No. 1	No. 2	No. 3	No. 4	No. 5	No. 6	No. 7	No. 8	No. 9	No. 10	No. 11	No. 12	No. 13	No. 14	No. 15	
PMS	$V_{equal}^T$	2	3	2	3	<b>3</b>	3	3	3	3	2	3	2	1	3	3	39
	$S_{FMI}^T$	2.00	3.00	2.00	3.00	<b>3.00</b>	3.00	3.00	3.00	3.00	2.00	3.00	2.00	1.00	3.00	3.00	39.00
	$S_{JC}^T$	2.00	3.00	2.00	3.00	<b>3.00</b>	3.00	3.00	3.00	3.00	2.00	3.00	2.00	1.00	3.00	3.00	39.00
	$S_{PR}^T$	2.00	3.00	2.00	3.00	<b>3.00</b>	3.00	3.00	3.00	3.00	2.00	3.00	1.96	1.00	3.00	3.00	38.96
	$S_{RR}^T$	2.00	3.00	2.00	3.00	<b>3.00</b>	3.00	3.00	3.00	3.00	2.00	3.00	1.93	1.00	3.00	3.00	38.93
	Average Time (s)	5.41	3.08	1.50	6.22	<b>1.17</b>	3.24	1.39	3.42	4.43	5.63	3.11	2.62	3.02	2.33	2.33	3.26
Ranking list	$V_{equal}^T$	2	0	3	2	<b>0</b>	1	0	1	2	2	1	1	1	0	3	19
	$S_{FMI}^T$	1.99	0.00	3.00	1.99	<b>0.00</b>	1.00	0.00	0.99	2.00	1.99	0.99	1.00	1.00	0.00	3.00	18.95
	$S_{JC}^T$	1.99	0.00	3.00	1.99	<b>0.00</b>	1.00	0.00	0.98	2.00	1.99	0.98	1.00	1.00	0.00	3.00	18.93
	$S_{PR}^T$	1.81	0.00	3.00	1.62	<b>0.00</b>	1.00	0.00	0.46	2.00	1.73	0.55	1.00	1.00	0.00	3.00	17.17
	$S_{RR}^T$	1.81	0.00	3.00	1.66	<b>0.00</b>	1.00	0.00	0.40	2.00	1.71	0.46	1.00	1.00	0.00	3.00	17.04
	Average Time (s)	14.33	25.37	18.50	23.21	<b>34.92</b>	21.08	27.38	44.98	13.92	13.11	20.63	16.83	4.69	33.35	6.33	21.24
Coverage	$V_{equal}^T$	0	2	1	1	<b>1</b>	0	2	2	0	2	2	1	3	2	0	19
	$S_{FMI}^T$	0.00	1.99	1.00	1.00	<b>1.00</b>	0.00	2.00	1.99	0.00	2.00	1.99	0.99	3.00	2.00	0.00	18.96
	$S_{JC}^T$	0.00	1.99	1.00	1.00	<b>1.00</b>	0.00	2.00	1.99	0.00	2.00	1.98	0.99	3.00	2.00	0.00	18.95
	$S_{PR}^T$	0.00	1.72	1.00	1.00	<b>1.00</b>	0.00	2.00	1.72	0.00	2.00	1.62	0.48	3.00	1.87	0.00	17.41
	$S_{RR}^T$	0.00	1.66	1.00	1.00	<b>1.00</b>	0.00	2.00	1.72	0.00	2.00	1.66	0.47	3.00	1.87	0.00	17.38
	Average Time (s)	12.00	16.86	38.59	20.17	<b>14.26</b>	46.63	20.33	35.40	34.31	24.72	14.78	21.90	3.33	12.93	12.71	21.93

PMS, five are based on suspiciousness ranking lists by GP19, and five are based on coverage vectors of execution frequency. For example, we can find that for Task-1, Task-2, and Task-3, Participants 1, 4, 7, 10, 13 will be based on PMS to perform failure indexing, Participants 2, 5, 8, 11, 14 will be based on suspiciousness ranking lists by GP19, and Participants 3, 6, 9, 12, 15 will be based on coverage vectors of execution frequency. Additionally, **from the perspective of participants**, each developer will handle all of the 9 tasks, among them, three for each of the three forms of representation. For example, we can find that Participant-1 handles Task-1, Task-2, and Task-3 based on PMS, handles Task-4, Task-5, and Task-6 based on suspiciousness ranking lists by GP19, and handles Task-7, Task-8, and Task-9 based on coverage vectors of execution frequency, as shown by the bold numbers in Table 11. As such, each form of representation will be carried out for 45 times, as shown in the columns “PMS”, “Ranking list”, and “Coverage” in Table 11.

