

# AGA: An Accelerated Greedy Additional Algorithm for Test Case Prioritization

Feng Li<sup>✉</sup>, Jianyi Zhou, Yinzhu Li, Dan Hao<sup>✉</sup>, and Lu Zhang

**Abstract**—In recent years, many test case prioritization (TCP) techniques have been proposed to speed up the process of fault detection. However, little work has taken the efficiency problem of these techniques into account. In this paper, we target the Greedy Additional (GA) algorithm, which has been widely recognized to be effective but less efficient, and try to improve its efficiency while preserving effectiveness. In our Accelerated GA (AGA) algorithm, we use some extra data structures to reduce redundant data accesses in the GA algorithm and thus the time complexity is reduced from  $\mathcal{O}(m^2n)$  to  $\mathcal{O}(kmn)$  when  $n > m$ , where  $m$  is the number of test cases,  $n$  is the number of program elements, and  $k$  is the iteration number. Moreover, we observe the impact of iteration numbers on prioritization efficiency on our dataset and propose to use a specific iteration number in the AGA algorithm to further improve the efficiency. We conducted experiments on 55 open-source subjects. In particular, we implemented each TCP algorithm with two kinds of widely-used input formats, adjacency matrix and adjacency list. Since a TCP algorithm with adjacency matrix is less efficient than the algorithm with adjacency list, the result analysis is mainly conducted based on TCP algorithms with adjacency list. The results show that AGA achieves 5.95X speedup ratio over GA on average, while it achieves the same average effectiveness as GA in terms of Average Percentage of Fault Detected (APFD). Moreover, we conducted an industrial case study on 22 subjects, collected from Baidu, and find that the average speedup ratio of AGA over GA is 44.27X, which indicates the practical usage of AGA in real-world scenarios.

**Index Terms**—Test case prioritization, additional strategy, acceleration

## 1 INTRODUCTION

TEST case prioritization (abbreviated as TCP) [1], [2], [3], [4], [5], [6], is proposed to schedule the execution order of test cases so as to detect faults as early as possible. To address this problem, a large number of TCP techniques have been proposed in the literature.

Among these TCP techniques, the Greedy Additional (GA) algorithm has received much attention since it was proposed in 1999 [5] due to its widely recognized effectiveness [7], [8], [9], [10]. In particular, the GA algorithm iteratively selects the next test case which covers the largest number of elements (e.g., methods, branches, statements) that have not been covered by previously selected test cases. When the selected test cases cover all elements, this GA algorithm deals with the remaining unselected test cases with any prioritization technique (e.g., Greedy Total algorithm [5], which schedules these test cases based on the descendent order of the number of total covered program elements). Later in 2002, Elbaum *et al.* [3] slightly modified

this algorithm by reordering the remaining test cases with the GA strategy again after resetting all the elements to be “uncovered”. This GA algorithm repeats the GA strategy until all the test cases are selected and thus its effectiveness is no worse than that of the original GA algorithm [5]. Therefore, the GA algorithm proposed by Elbaum *et al.* [3] is taken as the default GA algorithm by most researchers in TCP and in this paper.<sup>1</sup> Moreover, the original GA algorithm is called the GA-first algorithm for distinction. Note that we target GA rather than GA-first in this paper because the former is more widely used in the literature. Although researchers have put dedicated efforts in TCP and have proposed a large number of TCP techniques since then, the GA approach [3] remains one of the most effective strategies in terms of fault-detection rate [7], [8], [10], which is usually measured by the average percentage of faults detected (abbreviated as APFD). In other words, none of the existing TCP techniques can always outperform GA [3] in terms of effectiveness.

Besides effectiveness, time cost is widely recognized as another important issue influencing the application of an approach [11], [12], [13], [14], especially considering the limited available time. In particular, the time cost of TCP, called TCP efficiency in this paper, refers to how much time a TCP approach consumes. As reported, Google [15] runs 800K builds and 150M tests every day (the same tests are run many times). If a TCP approach consumes much more time on prioritization, the time left for test running will be reduced to a large extent. Furthermore, software modification occurs dramatically frequently so that regression testing

- Feng Li, Jianyi Zhou, Dan Hao, and Lu Zhang are with the Institute of Software, School of Computer Science, Peking University, Beijing 100871, China, and also with the Key Laboratory of High Confidence Software Technologies, MoE, Peking University, Beijing 100871, China. E-mail: {lifeng2014, zhoujianyi, haodan, zhanglucs}@pku.edu.cn.
- Yinzhu Li is with Baidu Online Network Technology (Beijing) Company, Ltd, Beijing 100080, China. E-mail: liyinzhu@baidu.com.

Manuscript received 8 August 2020; revised 13 November 2021; accepted 20 December 2021. Date of publication 23 December 2021; date of current version 12 December 2022.

This work was supported in part by the National Natural Science Foundation of China under Grant 61872008.

(Corresponding author: Dan Hao.)

Recommended for acceptance by K. Sen.

Digital Object Identifier no. 10.1109/TSE.2021.3137929

1. Without further clarification, the GA algorithm used in this paper refers to the one proposed by Elbaum *et al.* [3].

consumes about 80% testing cost [16]. For example, Google developers modify source code one time per second on average [15]. To improve the efficiency of regression testing, it is necessary to apply TCP more than once because frequent code modification may hamper the effectiveness of TCP [17]. That is, considering the practical application of TCP, including the GA algorithm, both effectiveness and efficiency are important.

However, existing TCP approaches, including the GA algorithm, suffer from the efficiency problem, e.g., the previous work shows that most existing TCP approaches cannot deal with large-scale application scenarios [13], [15], [18]. Furthermore, some work [13], [15], [18] points out that the GA algorithm spends dramatically long time on prioritization. Note that in the 20-year history of GA, there is no approach proposed to improve its efficiency while preserving the high effectiveness.

In this paper, we make the first attempt to accelerate the GA algorithm and maintain the effectiveness. In particular, we analyze the efficiency problem of the GA algorithm and propose to accelerate the GA algorithm through two enhancements. The proposed algorithm is called the Accelerated Greedy Additional (abbreviated as AGA) algorithm. First, many redundant data accesses occur during prioritization in GA. Whenever a test case is selected, the GA algorithm scans the coverage information of all test cases to mark elements covered by this selected test case and calculates the number of unmarked elements covered by each unselected test case. Such scanning is less efficient and may contain many redundant data accesses. Therefore, we design some extra data structures (e.g., indices) to summarize the coverage information of each test case in the AGA algorithm. Supposed that  $m, n, k$  are the number of test cases, the number of elements and the number of iterations to repeat GA strategy (which is called iteration number in this paper), and given  $n > m$  (which is true in most cases), the time complexity of our AGA algorithm is  $\mathcal{O}(kmn)$ , while the time complexity of the GA algorithm is  $\mathcal{O}(m^2n)$ . The value of  $k$  determines to what extent the former is superior to the latter. In practice,  $k$  is usually much smaller than  $m$ , and in our approach,  $k$  is fixed as a constant (by the second part below), so, our  $\mathcal{O}(kmn)$  is superior to  $\mathcal{O}(m^2n)$ . Second, the GA algorithm proposed by Elbaum *et al.* [3] repeats the GA strategy multiple times in TCP and thus the iteration number is usually larger than 1. Intuitively, when an element is covered for enough times, the probability that it still contains faults is low, so the remaining iterations may not contribute to the effectiveness but only decrease TCP efficiency. Therefore, we investigated their relation empirically and applied it to modify the GA algorithm to improve efficiency but preserving effectiveness. To sum up, our AGA algorithm consists of two parts, time complexity reduction and iteration number reduction. Note that theoretical improvement is rather important and gives clear assurance for high-efficiency under any situations (especially in the first part of AGA). Also, our simple technique with theoretical improvement is meaningful in practice and can illustrate the simple nature of the problem.