The performance of human developers is given in Table 12<sup>12</sup>. Take Participant-5 as an example, as shown in the column “Participant-5” in Table 12. Participant-5 carries out 9 failure indexing tasks with 3 PMS-based, 3 Ranking list-based, and 3 Coverage-based ones. As for the 3 PMS-based faulty versions, Participant-5 correctly estimate the number of faults on all of them, i.e., the value of the cell (“PMS”, “ $V_{equal}^T$ ”, “Participant-5”) is 3. The values of the four clustering metrics, i.e.,  $S_{FMI}^T$ ,  $S_{JC}^T$ ,  $S_{PR}^T$ ,  $S_{RR}^T$ , are both 3.00. Considering that the maximum values of FMI, JC, PR, and RR are 1, this result indicates that Participant-5 can achieve perfect failure indexing based on PMS. But if based on the other forms of failure representation, Participant-5 performs much worse. Specifically, “ $k == r$ ” cannot be obtained on any faulty version based on the representation of ranking lists, and only one faulty version can be properly handled based on the representation of coverage. Moreover, the time cost of manual failure indexing is also different among the three forms of representation. For example, Participant-5 takes 1.17s on average to index each failure when it is represented as PMS images, while takes 34.92s and 14.26s on average when failures are represented as ranking lists and coverage, respectively. That is to say, the proposed failure representation of PMS images enables developers to identify failures with the same/different root cause(s) at a glance, further demonstrating its intuition and friendliness towards human developers.

<sup>12</sup>For a faulty version, consistent with the strategy in Section 5.4, we analyze the clustering effectiveness only when “ $k == r$ ”, namely, the number of faults is correctly predicted.



be fully dominated by others, indicating a potential future direction of combining the advantages of different classes of failure proximities together. Despite the respective superiority of various techniques, SURE can solely handle more faulty versions than the other four ones. In particular, it exceeds 76.92% and 225.00% by  $MSeer_{GP19}$ , the best-performed baseline technique, in simulated and real-world environments, respectively.

## 7.2 Overhead of SURE

The time costs of SURE mainly involve three phases, namely, the failure representation, the distance measurement (including the model training and the deployment), and the clustering. According to our analyses, as for the failure representation, SURE takes 3.99 minutes and 5.90 minutes on average to collect run-time memory information during the execution of a failed test case, and then takes 0.18s and 0.27s on average to convert a set of memory information to a PMS image, on SIR and Defects4J faulty versions, respectively. As for the distance measurement, SURE takes 23 minutes and 35 seconds to train the Siamese-based model. Once the model is trained, 0.01s and 0.06s on average are taken to predict the distance between a pair of PMS images on SIR and Defects4J faulty versions, respectively. After these two steps are ready, the clustering process typically takes only a few seconds. To summarize, the overhead of SURE mainly lies in querying run-time memory information at preset breakpoints.

As pioneers have pointed out, failure indexing is essential yet very costly when it comes to manual debugging in the real world. “*Experienced developers can manually examine every failure and determine the culprit fault, but this is apparently too expensive*”, pioneers concluded [55]. In contrast, SURE can complete this task automatically and hence can save a great deal of manual effort, which is far less expensive than manual jobs. In particular, SURE represents a failed test case as a PMS image, which can be smoothly comprehended by human developers at a glance and thus can boost the comprehensibility of failure indexing outcomes. This can be of great importance because if developers are not convinced by failure indexing outcomes, they will take much time to verify their correctness. Even though coverage vectors (the failure representation of the CC-based strategy) and statement ranking lists (the failure representation of the SD-based strategy) can also provide insights for developers to comprehend failure indexing outcomes to an extent, they can be human labor-intensive. To put it another way, we transfer the costs that would have been borne by humans to machines. It is true that as compared with CC and SD-based techniques, SURE needs higher costs. However, in return, SURE delivers better performance and convenience. Therefore, the cost of SURE is acceptable: more sophisticated fingerprinting is naturally accompanied by higher overhead, i.e., “*no free lunch*” [55].

## 8 THREATS TO VALIDITY

Our experiments are subject to several threats to validity.

Threats to internal validity relate to the faulty versions we focus on to determine the effectiveness of a failure indexing technique. In the experimental evaluation, we only consider those faulty versions that satisfy the condition of “ $k == r$ ”, i.e., where the predicted number of faults  $k$  is equal to the real number of faults  $r$ . Actually, “ $k == r$ ” is ideal but not necessary for parallel debugging to work in practice. This is because if  $k$  exceeds  $r$ ,  $k$  developers will be employed to locate  $r$  faults, and parallel debugging can work at the cost of human labor ( $k - r$  developers are redundant). On the contrary, if  $k$  is less than  $r$ , parallel debugging can also work at the cost of more than one iteration of debugging. Filtering these faulty versions is because comparing  $k$  generated clusters with  $r$  oracle clusters when “ $k != r$ ” is difficult, and moreover, prior studies have pointed out that the performance of a failure indexing technique can be mainly determined by those “ $k == r$ ” faulty versions [26, 69], that is, the contribution of “ $k != r$ ” ones is marginal. Therefore, the threat is

acceptable. Nonetheless, we plan to further take this imperfect situation into account in our future work, for more comprehensive evaluation of failure indexing techniques.

Threats to external validity relate to the generalization capability of SURE. Specifically, given that the evaluation of the effectiveness of SURE is carried out empirically, our results may not be extended to all programs. In the experiments, we use nine projects from two open-source platforms, which are written in different languages (C and Java) and with various functionalities, thus mitigating the threat to an extent. In particular, we use only 30% of the simulated faults to train the model, and use the remaining 70% of the simulated as well as real-world faults to test. This allows us to have higher confidence with respect to the applicability of SURE.

Threats to construct validity relate to the adopted evaluation metrics. We use external metrics (i.e., FMI, JC, PR, and RR) to measure the clustering effectiveness. Despite the fact that these metrics have been extensively used by previous works, we consider using more diverse metrics in our future work for further evaluation.

## 9 RELATED WORK

It is well-recognized that the failure proximity essentially underpins failure indexing techniques. Liu et al. systematically summarized the status of the community of failure proximities, pointing out that CC-based and SD-based strategies can deliver good results [55]. In fact, these two are mainstream tactics in failure indexing to date, and numerous studies have emerged around them in the last two decades.

Let us first focus on the CC-based strategy. As a very early work, Podgurski et al. suggested using execution profiles, including code coverage, as the signature of failures [62]. Since then, this tactic has become popular gradually. For example, Liu et al. formally presented the definition of the CC-based failure proximity, and demonstrated its capability by experiments [55]. Högerle et al. constructed failed test coverage matrices, and used the Weil-Kettler algorithm to rearrange the mentioned coverage matrix into a block diagonal matrix thus achieving failure clustering [31]. Also based on failed test coverage matrices, Steimann and Frenkel utilized two partitioning procedures from integer linear programming to cluster failures, for breaking down the fault localization problem into smaller ones [71]. Huang et al. were concerned with the impact of failure indexing on the effectiveness of multi-fault localization. They conducted an empirical study to explore this topic on the basis of taking code coverage as the representation of failures [33]. Wu et al. also employed code coverage as the signature of failures to perform failure clustering. They iteratively chose the cluster of failures with the highest density to conduct single-fault localization, until all faults were fixed [87].

Later, the SD-based strategy, a more sophisticated solution of failure proximity, was proposed by Liu and Han [53], which regards two failures as similar if they suggest roughly the same fault location. Specifically, they used SOBER, a statistical model-based fault localization technique at the predicate granularity [54], to complete the fault location suggestion process. From then on, the SD-based strategy attracted more and more attention from academia. For example, Jones et al. used Tarantula [37], an SBFL technique working at the statement granularity, to suggest finer-grained fault location. As a result, the suspiciousness ranking list that reflects the possibility of each program statement being faulty is utilized to represent failures [36]. Following the workflow of the SD-based strategy, Gao and Wong proposed MSeer, which employed Crosstab [84], an empirically promising SBFL technique, to produce the suspiciousness ranking list that represents failures [26]. MSeer is also the state-of-the-art technique to date in the field of failure indexing. Song et al. investigated the impact of the adopted SBFL techniques on the effectiveness of SD-based failure indexing. Their conclusion demonstrated that MSeer can be further enhanced by replacing Crosstab with other SBFL formulas, such as the one evolved by genetic programming, GP19 [69].



Despite the integration of fault localization techniques, the resource on which the SD-based strategy relies is still only coverage. If failures that are triggered by distinct faults have the same coverage, neither the CC-based nor the SD-based strategy can work well. However, this situation has been demonstrated to be very common in practice, which causes performance degradation in CC and SD strategies. Moreover, the failure representation (code coverage vectors for CC and suspiciousness ranking lists for SD) is also tough and time-consuming to comprehend by human developers, which hinders human comprehension of failure indexing results. In this paper, we propose SURE, a novel failure indexing technique, to tackle the aforementioned threats to existing approaches. With the support of run-time memory information, SURE remarkably improves the effectiveness of failure indexing, and PMS images, the novel form of failure representation adopted by SURE, help human developers comprehend the result of failure indexing better.

There are some recent works that introduce external profiles to support failure indexing, such as code-independent features in regression testing [28], as well as code features and historical features in continuous integration [4]. We do not consider such types of studies since they go beyond our research scope: 1) their source information cannot be always available, and 2) This paper focuses on failure indexing in the context of multi-fault debugging.