We conducted controlled experiments by using 55 open-source projects from GitHub (whose total lines of code are from 1,621 to 177,546). Because the algorithm input (program coverage) has two kinds of format, adjacency matrix and

adjacency list, we conduct our experiments on both of them, which is discussed in Section 2.1. In the experiments, we studied the contributions of the two parts of AGA separately, and found that both of them improve the efficiency to a large extent. Furthermore, we investigated the effectiveness and efficiency of AGA by comparing it with GA. The results showed that on average the speedup ratio of AGA over GA is 5.95X and 27.72X on two input formats, which is a very large improvement. We also find that the average APFD of AGA and GA is the same, and Analysis of Covariance (ANCOVA) [19] shows no significant difference between them. Moreover, the effect size (Cohen's  $d$ ) also indicates small effect.

We also empirically compared AGA with *FAST* [18], which focuses on the TCP efficiency problem. As *FAST* [18] targets a different problem, improving the time efficiency by sacrificing effectiveness, such a comparison in terms of efficiency may be a bit unfair to our AGA approach. Surprisingly, the results showed that the average speedup ratio of AGA over *FAST* is 4.29X (with significant difference and medium effect), which means AGA even outperforms the technique that sacrifices effectiveness to achieve high efficiency. Also, the average APFD difference that AGA exceeds *FAST* is 0.1702, and ANCOVA shows that the difference is statistically significant. Moreover, the effect size (Cohen's  $d$ ) also indicates huge effect.

We further performed an industrial case study in Baidu, a famous Internet service provider with over 600M monthly active users. In particular, we compared the performance of AGA and GA in 22 subjects of Baidu. In this industrial case study, the average speedup ratio of AGA over GA is 44.27X and 61.43X on two input formats, which indicates the usefulness of AGA in real-world large-scale scenarios. Also, AGA is faster than *FAST* on all 22 subjects and achieves 4.58X speedup ratio on average, and the difference is statistically significant with very large effect. Due to the commercial constraints, we cannot access the source code of these projects, and the developers in Baidu also do not record the fault positions in the history, which are necessary to calculate the APFD results. So, we did not compare the effectiveness of these approaches in this study.

The contributions of this work are summarized as below.

- The first attempt to improve the efficiency of GA while preserving its effectiveness, since GA is believed to have high effectiveness. In particular, we resolve the efficiency issue of GA through theoretical improvement, which gives clear assurance for high-efficiency under any situations.
- An approach to accelerating the widely-known GA algorithm through two parts, including time complexity reduction and iteration number reduction. With the former, the complexity is reduced from  $\mathcal{O}(m^2n)$  to  $\mathcal{O}(kmn)$  given  $n > m$ , which is theoretically proved; with the latter, the corresponding AGA algorithm is more efficient and can be as competitive as GA regarding to effectiveness, which is empirically shown. In fact, although it seems like an easy-to-implement algorithm, in the broad literature, nobody realizes this optimization and the subsequent reduction of complexity. Therefore, this paper is the first to

TABLE 1  
An Example

(a) Adjacency Matrix						
Cover or Not		Elements				
		E1	E2	E3	E4	E5
Test Cases	T1	○	○	○	×	×
	T2	×	×	○	○	○
	T3	○	○	×	×	×
	T4	×	×	○	○	×
	T5	×	×	×	×	○

(b) Adjacency List						
Test Cases	Covered Elements					
T1	E1	E2	E3			
T2	E3	E4	E5			
T3	E1	E2				
T4	E3	E4				
T5	E5					

systematically analyze this problem and propose and evaluate the optimization approach, which is helpful for the community.

- Large scale experiments on 55 open-source projects demonstrating the effectiveness and efficiency of our AGA approach, compared with the GA algorithm.
- An empirical comparison of AGA with *FAST*, which improves time efficiency but decreases effectiveness.
- An industrial case study on 22 subjects from Baidu, which indicates the practical usage of AGA in real-world scenarios.

## 2 TIME COMPLEXITY REDUCTION

In this section, we review the Greedy Additional (GA) algorithm by an example (in Section 2.1). By analyzing its time complexity (in Section 2.2), we propose to accelerate GA through extra-defined data structures (in Section 2.3). Such modification improves the efficiency of GA so that the time complexity becomes  $\mathcal{O}(kmn)$  (given  $n > m$ ), whereas the complexity of GA is  $\mathcal{O}(m^2n)$ , where  $n$  is the number of program elements (e.g., statements, branches, methods) covered by the test suite,  $m$  is the number of test cases in the test suite, and  $k$  is the iteration number.

### 2.1 Example

Table 1 presents an example showing the coverage information of a test suite. This test suite consists of five test cases (i.e., T1, T2, ..., and T5) and the test suite covers five program elements (i.e., E1, E2, ..., and E5). A common representation form of coverage information is adjacency matrix, which is shown in Table 1 a. ○ represents that the test case covers the corresponding program element, while × represents the opposite. Another representation form of coverage information is adjacency list, which is shown in Table 1 b. In our example, the two forms represent totally the same information.

If we take the adjacency matrix as input, the GA algorithm runs as follow. First, no element has been covered before and this algorithm scans the whole table to calculate the number of elements covered by each test case. Then it chooses T1 or T2 since both of them cover the most elements. Supposed that this algorithm chooses T1, then T2, T3, T4, and T5 remain

unselected. As the selected test case T1 covers elements E1, E2, and E3, the rest elements E4 and E5 remain uncovered. The algorithm scans the whole table again to find that T2, T3, T4, and T5 covers 2, 0, 1, and 1 of the 2 uncovered elements, respectively. So, the GA algorithm chooses T2 as the next test case. Now, all elements have been covered and the GA algorithm [3] starts another iteration by resetting all elements to “uncovered”. Finally, the test execution sequence produced by the GA algorithm is “T1, T2, T3, T4, T5”. On the other hand, provided the adjacency list as input, GA runs similarly and produces the same output.

### 2.2 Analysis of the GA Algorithm

In this section, we analyze the time complexity of the GA algorithm through its general implementation. Suppose the coverage information is recorded in a table like Table 1 a, the GA algorithm first scans the whole table to find the line with the most “○” entries and selects the corresponding test case into the prioritized sequence. When a test case is selected and added to the sequence, the GA algorithm scans the whole table to find the “○”s whose corresponding element is covered by the latest selected test case. These “○”s are replaced by “×”s. The GA algorithm repeats the proceeding process until all the entries in the table are “×”s or all the test cases have been selected. In the latter case the termination condition is satisfied and the GA algorithm ends by producing a prioritized test suite; otherwise, GA reuses the initial table by replacing “○”s with “×”s for each selected test case and repeats the proceeding process again.

Supposed that there are  $m$  test cases in the given test suite to be prioritized and  $n$  program elements are covered by the test suite, the GA algorithm needs to scan the whole table for  $m$  times and thus the time complexity is  $\mathcal{O}(m^2n)$ , as shown by previous work [3], [7], [8]. However, lots of accesses of the table are redundant. First and the most importantly, every time the coverage table is updated, the GA algorithm recalculates the total “○” entries of each unselected test case, without reusing previous calculation. Second, none of the accesses to “×”s in the table is necessary because the GA algorithm does not want to update them in the process. Third, in order to find the elements covered by the latest selected test case, the GA algorithm scans all elements in the table, which is also unnecessary. Let us illustrate the preceding redundant accesses by the example. When T1 is selected first, the GA algorithm scans Row T1 and finds three “○”s. Among the five accesses (i.e., E1, E2, ..., E5), the accesses of E4 and E5 are redundant. Then, the GA algorithm changes the state of E1, E2, and E3 in other four test cases from “○” to “×”. During this process, it is also not necessary to access the state “×”. Then, the GA algorithm scans the whole table to select the next test case, but this process can be optimized by analyzing updated columns and the previous calculation on total number of “○” covered by each test case. To sum up, due to such a large number of redundant accesses in the GA algorithm, it is possible to reduce its time cost and improve its efficiency.