## 10 CONCLUSION

In this paper, we propose the program memory-based failure proximity, and based on that propose a novel failure indexing technique, SURE. Experimental results demonstrate the high competitiveness of SURE: it can achieve 101.20% and 41.38% improvements in faults number estimation, as well as 105.20% and 35.53% improvements in clustering, compared with the state-of-the-art technique in this field to date, in simulated and real-world environments, respectively. SURE also provides human developers with insights to better comprehend the failure indexing results. Our human study involving 15 participants shows that the form of failure representation of SURE, i.e., PMS images, can boost human developers' comprehension of failure indexing results by 105.26% compared with the two most prevalent and advanced strategies, while the time cost is reduced by 84.65% and 85.13%, respectively.

In the future, we plan to dig deeper into the characteristics of various forms of failure proximities, and try to combine their advantages for better extracting the signature of failures. A more extensive experiment with a larger scale of benchmarks as well as more diverse evaluation metrics is also being considered.

## REFERENCES

- [1] Rui Abreu, Peter Zoetewij, and Arjan JC Van Gemund. 2009. Spectrum-based multiple fault localization. In *2009 IEEE/ACM International Conference on Automated Software Engineering*. IEEE, 88–99.
- [2] Hiralal Agrawal, Joseph R Horgan, Saul London, and W Eric Wong. 1995. Fault localization using execution slices and dataflow tests. In *Proceedings of Sixth International Symposium on Software Reliability Engineering. ISSRE'95*. IEEE, 143–151.
- [3] Higor Amario de Souza, Marcelo de Souza Lauretto, Fabio Kon, and Marcos Lordello Chaim. 2023. Understanding the use of spectrum-based fault localization. *Journal of Software: Evolution and Process* (2023), e2622. <https://doi.org/10.1002/smr.2622>
- [4] Gabin An, Juyeon Yoon, Jeongju Sohn, Jingun Hong, Dongwon Hwang, and Shin Yoo. 2022. Automatically identifying shared root causes of test breakages in SAP HANA. In *Proceedings of the 44th International Conference on Software Engineering: Software Engineering in Practice*. 65–74.
- [5] Gabin An, Juyeon Yoon, and Shin Yoo. 2021. Searching for multi-fault programs in defects4j. In *International Symposium on Search Based Software Engineering*. Springer, 153–158.
- [6] James H Andrews, Lionel C Briand, and Yvan Labiche. 2005. Is mutation an appropriate tool for testing experiments?. In *Proceedings of the 27th international conference on Software engineering*. 402–411.