If we take the adjacency list as input, similar analysis can be done. First, the accesses of “×”s to find covered elements in one row is reduced, while more time is spent on finding all test cases that cover a specific element (through scanning of the whole list). As a result, the overall time complexity remains

$\mathcal{O}(m^2n)$ . Second, lots of accesses of the list are redundant, too. Following our previous analysis, the time efficiency can be improved through reducing the unnecessary operations.

### 2.3 Improvement of Time Complexity

To reduce such redundant accesses, we propose the **AGA\_C** approach that defines extra data structures, which reuse previous information collected during its execution. In particular, we use a list to record the total number of elements covered by each test case and dynamically update it during prioritization, in order to alleviate the scanning of the coverage table. We also use forward and inverted indices to save the data accesses of “×” entries in the table.

Our **AGA\_C** algorithm is shown in Algorithm 1. Line 1 initializes several data structures. *TC* is a list of length  $m$  recording the number of elements covered by each test case. In our example, *TC* is [3, 3, 2, 2, 1] from Table 1 by definition. *HS* is a list of length  $n$  recording whether each test case has been selected. *HC* is a list of length  $n$  recording whether each element has been covered by previous test cases. *FI* are forward indices that index all elements covered by each test case, while *II* are inverted indices that index all test cases that cover each element. From Table 1, in our example, *FI* records that T1 covers [E1, E2, E3], T2 covers [E3, E4, E5], etc. *II* records that E1 is covered by [T1, T3], E2 is covered by [T1, T3], etc. Line 2 initializes **P** as the empty list. Then, in Line 3 to Line 21, the algorithm selects  $m$  test cases in turn. First, it chooses the largest value in *TC* whose test case  $t$  is marked unselected in *HS*. The algorithm adds  $t$  to the prioritized list **P** and marks it in *HS*. In our example, in the first loop, T1 is selected (since it covers the most program elements), marked in *HS*, and added to **P**. Then, for every element  $j$  in *FI*[ $t$ ] that is marked uncovered in *HC*, the algorithm marks it as covered and for every test case  $i$  in *II*[ $j$ ], the algorithm subtracts *TC*[ $i$ ] by 1. In our example, in the first loop, E1 and E2 are marked covered and the updated *TC* is [0, 2, 0, 1, 1]. Finally, the algorithm continues to select the next test case by repeating the process. As shown from Line 5 to Line 9, if all elements have been covered by selected test cases, the algorithm completes current iteration and restores the original *TC* to start the next iteration. In our example, after T1 and T2 are selected, the original *TC* is restored. The total number of iterations is called *iteration number*.

Furthermore, we analyze the time complexity of our **AGA\_C** algorithm. All initialization operations consume  $\mathcal{O}(mn)$  time. Each calculation of maximum value in *TC* consumes  $\mathcal{O}(m)$  time, which leads to  $\mathcal{O}(m^2)$  time in total. The number of times to update *TC* is equal to the elements in *FI* (also equal to the test cases in *II*) in an iteration, which is the number of “○” entries in the coverage matrix. So, in each iteration, the algorithm updates *TC* for up to  $\mathcal{O}(mn)$  times, and the total time for updating *TC* is  $\mathcal{O}(kmn)$ , where  $k$  is the iteration number. Generally speaking, the number of elements is often larger than the number of test cases, which means  $n > m$ . So, according to the definition of Big O notation, the total time complexity  $\mathcal{O}(kmn + m^2)$  can be simplified as  $\mathcal{O}(kmn + m^2) = \mathcal{O}(kmn + mn) = \mathcal{O}((k+1)mn) = \mathcal{O}(kmn)$ , where  $k$  is the iteration number. Note that in most cases,  $n > m$  obviously holds, and can also be verified by the subject statistics in this paper (given by Table 6). For other special cases, the original time complexity  $\mathcal{O}(kmn + m^2)$  is still a large improvement.

---

### Algorithm 1. AGA\_C Algorithm

---

**Input:** Coverage information **M**;  
**Output:** Prioritized test cases **P**;  
1 Initialize *TC*, *HS*, *HC*, *FI*, and *II* from **M**,  $t = 0$ ;  
2 Set **P** as empty list;  
3 **while**  $t < m$  **do**  
4   Find the largest value in *TC* that the corresponding test case  $t$  has not been selected (take the use of *HS*);  
5   **if** No test case can be selected **then**  
6     Change *TC* to the original value;  
7     Change *HC* to the original value;  
8     Continue;  
9   **end**  
10   Add  $t$  into **P**;  
11   Mark  $t$  as selected in *HS*;  
12   **forall**  $j$  in *FI*[ $t$ ] **do**  
13     **if** *HC*[ $j$ ] is “uncovered” **then**  
14       Mark  $j$  as “covered” in *HC*;  
15       **forall**  $i$  in *II*[ $j$ ] **do**  
16         Decrease *TC*[ $i$ ] by 1;  
17       **end**  
18     **end**  
19   **end**  
20    $t = t + 1$ ;  
21 **end**  
22 Return **P**;

---

In addition, in our algorithm, we use more storage space to maintain the extra data structures in order to improve time complexity. So, it is necessary to analyze the space complexity, too. In GA, the coverage information (adjacency matrix/list) takes  $\mathcal{O}(mn)$  space, and additional  $\mathcal{O}(1)$  space is used to store temporary variables in the algorithm, which means the overall space complexity of GA is  $\mathcal{O}(mn)$ . In AGA, the same  $\mathcal{O}(mn)$  space is used to store the coverage table, while *TC*, *HS*, *HC*, *FI*, and *II* need  $\mathcal{O}(m)$ ,  $\mathcal{O}(m)$ ,  $\mathcal{O}(n)$ ,  $\mathcal{O}(mn)$ , and  $\mathcal{O}(mn)$  spaces, respectively. So, the overall space complexity of AGA is  $\mathcal{O}(mn)$ , which is the same as GA, with the only difference lying in the constant factor.

## 3 ITERATION NUMBER REDUCTION

From Section 2, we obtain a new approach with time complexity  $\mathcal{O}(kmn)$ , where  $k$  is the iteration number. In practice,  $k$  is often much smaller than  $m$  in most projects because usually many test cases are needed to cover all elements in an iteration. However, in the worst case,  $k$  may be equal to  $m$ , indicating the worst time complexity of our AGA algorithm is still the same as that of the GA algorithm.

To further improve the efficiency of the GA algorithm, especially in the worst case, in this section, we discuss the impact of the iteration number and introduce another modification adopted in our AGA algorithm. Finally, we present an experiment to evaluate the impact of iteration number on the GA algorithm.

### 3.1 Modification With Iteration Number Reduction

Let us re-examine the definition of “an iteration”. In this paper, the process of selecting some test cases from covering 0 element to covering all possible elements and then resetting them to be “uncovered” is called “an iteration”.

TABLE 2  
Statistics of the Impact of Iteration Number

Subjects	#Projects	Iteration Number						Time_GA / Time_GA-first*						APFD_range**					
		min.	Q1	Q2	Q3	max.	ave.	min.	Q1	Q2	Q3	max.	ave.	min.	Q1	Q2	Q3	max.	ave.
Open-Source	55	1	6	10	16.5	679	29.20	1.00	1.21	1.57	1.84	17.56	2.14	0.0000	0.0004	0.0013	0.0039	0.1328	0.0085

\* The time ratio between the GA algorithm and the GA-first algorithm.

\*\* The highest APFD subtracts the lowest APFD among all iteration numbers.

Intuitively, the iteration number may have large impact on time cost of the GA algorithm. The difference between the GA-first algorithm [5] and the GA algorithm [3] also indicates the influence of such an iteration number. Moreover, between these two algorithms, there exist many other potential algorithms, depending how many times the GA strategy is used (i.e., iteration number of the GA strategy) and what strategy is used to deal with the remaining unselected test cases (e.g., Greedy Total strategy, which schedules test cases based on the descendent order of the number of total covered program elements)

$$n = k \times l. \quad (1)$$

Here, we define the average number of test cases selected in one iteration as  $l$ , so we can deduce Formula (1). According to Formula (1), if a large number of test cases are selected in one iteration, the total iteration number of this project is small; if few test cases are selected in one iteration, the total iteration number of this project is large. As our goal is to improve efficiency while preserving effectiveness, projects with small iteration number have already been efficient enough, and the time complexity  $\mathcal{O}(kmn)$  can be reduced to  $\mathcal{O}(mn)$ . For those projects with large iteration number, it is necessary to optimize the iteration number to some extent.

In fact, everytime a program element is covered, the probability that it still contains faults decreases. After many iterations, all elements have been covered for enough times. On one hand, if all faults have been revealed after these iterations, the remaining iterations are useless for detecting faults but only increase the time cost. On the other hand, if there are still several faults existed after many iterations, they are supposed to be hard to reveal and the remaining iterations may only reveal them by chance, intuitively. So, we conjecture that after some iterations, the effectiveness of GA just fluctuates along with the remaining iterations.

Based on the above reasoning, we introduce another component of the proposed AGA algorithm, **AGA\_I**. AGA\_I reduces the time cost by reducing the iteration number. Different from the GA algorithm, AGA\_I does not repeat applying the GA strategy until all the test cases are prioritized, but stops when the specified iteration number is achieved. Regarding to the remaining unselected test cases, AGA\_I applies other less costly techniques (e.g., the Greedy Total technique (GT) [5], which is usually used in previous work and also in this paper). Take Table 1 as an example, the original iteration number is 2. If we reduce it to be 1, AGA\_I does not repeat the additional strategy after selecting T1 and T2 and prioritizes the remaining test cases using GT.

### 3.2 Experiment

We conjecture that AGA\_I does not influence the effectiveness (e.g., APFD) much but can improve efficiency (i.e.,

time cost) a lot. To verify our conjecture, we design an experiment to investigate how the iteration number impacts TCP in terms of both effectiveness and efficiency.

Specifically, we use the same setup as the comprehensive experiments in Section 5. More details about the subjects, faults, implementation and supporting tools, and measurement are given in Section 5.

We applied the GA algorithm to all subjects, and recorded the total number of iterations the GA strategy is applied during the process for each project, which is denoted as  $k$ . Then we applied to each project  $k$  modified GA algorithms, each of which is denoted as algorithm  $algo_i$  ( $1 \leq i \leq k$ ), recording their APFD values and time spent during prioritization. In particular, algorithm  $algo_i$  repeats the GA strategy  $i$  times and prioritizes the remaining unselected test cases by the Greedy Total strategy [5]. Note that algorithm  $algo_1$  is actually the GA-first algorithm, whereas algorithm  $algo_k$  is actually the GA algorithm.

Due to space limit, we only present some statistics of the experimental results in Table 2, that is minimum, maximum, average, quartiles (Q1, Q2, Q3), and the detailed results are given on the website of this project. From the eighth column, the average iteration number among all open-source subjects is 29.20. The ninth to the fourteenth columns present the ratio between the time cost of the GA approach and that of the GA-first approach [5]. The big gap between the maximal and minimal time ratio indicates the influence of the iteration number. To better analyze the relationship between iteration number and time cost, we put detailed results in Appendix A, which can be found on the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TSE.2021.3137929>. We draw a line chart of iteration number and time cost for each project. Note that in order to see the trend, we only present the projects whose iteration number is no less than 20 ( $k \geq 20$ ). The plots also support our claim that the iteration number contributes much to the time cost. As  $k$  is the coefficient of time complexity, it largely determines the actual efficiency in practice, so, we think there is a large space to reduce time complexity.

The last six columns in Table 2 present the APFD ranges of each project with different iteration numbers, that is, the highest APFD value minus the lowest APFD value. From the quartiles, we conclude that although some outliers exist, most of the APFD ranges are very small. And the average APFD range is only 0.0085 among all open-source subjects, indicating that little fluctuation of APFD occurs as the iteration number varies.

To sum up, we have two main observations. First, along with the increase of the iteration number, the time cost also increases, indicating that the iteration number contributes much to the time cost. Second, the APFD value varies a little when the iteration number varies, which means a too large

iteration number contributes little to the APFD value. These two observations also verify our conjectures in Section 3.

As we discuss in Section 3, projects with small iteration numbers are efficient enough by using AGA\_C, so, we need to decide a proper reduced iteration number for projects with a large iteration number. In fact, this reduced iteration number is not fixed, which means it can be adjusted for specific usage. In this paper, we determine this value from some heuristics. On one hand, although we conjecture that there is no need to conduct too many iterations to detect faults, we still prefer to choose a relatively high value to ensure the effectiveness. On the other hand, if we assume that every time an element is covered, the probability that it still contains faults decreases to half of the original probability, given that the initial probability is 1, we need to cover an element 10 times to reduced the probability to be less than 1% ( $(1/2)^{10} = 1/1024$ ). As a result, in the remaining of this paper, we implement our AGA approach by using 10 as the reduced iteration number.

*Finding:* The iteration number has large influence on the efficiency of the GA algorithm, while it impacts little on effectiveness. In this paper we set the iteration number to be 10 in implementing the AGA approach.

Note that this finding is confirmed on our dataset empirically and may have bias considering the diversity of different datasets. However, the constraint on  $k$  does reduce the overall time complexity from  $\mathcal{O}(kmn + m^2)$  to  $\mathcal{O}(mn + m^2)$ . When  $n > m$ , which is general in most cases, the reduction is from  $\mathcal{O}(kmn)$  to  $\mathcal{O}(mn)$ .

### 3.3 Discussion on the Chosen Iteration Number

In this paper, we set the iteration number to be 10 in implementing AGA through some heuristics. Here we discuss the influence of this choice. First, we analyzed the APFD results of the GA algorithm with various iteration number (i.e., algorithm  $algo_i$  ( $1 \leq i \leq k$ ) in Section 6.1). In particular, for each project we recorded the highest APFD value (denoted as  $APFD_{\max}$ ) among these algorithms, and found the smallest iteration number  $r$  whose corresponding APFD value is no smaller than  $APFD_{\max} * 99\%$ . Surprisingly, the smallest iteration number  $r$  for all projects are no larger than 10, which indicates that only several iterations is enough for maintaining original effectiveness, even in projects with the iteration number up to 679. Second, although we set the iteration number to be 10 in this paper, it may not be the best choice. We respectively applied  $algo_8$ ,  $algo_9$ ,  $algo_{10}$ ,  $algo_{11}$ , and  $algo_{12}$  to all projects with  $k > 8$  as Section 6.1, and found that the gap between the maximum and minimum APFD value of these  $algos$  is 0.0006 on average, which means that there might be many possible choices of the reduced iteration number in practice. In other words, the value of  $k$  in our evaluation is decided by reasoning, but it can have various values, depending on the choices of developers. For example, they can use historical faults or seeded faults to empirically decide the value of  $k$ .