- [7] James H Andrews, Lionel C Briand, Yvan Labiche, and Akbar Siami Namin. 2006. Using mutation analysis for assessing and comparing testing coverage criteria. *IEEE Transactions on Software Engineering* 32, 8 (2006), 608–624.
- [8] Aitor Arrieta, Sergio Segura, Urtzi Markiegi, Goiuria Sagardui, and Leire Etxeberria. 2018. Spectrum-based fault localization in software product lines. *Information and Software Technology* 100 (2018), 18–31.
- [9] Arun Babu, Qingkai Shi, and Muhammad Ashfaq. 2020. Python script for performing mutation testing. Github Repository. <https://github.com/arun-babu/mutate.py>
- [10] Antonia Bertolino, Breno Miranda, Roberto Pietrantuono, and Stefano Russo. 2017. Adaptive coverage and operational profile-based testing for reliability improvement. In *2017 IEEE/ACM 39th International Conference on Software Engineering (ICSE)*. IEEE, 541–551.
- [11] Jane Bromley, Isabelle Guyon, Yann LeCun, Eduard Säckinger, and Roopak Shah. 1993. Signature verification using a "siamese" time delay neural network. *Advances in neural information processing systems* 6 (1993).
- [12] Dylan Callaghan and Bernd Fischer. 2023. Improving Spectrum-Based Localization of Multiple Faults by Iterative Test Suite Reduction. *arXiv preprint arXiv:2306.09892* (2023).
- [13] Davide Chicco. 2021. Siamese neural networks: An overview. *Artificial neural networks* (2021), 73–94.
- [14] Stephen L Chiu. 1994. Fuzzy model identification based on cluster estimation. *Journal of Intelligent & fuzzy systems* 2, 3 (1994), 267–278.
- [15] Sumit Chopra, Raia Hadsell, and Yann LeCun. 2005. Learning a similarity metric discriminatively, with application to face verification. In *2005 IEEE computer society conference on computer vision and pattern recognition (CVPR'05)*, Vol. 1. IEEE, 539–546.
- [16] William Dickinson, David Leon, and A Fodgurski. 2001. Finding failures by cluster analysis of execution profiles. In *Proceedings of the 23rd International Conference on Software Engineering, ICSE 2001*. IEEE, 339–348.
- [17] Nicholas DiGiuseppe and James A Jones. 2011. On the influence of multiple faults on coverage-based fault localization. In *Proceedings of the 2011 international symposium on software testing and analysis*. 210–220.
- [18] Nicholas DiGiuseppe and James A Jones. 2012. Concept-based failure clustering. In *Proceedings of the ACM SIGSOFT 20th international symposium on the foundations of software engineering*. 1–4.
- [19] Nicholas DiGiuseppe and James A Jones. 2012. Software behavior and failure clustering: An empirical study of fault causality. In *2012 IEEE Fifth International Conference on Software Testing, Verification and Validation*. IEEE, 191–200.
- [20] Nicholas DiGiuseppe and James A Jones. 2015. Fault density, fault types, and spectra-based fault localization. *Empirical Software Engineering* 20 (2015), 928–967.
- [21] Hyunsook Do, Sebastian Elbaum, and Gregg Rothermel. 2005. Supporting controlled experimentation with testing techniques: An infrastructure and its potential impact. *Empirical Software Engineering* 10 (2005), 405–435.
- [22] Hyunsook Do and Gregg Rothermel. 2006. On the use of mutation faults in empirical assessments of test case prioritization techniques. *IEEE Transactions on Software Engineering* 32, 9 (2006), 733–752.
- [23] Hugh S Fairman, Michael H Brill, and Henry Hemmendinger. 1997. How the CIE 1931 color-matching functions were derived from Wright-Guild data. *Color Research & Application* 22, 1 (1997), 11–23.
- [24] Wei Fu and Patrick O Perry. 2020. Estimating the number of clusters using cross-validation. *Journal of Computational and Graphical Statistics* 29, 1 (2020), 162–173.
- [25] Meng Gao, Pengyu Li, Congcong Chen, and Yunsong Jiang. 2018. Research on software multiple fault localization method based on machine learning. In *MATEC web of conferences*, Vol. 232. EDP Sciences, 01060.
- [26] Ruizhi Gao and W Eric Wong. 2019. MSeer—An Advanced Technique for Locating Multiple Bugs in Parallel. *IEEE Transactions on Software Engineering* 45, 03 (2019), 301–318.
- [27] Laleh Sh Ghandehari, Yu Lei, Raghu Kacker, Richard Kuhn, Tao Xie, and David Kung. 2018. A combinatorial testing-based approach to fault localization. *IEEE Transactions on Software Engineering* 46, 6 (2018), 616–645.
- [28] Mojdeh Golagha, Constantin Lehnhoff, Alexander Pretschner, and Hermann Ilmberger. 2019. Failure clustering without coverage. In *Proceedings of the 28th ACM SIGSOFT International Symposium on Software Testing and Analysis*. 134–145.
- [29] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [30] Robert Hirsch. 2005. *Exploring colour photography: a complete guide*. Laurence King Publishing.
- [31] Wolfgang Högerle, Friedrich Steimann, and Marcus Frenkel. 2014. More debugging in parallel. In *2014 IEEE 25th International Symposium on Software Reliability Engineering*. IEEE, 133–143.
- [32] Rubing Huang, Dave Towey, Yinyin Xu, Yunan Zhou, and Ning Yang. 2022. Dissimilarity-based test case prioritization through data fusion. *Software: Practice and Experience* 52, 6 (2022), 1352–1377.
- [33] Yanqin Huang, Junhua Wu, Yang Feng, Zhenyu Chen, and Zhihong Zhao. 2013. An empirical study on clustering for isolating bugs in fault localization. In *2013 IEEE International Symposium on Software Reliability Engineering Workshops (ISSREW)*. IEEE, 138–143.