## 4 RESEARCH METHOD

To investigate the performance of our proposed AGA approach, we design comprehensive experiments. In this section, we briefly introduce each component of our experiments and their intentions.

- 1) The main experiment of this paper is designed to confirm the contributions of our approach, and thus we investigate the improvement of AGA and its component (i.e., AGA\_I and AGA\_C) over the GA algorithm. In particular, this experiment is conducted on 55 open-source subjects. Details of this experiment are referred to Sections 5 and 6. Note that the experiment in Section 3.2 also shares the same setup and RQ1 complements the experiment in Section 3.2. This part of experiment can show the superiority of AGA in widely-used open-source subjects.
- 2) Although we aim to improve the efficiency of GA, we are also curious about how AGA performs compared with other TCP techniques. Specifically, *FAST* targets the TCP efficiency problem and its goal is close to ours. Therefore, we first compare AGA with *FAST*, and then with other representative TCP techniques, including ART-D, GA-S, and GE. This experiment is in Section 7. This part of experiment can show that AGA even outperforms techniques that aim to reduce TCP time cost while sacrificing effectiveness.
- 3) To show the practical usage of our approach, we conduct an industrial case study on Baidu, which is a famous Internet service providers with over 600M monthly active users. Specifically, we compare AGA with GA, *FAST*, ART-D, GA-S, and GE, respectively and the experiment is in Section 8. This part of experiment can show that AGA works also well in real-world industrial applications and we receive positive feedback from Baidu.

## 5 EVALUATION DESIGN

We conducted experiments to evaluate our AGA approach. The experiments was performed on a server whose CPU is Intel(R) Xeon(R) E5-2683 2.10GHz with 132GB memory and whose operating system is Ubuntu 16.04.5 LTS. To make a fair comparison of time cost, we conducted all experiments on a single thread without parallel execution.

In order to make our results more reliable and let readers reuse the artefacts, we share our data, analysis scripts, and detailed data tables online. They are publicly available on our website: <https://github.com/Spiridempt/AGA>, and also on figshare: <https://figshare.com/s/cf8cc6ba9259c0e0754d>.

### 5.1 Research Questions

As our AGA approach consists of two parts, time complexity reduction (AGA\_C) and iteration number reduction (AGA\_I), the first two research questions are to investigate their impacts, separately. Note that the first research question also complements the experiment in Section 3.2. The third research question is designed to investigate the performance of the whole AGA approach by comparing it with the GA algorithm. To investigate the influence of coverage type, the fourth research question is designed to investigate



whether AGA can also improve the efficiency of GA with method coverage.

To sum up, this experiment is to answer the following three research questions.

**RQ1:** How does our reduction of iteration number perform compared with the GA algorithm in terms of efficiency?

**RQ2:** How does our reduction of the time complexity perform compared with the GA algorithm in terms of efficiency?

**RQ3:** How does our AGA approach perform compared with the GA algorithm in terms of effectiveness and efficiency?

**RQ4:** Can our AGA approach also improve the efficiency when method coverage is used?

## 5.2 Subjects and Faults

**Subjects.** In this work, we use 55 open-source projects in total. Among these projects, 33 are widely used in prior work [17], [20], [21], the others are the most popular subjects selected from GitHub according to the number of stars. Specifically, we target Github subjects whose primary programming language is Java and order them according to the number of stars in Jan 2019. Then, we check the first 100 subjects and keep only the ones that are code repository and the required tools (e.g., Maven, Clover, PIT, which is explained in Section 5.3) could work. All the open-source projects used in this work are written in *Java*, whose number of lines of code is from 1,621 to 254,284. Each of these projects has a test suite written in JUnit Testing Framework. The detailed information is given in Appendix B (Table 6), available in the online supplemental material. It is worth noting that compared with the experimental dataset used in recent TCP work [18], [22], [23], our dataset is larger and contains more large-scale projects, which can make our experimental results more reliable and convincing.

**Faults.** As existing work [24], [25], [26] have demonstrated mutation faults to be suitable for software testing experimentation and mutation faults are widely used in prior work [7], [17], [22], [27], [28], [29], [30], [31] to evaluate test case prioritization, we use a widely-used mutation testing tool PIT [32] to generate mutants for all open-source subjects. In particular, for each subject, first, we generate all mutants. Second, we keep the mutants that are killed by at least one failing test case.<sup>2</sup> Third, we construct one mutation group for each subject by containing all the remaining mutation faults, which is also consistent with previous work [12].

## 5.3 Implementation and Supporting Tools

We used Clover [33] to collect code coverage information including both statement coverage and method coverage for each open-source subject. In this work, most experiments are conducted on statement coverage because it is the mostly studied test case prioritization granularity and its low-efficiency problem is severe. In other word, the number of statements is larger than the number of methods and branches. Additionally, we also design a research question

to investigate whether AGA still improves efficiency in the scenario of method coverage. The implementation code and all scripts used in this work are written with *Python*.

In prior work on coverage based test case prioritization, some takes the adjacency matrix as input [12], while some uses the adjacency list [18]. In this work, in order to make a more general comparison, on one hand, we utilize the GA implementation in [18], which is a relatively efficient implementation and uses adjacency list as input, and we implement AGA based on adjacency list. On the other hand, we implement GA and AGA based on adjacency matrix, too. Due to the space limit, in the experimental results, we only report the results based on adjacency list [18], which can be more reliable, and the detailed results based on adjacency matrix are put on the website.

It is worth mentioning that in our experiments, when ties happen (i.e., more than one test case has the same number of covered elements), AGA/GA selects the topmost test case in the test-list (given by developers).

## 5.4 Compared Prioritization Approaches

Besides the proposed AGA approach and the GA approach [3], in this study we also implemented the GA-first approach proposed by Rothermel *et al.* [5]. The GA-first approach [5] applies the greedy additional strategy only in the first iteration, and deals with the remaining test cases by other prioritization approach, e.g., the Greedy Total approach in this paper, which schedules these test cases based on the descendent order of the number of covered program elements.

## 5.5 Measurement

In this study, similar to existing work [3], [5], we used the Average Percentage of Fault Detected (APFD) to measure the effectiveness of TCP approaches. Formula (2) presents how to calculate APFD values for a subject with  $n$  tests and  $m$  faults. Typically,  $TF_i$  represents the first test case's position in the test suite that detects the  $i$ th fault

$$APFD = 1 - \frac{TF_1 + TF_2 + \dots + TF_m}{nm} + \frac{1}{2n}. \quad (2)$$

Besides, we used the total time spent during the TCP process to measure the efficiency of a TCP approach. For fair comparison, we included the preparation time for a TCP approach, i.e., the time spent in constructing extra data structures in the AGA approach.

## 5.6 Threats

The internal threats to validity mainly lie in the implementation of studied approaches and scripts used in the experiments. To reduce this threat, the first two authors reviewed all the implementation and scripts used in this work. Also, to improve the reliability of our work, we reuse some implementation code in previous work [18] to reduce the threats.

The external threats to validity mainly lie in the subjects and faults. To reduce the former threat, we used 55 widely used open-source subjects in our study, which consist of 33 previously used subjects [17], [20] and 22 popular subjects selected from GitHub. At the same time, as AGA is a general approach, it is not biased towards the chosen projects.