- [34] Anil K Jain, M Narasimha Murty, and Patrick J Flynn. 1999. Data clustering: a review. *ACM computing surveys (CSUR)* 31, 3 (1999), 264–323.
- [35] Dennis Jeffrey, Neelam Gupta, and Rajiv Gupta. 2008. Fault localization using value replacement. In *Proceedings of the 2008 international symposium on Software testing and analysis*. 167–178.
- [36] James A Jones, James F Bowring, and Mary Jean Harrold. 2007. Debugging in parallel. In *Proceedings of the 2007 international symposium on Software testing and analysis*. 16–26.
- [37] James A Jones and Mary Jean Harrold. 2005. Empirical evaluation of the tarantula automatic fault-localization technique. In *Proceedings of the 20th IEEE/ACM international Conference on Automated software engineering*. 273–282.
- [38] James A Jones, Mary Jean Harrold, and John Stasko. 2002. Visualization of test information to assist fault localization. In *Proceedings of the 24th international conference on Software engineering*. 467–477.
- [39] René Just, Darioush Jalali, and Michael D Ernst. 2014. Defects4J: A database of existing faults to enable controlled testing studies for Java programs. In *Proceedings of the 2014 international symposium on software testing and analysis*. 437–440.
- [40] René Just, Darioush Jalali, Laura Inozemtseva, Michael D Ernst, Reid Holmes, and Gordon Fraser. 2014. Are mutants a valid substitute for real faults in software testing?. In *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering*. 654–665.
- [41] Leonard Kaufman and Peter J Rousseeuw. 2009. *Finding groups in data: an introduction to cluster analysis*. John Wiley & Sons.
- [42] Fabian Keller, Lars Grunske, Simon Heiden, Antonio Filieri, Andre van Hoorn, and David Lo. 2017. A critical evaluation of spectrum-based fault localization techniques on a large-scale software system. In *2017 IEEE International Conference on Software Quality, Reliability and Security (QRS)*. IEEE, 114–125.
- [43] Maurice George Kendall. 1948. Rank correlation methods. (1948).
- [44] Yunho Kim and Shin Hong. 2022. Learning-based mutant reduction using fine-grained mutation operators. *Software Testing, Verification and Reliability* 32, 7 (2022), e1786.
- [45] Suneel Kumar Kingrani, Mark Levene, and Dell Zhang. 2018. Estimating the number of clusters using diversity. *Artificial Intelligence Research* 7, 1 (2018), 15–22.
- [46] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems* 25 (2012).
- [47] Si-Mohamed Lamraoui and Shin Nakajima. 2016. A formula-based approach for automatic fault localization of multi-fault programs. *Journal of Information Processing* 24, 1 (2016), 88–98.
- [48] Tien-Duy B Le, David Lo, and Ferdian Thung. 2015. Should i follow this fault localization tool’s output? automated prediction of fault localization effectiveness. *Empirical Software Engineering* 20 (2015), 1237–1274.
- [49] Shangru Li. 2020. *Research on fault location based on memory data*. Master’s thesis. Northwestern Polytechnical University.
- [50] Yihao Li and Pan Liu. 2022. A preliminary investigation on the performance of SBFL techniques and distance metrics in parallel fault localization. *IEEE Transactions on Reliability* 71, 2 (2022), 803–817.
- [51] Zheng Li, Yonghao Wu, Haifeng Wang, Xiang Chen, and Yong Liu. 2022. Review of Software Multiple Fault Localization Approaches. *Chinese Journal of Computers* 45, 2 (02 2022), 256–288.
- [52] Chao Liu, Long Fei, Xifeng Yan, Jiawei Han, and Samuel P Midkiff. 2006. Statistical debugging: A hypothesis testing-based approach. *IEEE Transactions on software engineering* 32, 10 (2006), 831–848.
- [53] Chao Liu and Jiawei Han. 2006. Failure proximity: a fault localization-based approach. In *Proceedings of the 14th ACM SIGSOFT international symposium on Foundations of software engineering*. 46–56.
- [54] Chao Liu, Xifeng Yan, Long Fei, Jiawei Han, and Samuel P Midkiff. 2005. SOBER: statistical model-based bug localization. *ACM SIGSOFT Software Engineering Notes* 30, 5 (2005), 286–295.
- [55] Chao Liu, Xiangyu Zhang, and Jiawei Han. 2008. A systematic study of failure proximity. *IEEE Transactions on Software Engineering* 34, 6 (2008), 826–843.
- [56] Muhammad Azhar Mushtaq, Abid Sultan, Muhammad Afrasayab, and Taimoor Zubair. 2019. New cryptographic algorithm using ASCII values and gray code (AGC). In *Proceedings of the 4th International Conference on Big Data and Computing*. 242–246.
- [57] Lee Naish, Hua Jie Lee, and Kotagiri Ramamohanarao. 2011. A model for spectra-based software diagnosis. *ACM Transactions on software engineering and methodology (TOSEM)* 20, 3 (2011), 1–32.
- [58] Farah Naz, Ijaz Ali Shoukat, Rehan Ashraf, Umer Iqbal, and Abdul Rauf. 2020. An ASCII based effective and multi-operation image encryption method. *Multimedia Tools and Applications* 79 (2020), 22107–22129.
- [59] Chris Parnin and Alessandro Orso. 2011. Are automated debugging techniques actually helping programmers?. In *Proceedings of the 2011 international symposium on software testing and analysis*. 199–209.
- [60] Spencer Pearson, José Campos, René Just, Gordon Fraser, Rui Abreu, Michael D Ernst, Deric Pang, and Benjamin Keller. 2017. Evaluating and improving fault localization. In *2017 IEEE/ACM 39th International Conference on Software*

*Engineering (ICSE)*. IEEE, 609–620.

- [61] Hanyu Pei, Beibei Yin, Min Xie, and Kai-Yuan Cai. 2021. Dynamic random testing with test case clustering and distance-based parameter adjustment. *Information and Software Technology* 131 (2021), 106470.
- [62] Andy Podgurski, David Leon, Patrick Francis, Wes Masri, Melinda Minch, Jiayang Sun, and Bin Wang. 2003. Automated support for classifying software failure reports. In *25th International Conference on Software Engineering, 2003. Proceedings*. IEEE, 465–475.
- [63] Michael Pradel and Koushik Sen. 2018. Deepbugs: A learning approach to name-based bug detection. *Proceedings of the ACM on Programming Languages* 2, OOPSLA (2018), 1–25.
- [64] Henrique L Ribeiro, Roberto PA de Araujo, Marcos L Chaim, Higor A de Souza, and Fabio Kon. 2018. Jaguar: A spectrum-based fault localization tool for real-world software. In *2018 IEEE 11th International Conference on Software Testing, Verification and Validation (ICST)*. IEEE, 404–409.
- [65] Ting Shu, Tiantian Ye, Zuohua Ding, and Jinsong Xia. 2016. Fault localization based on statement frequency. *Information Sciences* 360 (2016), 43–56.
- [66] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).
- [67] SIR. 2020. *The Software Infrastructure Repository*. Retrieved September 5, 2020 from <https://sir.csc.ncsu.edu/portal/index.php>
- [68] Ian Sommerville. 2011. *Software Engineering, 9/E*. Pearson Education India.
- [69] Yi Song, Xiaoyuan Xie, Quanming Liu, Xihao Zhang, and Xi Wu. 2022. A comprehensive empirical investigation on failure clustering in parallel debugging. *Journal of Systems and Software* 193 (2022), 111452.
- [70] Yi Song, Xiaoyuan Xie, Xihao Zhang, Quanming Liu, and Ruizhi Gao. 2022. Evolving Ranking-Based Failure Proximities for Better Clustering in Fault Isolation. In *Proceedings of the 37th IEEE/ACM International Conference on Automated Software Engineering*. 1–13.
- [71] Friedrich Steimann and Marcus Frenkel. 2012. Improving coverage-based localization of multiple faults using algorithms from integer linear programming. In *2012 IEEE 23rd International Symposium on Software Reliability Engineering*. IEEE, 121–130.
- [72] Pang-Ning Tan, Michael Steinbach, and Vipin Kumar. 2016. *Introduction to data mining*. Pearson Education India.
- [73] Robert Tibshirani, Guenther Walther, and Trevor Hastie. 2001. Estimating the number of clusters in a data set via the gap statistic. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 63, 2 (2001), 411–423.
- [74] Béla Vancsics. 2023. *New Algorithms and Benchmarks for Supporting Spectrum-Based Fault Localization*. Ph.D. Dissertation. University of Szeged.
- [75] A Vijayan, T Gobinath, and M Saravanakarhikeyan. 2016. ASCII value based encryption system (AVB). *International Journal of Engineering Research and Applications* 6, 4 (2016), 8–11.
- [76] MK Vijaymeena and K Kavitha. 2016. A survey on similarity measures in text mining. *Machine Learning and Applications: An International Journal* 3, 2 (2016), 19–28.
- [77] Jeffrey M. Voas. 1992. PIE: A dynamic failure-based technique. *IEEE Transactions on software Engineering* 18, 8 (1992), 717.
- [78] Qing Wang, Shujian Wu, and Ming-Shu Li. 2008. Software defect prediction. *Journal of software* 19, 7 (2008), 1565–1580.
- [79] Xingya Wang, Shujuan Jiang, Pengfei Gao, Kai Lu, Bo Lili, Xiaolin Ju, and Yanmei Zhang. 2019. Fuzzy C-means clustering based multi-fault localization. *Chinese Journal of Computers* 42 (2019).
- [80] Wanzhi Wen. 2012. Software fault localization based on program slicing spectrum. In *2012 34th International Conference on Software Engineering (ICSE)*. IEEE, 1511–1514.
- [81] Ratnadira Widayarsi, Gede Artha Azriadi Prana, Stefanus Agus Haryono, Shaowei Wang, and David Lo. 2022. Real world projects, real faults: evaluating spectrum based fault localization techniques on Python projects. *Empirical Software Engineering* 27, 6 (2022), 147.
- [82] W Eric Wong, Vidroha Debroy, Ruizhi Gao, and Yihao Li. 2013. The DStar method for effective software fault localization. *IEEE Transactions on Reliability* 63, 1 (2013), 290–308.
- [83] W Eric Wong, Vidroha Debroy, Richard Golden, Xiaofeng Xu, and Bhavani Thuraisingham. 2011. Effective software fault localization using an RBF neural network. *IEEE Transactions on Reliability* 61, 1 (2011), 149–169.
- [84] W Eric Wong, Vidroha Debroy, and Dianxiang Xu. 2012. Towards Better Fault Localization: A Crosstab-Based Statistical Approach. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)* 3, 42 (2012), 378–396.
- [85] W Eric Wong, Ruizhi Gao, Yihao Li, Rui Abreu, and Franz Wotawa. 2016. A survey on software fault localization. *IEEE Transactions on Software Engineering* 42, 8 (2016), 707–740.
- [86] Junjie Wu, Hui Xiong, and Jian Chen. 2009. Adapting the right measures for k-means clustering. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*. 877–886.