2. That is, the subject and the mutant produce different outputs on at least one test case.

Note that because the second part of our approach (iteration number reduction) is empirically verified on our dataset, the large dataset itself also addresses the threat that our approach may be biased. Also, some prior work [34], [35] shows that the relative performance of different test case prioritization techniques on mutation faults may not strongly correlate with the performance on real faults, depending upon the attributes of the studied subjects, but we follow the common practice to use mutation faults for open-source projects following the preceding TCP work [7], [17], [22], [27], [28], [29], [30], [31]. Additionally, to complement this experiment, in Section 7, we also evaluate our approach on real faults. In the future, we plan to conduct an extensive study by using more projects with more real faults. In addition, in this paper, we only target the GA algorithm and compare AGA with it. On one hand, it is widely accepted that GA remains one of the most effective strategies in terms of fault-detection rate [7], [8], [10]. On the other hand, the results of a recent work [12] shows that other black-box techniques that do not use coverage information (e.g., [36], [37]) are often less effective than GA. At the same time, we also design another experiment in Section 7 to compare AGA with other representative prioritization techniques. Additionally, most of our experiments are conducted on statement coverage because its wide usage and severe low-efficiency problem. In fact, our analysis of AGA is regardless of the scale of coverage matrix, and our theoretical improvement is general for all types of coverage. We also include RQ4 to empirically verify our improvement on method coverage. Another minor threat is induced by the diversity of used subjects, which may lead to misleading statistics of our results. To address this threat, besides reporting the mean and median values, we also draw violin plots to learn the data distribution, which are shown on our website.

## 6 RESULTS AND ANALYSIS

In this section, we analyze the experimental results on open-source projects and answer the four research questions.

### 6.1 RQ1: Efficiency of Iteration Number Reduction

In this section, we further investigate the efficiency improvement of the iteration number reduction. According to Section 3.2, we implement our approach with iteration number reduction alone by setting  $k = 10$  and call this implementation AGA\_I. In other words, in this subsection, we assess the contribution of iteration number reduction alone (without the time complexity reduction).

The results on the 55 open-source projects are given in Table 7 (Appendix C)<sup>3</sup>, available in the online supplemental material, where the projects are sorted in ascending order of source lines of code (SLOC) and the first two columns present the results for RQ1.  $\text{Time}_{\text{GA}}$  presents the time cost of the GA approach, whereas  $\text{Time}_I$  represents that of AGA\_I. The speedup ratio of AGA\_I over GA is 1.08X. It is apparent that most subjects have a small iteration number in GA (less than or slightly more than 10). Therefore, AGA\_I does not improve the efficiency much for them. However, for those

subjects with a large iteration number, AGA\_I could reduce their time cost.

To statistically check the differences between AGA\_I and GA, we adopt hypothesis test. We first use Shapiro-Wilk test [38] to check the normality of residuals, and the  $p$ -value in AGA\_I and GA is  $9.416 \times 10^{-16}$  and  $5.239 \times 10^{-16}$ , which reject the hypothesis that they are normally distributed. Therefore, we need to adopt a non-parametric test. As we need to include project size as a control variable, Wilcoxon rank sum test [39] cannot be used. We seek for the proportional odds regression [40], which is a class of generalized linear models and is equivalent to Wilcoxon rank sum test when there is a single binary covariate. We introduce a variable “group” representing AGA\_I and GA and take project size as a control variable. The results show that the  $p$ -value of “group” is  $1.380 \times 10^{-6}$ , indicating significant difference between AGA\_I and GA, and the effect size (Cohen’s  $d$  [41]) is 0.274 (medium effect). Here, because statistical tests of normality (e.g., Shapiro-Wilk test) might be impacted by characteristics of the data, we draw the normal probability plots additionally and put them on our website. Note that this applies to all normality checks in the following of the paper.

### 6.2 RQ2: Efficiency of Time Complexity Reduction

The extra data structures defined in our AGA\_C approach do not affect the prioritization results, but reduce the time complexity of prioritization. In this section, we compared AGA\_C with GA only in terms of time cost. Note that we did not implement AGA\_I in this research question.

The results are given by the first five columns (except the third column) of Table 7 (Appendix C), available in the online supplemental material, where  $\text{Time}_{\text{GA}}$  presents the time cost of the GA approach and  $\text{Time}_C$  represents that of AGA\_C. Moreover, we mark the results of  $\text{Time}_C$  with ✓ only if  $\text{Time}_C < \text{Time}_{\text{GA}}$ .

The last row summarizes the total number of subjects where AGA\_C outperforms the GA approach. The results show that in most projects (48 out of 55 open-source subjects), the time cost of AGA\_C is lower than the GA approach [3], which confirms our previous theoretical analysis in Section 2. As we can see, in smaller subjects, the differences between GA and AGA\_C are very small, which may be caused by precision errors resulting from calculation or the operating system. In larger subjects, their differences are very large, which indicates the efficiency of AGA\_C.

In order to make our experiments comprehensive, we compared AGA\_C with the GA-first approach, whose time cost is given by the fourth column  $\text{Time}_{\text{GAF}}$  of Table 7 (Appendix C), available in the online supplemental material. In 36 open-source subjects, AGA\_C is even more efficient than the GA-first approach, which applies the time-consuming additional strategy for only one iteration.

In general, the efficiency improvement of AGA\_C is usually very large. In particular, if we define  $\text{Time}_{\text{GA}}/\text{Time}_C$  as the speedup ratio of AGA\_C over GA for a project, the average speedup ratio is 4.37X. As small time cost may yield biased speedup ratio, also in order to show the performance of AGA in projects with different sizes, we classify all 55 projects into small-size, middle-size, and large-size, according to the SLOC. The small-size projects ( $\mathcal{S}_1$  to  $\mathcal{S}_{22}$  in Table 7 (Appendix C), available in the online supplemental material, all have less

3. Due to the space limit, we put the results of several research questions into one table and put the table in Appendix C, available in the online supplemental material.



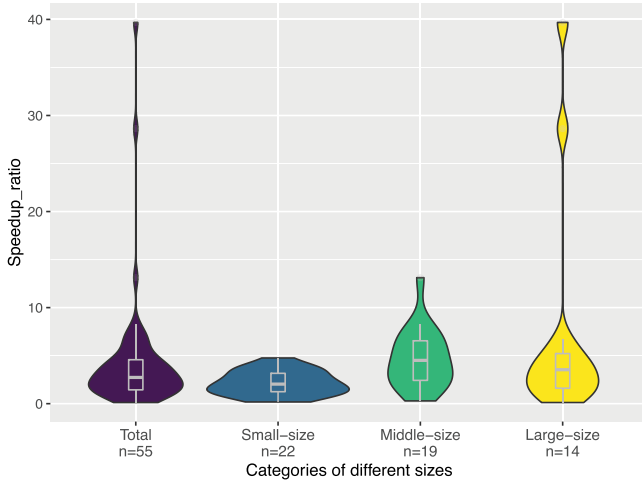


Fig. 1. Speedup ratios distribution of AGA\_C over GA on open-source projects.

than 5,000 SLOC, the middle-size projects ( $S_{23}$  to  $S_{41}$ ) all have 5,000-20,000 SLOC, and the other large-size projects have more than 20,000 SLOC. The results show that the average speedup ratio in the three categories is 2.16X, 4.65X, and 7.44X, respectively. So, the reduction of time complexity (AGA\_C) performs well, especially in projects with large sizes. In order to give a more deep view into the distribution and variation of speedup ratios, we further present the violin plot with included box plot in Fig. 1. The X-axis represents all projects and projects in three categories, respectively. We put the violin plots and box plots together to better present the distributions. From the plots, the speedup ratio of large-size projects tends to be slightly larger than that of small-size projects. Moreover, from the plot of large-size projects, several projects have very large speedup ratio because their scale is also large.

To statistically check the differences between AGA\_C and GA, we perform hypothesis testing similar to the above. We first use Shapiro-Wilk test [38] to check the normality of residuals, and the  $p$ -value in AGA\_C and GA is  $4.207 \times 10^{-15}$  and  $5.239 \times 10^{-16}$ , which reject the hypothesis that they are normally distributed. We also use the proportional odds regression [40] and include project size as a control variable. The results show that the  $p$ -value of “group” is 0.038, indicating significant difference between AGA\_C and GA, and the effect size (Cohen’s  $d$  [41]) is 0.234 (medium effect).