- [87] Yong-Hao Wu, Zheng Li, Yong Liu, and Xiang Chen. 2020. Fatoc: Bug isolation based multi-fault localization by using optics clustering. *Journal of Computer Science and Technology* 35 (2020), 979–998.
- [88] Yan Xiaobo, Liu Bin, and Wang Shihai. 2021. A Test Restoration Method based on Genetic Algorithm for effective fault localization in multiple-fault programs. *Journal of Systems and Software* 172 (2021), 110861.
- [89] J Xie, Y Zhou, M Wang, and W Jiang. 2017. New criteria for evaluating the validity of clustering. *CAAI Transactions on Intelligent Systems* 12, 6 (2017), 873–882.
- [90] Xiaoyuan Xie, Tsong Yueh Chen, Fei-Ching Kuo, and Baowen Xu. 2013. A theoretical analysis of the risk evaluation formulas for spectrum-based fault localization. *ACM Transactions on software engineering and methodology (TOSEM)* 22, 4 (2013), 1–40.
- [91] Xiaoyuan Xie, Zicong Liu, Shuo Song, Zhenyu Chen, Jifeng Xuan, and Baowen Xu. 2016. Revisit of automatic debugging via human focus-tracking analysis. In *Proceedings of the 38th International Conference on Software Engineering*. 808–819.
- [92] Xiaofeng Xu, Vidroha Debroy, W Eric Wong, and Donghui Guo. 2011. Ties within fault localization rankings: Exposing and addressing the problem. *International Journal of Software Engineering and Knowledge Engineering* 21, 06 (2011), 803–827.
- [93] Jifeng Xuan, Matias Martinez, Favio Demarco, Maxime Clement, Sebastian Lamelas Marcote, Thomas Durieux, Daniel Le Berre, and Martin Monperrus. 2016. Nopol: Automatic repair of conditional statement bugs in java programs. *IEEE Transactions on Software Engineering* 43, 1 (2016), 34–55.
- [94] Ronald R Yager and Dimitar P Filev. 1994. Approximate clustering via the mountain method. *IEEE Transactions on systems, man, and Cybernetics* 24, 8 (1994), 1279–1284.
- [95] Shin Yoo. 2012. Evolving human competitive spectra-based fault localisation techniques. In *Search Based Software Engineering: 4th International Symposium, SSBSE 2012, Riva del Garda, Italy, September 28-30, 2012. Proceedings 4*. Springer, 244–258.
- [96] Shin Yoo and Mark Harman. 2010. Using hybrid algorithm for pareto efficient multi-objective test suite minimisation. *Journal of Systems and Software* 83, 4 (2010), 689–701.
- [97] Shin Yoo, Xiaoyuan Xie, Fei-Ching Kuo, Tsong Yueh Chen, and Mark Harman. 2017. Human competitiveness of genetic programming in spectrum-based fault localisation: Theoretical and empirical analysis. *ACM Transactions on Software Engineering and Methodology (TOSEM)* 26, 1 (2017), 1–30.
- [98] Juyeon Yoon and Shin Yoo. 2021. Enhancing Lexical Representation of Test Coverage for Failure Clustering. In *2021 36th IEEE/ACM International Conference on Automated Software Engineering Workshops (ASEW)*. IEEE, 232–238.
- [99] Zhongxing Yu, Chenggang Bai, and Kai-Yuan Cai. 2015. Does the failing test execute a single or multiple faults? An approach to classifying failing tests. In *2015 IEEE/ACM 37th IEEE International Conference on Software Engineering*, Vol. 1. IEEE, 924–935.
- [100] Abubakar Zakari and Sai Peck Lee. 2019. Parallel debugging: An investigative study. *Journal of Software: Evolution and Process* 31, 11 (2019), e2178.
- [101] Abubakar Zakari, Sai Peck Lee, Rui Abreu, Babiker Hussien Ahmed, and Rasheed Abubakar Rasheed. 2020. Multiple fault localization of software programs: A systematic literature review. *Information and Software Technology* 124 (2020), 106312.
- [102] Abubakar Zakari, Sai Peck Lee, and Ibrahim Abaker Targio Hashem. 2019. A community-based fault isolation approach for effective simultaneous localization of faults. *IEEE Access* 7 (2019), 50012–50030.
- [103] Lejun Zhang, Jinlong Wang, Weizheng Wang, Zilong Jin, Yansen Su, and Huiling Chen. 2022. Smart contract vulnerability detection combined with multi-objective detection. *Computer Networks* 217 (2022), 109289.