Besides, we also calculate the speedup ratios of AGA\_C over GA-first for a more complete comparison. The average speedup ratio is 3.01X, and the average speedup ratio in the three categories is 1.26X, 3.31X, and 5.35X, respectively. This shows our AGA approach is also superior to GA-first.

To statistically check the differences between AGA\_C and GA-first, we perform the similar procedure as above. We first use Shapiro-Wilk test to check the normality of residuals, and the  $p$ -value in AGA\_C and GA-first is  $4.207 \times 10^{-15}$  and  $3.828 \times 10^{-16}$ , which reject the hypothesis that they are normally distributed. We also use the proportional odds regression [40] and include project size as a control variable. The results show that the  $p$ -value of “group” is 0.399, indicating no significant difference between AGA\_C and GA-first, and the effect size (Cohen’s  $d$ ) is 0.208 (medium effect).

Provided adjacency matrix as input, we also implemented GA and AGA\_C, and the detailed results are on our website. Specifically, the average speedup ratio of AGA\_C over GA is 24.18X, and the average speedup ratio in the three categories is 5.47X, 28.16X, and 48.19X, respectively.

Besides, the speedup ratios of the AGA\_C approach vary a lot in different projects. On 19 open-source subjects AGA\_C is less efficient than GA-first. On the one hand, the iteration numbers of these projects are high so that AGA\_C becomes a bit costly. On the other hand, in the only iteration of GA-first, few test cases are needed to cover all statements and they are selected fast so that GA-first is efficient on these projects.

To sum up, AGA\_C addresses the high-complexity problem of GA well and successfully reduces its time complexity. For any project, any scale of coverage matrix, our approach could improve the efficiency a lot.

*Conclusion to RQ2:* The time complexity reduction strategy used in our AGA approach demonstrates great efficiency improvement compared to GA. Specifically, the average speedup ratio of AGA\_C over GA is 4.37X/24.18X on two types of input.

### 6.3 RQ3: Comparison With Greedy Additional Approaches

In this section, we compare the effectiveness and efficiency between the proposed AGA approach and two Greedy Additional approaches (including both GA and GA-first), whose results are given by the first ninth columns (except the third and fifth column) of Table 7 (Appendix C), available in the online supplemental material, where  $APFD_{AGA}$  and  $Time_{AGA}$  represent the APFD results and time cost of the AGA approach whose iteration number is set to be 10. Moreover, when the GA approach [3] does not outperform the corresponding AGA approach [3], i.e.,  $APFD_{AGA} \geq APFD_{GA}$  or  $Time_{AGA} < Time_{GA}$ , the corresponding results of the AGA approach is marked with ✓.

#### 6.3.1 Effectiveness

The proposed AGA approach has the same or better APFD performance as the GA approach in 51 out of 55 open-source subjects, and the average APFD value of AGA is 0.8870, which is the same as GA. On some subjects (e.g., the open-source project whose ID is  $S_{44}$ ), the AGA approach does not outperform the GA approach, but their APFD difference is usually very small (e.g., 0.0021 for this subject). We also make extra comparisons of AGA and GA-first and find that AGA has the same or better APFD performance as GA-first in 45 out of 55 open-source subjects and their average APFD values are the same. On 14 projects, neither the AGA approach nor the GA approach outperforms the GA-first approach, but their differences are small. Through our analysis, we suspect that after the first iteration, although all elements have been covered, the numbers of times that each element is covered still differ. This means test cases with a small number of times being covered should have higher priority, but in later iterations, this information is ignored.

Moreover, we statistically analyze whether the AGA approach and the Greedy Additional approaches have

significant difference on their APFD values. First, we conduct the Shapiro-Wilk test to check the normality of residuals. The  $p$ -value of AGA, GA, and GAF is 0.328, 0.298, and 0.283, indicating we cannot reject the hypothesis that they are normally distributed. We additionally perform Shapiro-Wilk test to check the normality of residuals, and the  $p$ -value of AGA, GA, and GAF is 0.328, 0.298, and 0.283, indicating we cannot reject the hypothesis that they are normally distributed. Therefore, we can use parametric test in the following. We use Bartlett's test [42] to check the homogeneity of variance, and the  $p$ -value is 0.880, indicating we cannot reject the hypothesis that they have equal variance. Then, as we need to take project size as a control variable (covariate), we use Analysis of Covariance (ANCOVA) [19], a parametric test that works on two or more groups to check whether different groups have the same means. The  $p$ -value is 0.641, indicating we cannot reject that they have the same means. Then, pairwise ANCOVA tests show that the  $p$ -values of AGA versus GA, AGA versus GAF, and GA versus GAF are 0.981, 0.427, and 0.414. In other words, the probability that AGA is as competitive as GA is more than 98%. Then, we employ Cohen's  $d$  [41] to compute the effect size (ES), and the results in AGA versus GA, AGA versus GAF, and GA versus GAF are 0.005, 0.151, and 0.156, which are all small effects. Furthermore, we conduct Tukey's range test [43] to check the 95% confidence intervals for all pairwise differences, and the results are [-0.022, 0.022], [-0.030, 0.015], and [-0.030, 0.014].

### 6.3.2 Efficiency

According to Table 7 (Appendix C), available in the online supplemental material, in almost all subjects (i.e., 44 out of 55), the time cost of AGA is much lower than the GA approach. On average, the speedup ratio of AGA over GA is 5.95X. Moreover, the speedup ratios in small-size, middle-size, large-size projects are 2.26X, 6.69X, and 10.76X, respectively. To learn the distribution of speedup ratios in small-size, middle-size, large-size projects, we also present the violin plot with included box plot in Fig. 2. From this figure, most medium-size and large-size projects achieve higher speedup ratios than small-size projects. Moreover, AGA achieves very large speedup ratios on some large-size projects. So, AGA scales up well in large-size projects. Furthermore, we compared the time cost of the AGA approach with the GA-first approach, which requires less time than the GA approach, and find that the AGA approach even outperforms the GA-first approach in 37 open-source subjects. The average speedup ratio is 3.95X, and the average speedup ratio in the three categories is 1.36X, 4.39X, and 7.44X, respectively. Here, we notice that the speedup ratio of AGA over other approaches is sometimes less than 1 (e.g.,  $S_3$ ,  $S_4$ ,  $S_7$ ). In fact, the overall time complexity analysis is meaningful only when the parameters are large enough. In our dataset, some projects have a relatively small  $m$  value. In this case, although  $\mathcal{O}(mn)$  seems to be small, its coefficient is not negligible compared to  $m$ . In other words, the preliminary data structure setup consumes much time and it impacts the overall running time in some cases. This is also consistent with the empirical results that AGA performs better on large projects. On the other hand, the adjacency lists in some projects

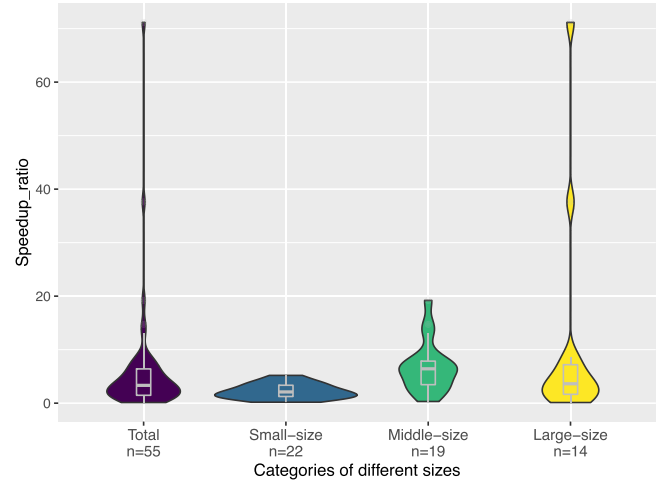


Fig. 2. Speedup ratios distribution of AGA over GA on open-source projects.

are very dense, which takes much time in the preparation of data structure, and further leads to a large coefficient. For example,  $S_7$  and  $S_{42}$  have relatively small  $m$  values (45 and 34) and dense adjacency lists.

Provided adjacency matrix as input, we also implemented GA and AGA. The average speedup ratio of AGA over GA is 27.72X, and the average speedup ratio in the three categories is 5.84X, 35.47X, and 51.59X, respectively.

To sum up, not surprisingly, the speedup ratio of AGA is higher than AGA\_C and AGA\_I. After combining AGA\_I and AGA\_C, our whole AGA approach obtains more efficient results while preserving high effectiveness. At the same time, the proposed AGA approach is demonstrated to be efficient especially on large-scale projects. In fact, the surprisingly high efficiency of the AGA approach also indicates the existence of many redundant accesses of data and it is ubiquitous in most projects.

*Conclusion to RQ3:* The AGA approach requires much less time in prioritization than the GA approach and the average speedup ratio is 5.95X and 27.72X on two types of input. Also, AGA is as competitive as the latter in terms of APFD values (with no significant difference). This means that we achieve our goal in this paper and it has promising use in practice.

### 6.4 RQ4: Performance on Method Coverage

In previous research questions, we focus on statement-level coverage because it is the mostly studied coverage criterion and its low-efficiency problem is more severe than other granularities. In this section, we collect the method-level coverage for each of our 55 subjects and compare the efficiency of AGA and GA. The results are shown in Table 3. For each subject, we report the running time (in seconds) of GA and AGA.

According to Table 3, in almost all subjects, the time cost of AGA is much lower than GA. On average, the speedup ratio of AGA over GA is 6.02X. Moreover, the speedup ratios in small-size, middle-size, large-size projects are 2.28X, 7.32X, and 10.13X, respectively. Compared to the results on statement coverage in Section 6.3, the speedup ratios are almost the same

TABLE 3  
Results of Open-Source Subjects (Method-Level)

Project	Time <sub>GA</sub>	Time <sub>AGA</sub>	Project	Time <sub>GA</sub>	Time <sub>AGA</sub>	Project	Time <sub>GA</sub>	Time <sub>AGA</sub>
$S_1$	0.0030	0.0031	$S_2$	0.0011	0.0010	$S_3$	0.0014	0.0013
$S_4$	0.0021	0.0116	$S_5$	0.0051	0.0014	$S_6$	0.0019	0.0007
$S_7$	0.0007	0.0012	$S_8$	0.0032	0.0028	$S_9$	0.2587	0.0504
$S_{10}$	0.0145	0.0081	$S_{11}$	0.0158	0.0080	$S_{12}$	0.0027	0.0020
$S_{13}$	0.0068	0.0048	$S_{14}$	0.0158	0.0076	$S_{15}$	0.0036	0.0008
$S_{16}$	0.0106	0.0024	$S_{17}$	0.0741	0.0146	$S_{18}$	0.0047	0.0013
$S_{19}$	0.2907	0.1007	$S_{20}$	0.0024	0.0023	$S_{21}$	0.0098	0.0064
$S_{22}$	0.0111	0.0055	$S_{23}$	0.0502	0.0147	$S_{24}$	0.0035	0.0031
$S_{25}$	1.1046	0.1654	$S_{26}$	0.0131	0.0043	$S_{27}$	0.0041	0.0020
$S_{28}$	0.1945	0.0347	$S_{29}$	1.2958	0.3499	$S_{30}$	0.0624	0.0177
$S_{31}$	0.4390	0.0495	$S_{32}$	0.3495	0.0589	$S_{33}$	0.1714	0.0162
$S_{34}$	0.8556	0.0824	$S_{35}$	0.2642	0.0330	$S_{36}$	31.6469	3.8975
$S_{37}$	0.7977	0.0594	$S_{38}$	0.5507	0.0445	$S_{39}$	6.6268	0.3601
$S_{40}$	0.6570	0.0963	$S_{41}$	0.3203	0.0463	$S_{42}$	0.0011	0.0006
$S_{43}$	0.2057	0.0211	$S_{44}$	1.3687	0.1088	$S_{45}$	0.0202	0.0034
$S_{46}$	0.0182	0.0044	$S_{47}$	0.0183	0.0045	$S_{48}$	0.0010	0.0003
$S_{49}$	0.0230	0.0183	$S_{50}$	0.0011	0.0007	$S_{51}$	1.0812	0.0883
$S_{52}$	0.2423	0.0652	$S_{53}$	0.1631	0.0369	$S_{54}$	15.7745	1.1856
$S_{55}$	190.9669	2.9963						

for all projects and projects in different sizes. This confirms that AGA also works well on method coverage.

In fact, the complexity analysis of GA and our AGA approach is based on a general (0,1) matrix, regardless of the meaning behind it. In other words, the type of program element (e.g., statement, method) does not affect any aspect of AGA, which means our approach works on any coverage and has a stable improvement.

*Conclusion to RQ4:* The AGA approach also works on method-level coverage. Specifically, the average speedup ratio of AGA over GA is 6.02X.

## 7 EMPIRICAL COMPARISON WITH REPRESENTATIVE PRIORITIZATION TECHNIQUES

In this section, we present an experiment comparing AGA with some representative prioritization techniques. In particular, as *FAST* targets the TCP efficiency problem and thus is closet to our goal, we first present the comparison study with *FAST* in Section 7.1. Then we present the comparison study with other representative TCP techniques in Section 7.2.

### 7.1 Comparison With *FAST*

In this section, we investigate the performance of AGA with its most related work *FAST* [18]. In particular, *FAST* is proposed as a TCP approach to address the general TCP efficiency problem by sacrificing the TCP effectiveness, and it is shown to be more efficient than other TCP techniques [18]. Note that there is no other work in the literature focusing on the same objective as ours, and thus we compare AGA against *FAST*. However, AGA and *FAST* target at slightly different goals: *FAST* approach focuses on the efficiency problem of test prioritization, not specific to GA approaches. Although *FAST* targets a different goal, it is still interesting to learn how AGA performs compared

with *FAST* in terms of time cost since both AGA and *FAST* can be viewed as addressing the efficiency problem. However, as *FAST* improves efficiency while sacrifices effectiveness, the comparison in terms of time cost is a bit “unfair” for AGA.

In this study, we compare the performance of AGA and *FAST* on both the 55 open-source projects used in Section 5 and Defects4J [44], which is the largest real-fault benchmark (i.e., a set of projects with reproducible real bugs) widely used in test case prioritization [35], [45], [46], [47], [48] and fault localization [49], [50], [51], [52], [53]. For ease of understanding, we present the results of the former subjects with seeded faults and the results of the latter subjects with real faults separately.

The *FAST* approach borrows algorithms commonly used in the big data domain to find similar items and contains a family of similarity-based test case prioritization approaches. In general, the authors proposed two categories of *FAST*, While-box (WB) and Black-box (BB). BB approaches take test code as input, while WB approaches take program coverage as input. As WB approaches have the same input as us and are much faster than BB approaches, we compare our work with WB approaches [18]. WB approaches include five algorithms *FAST-pw*, *FAST-all*, *FAST-1*, *FAST-log*, and *FAST-sqrt*, whose difference lies in how many test cases are randomly selected for prioritization at a time. In this section, we implemented this family, and for each subject, we compared the best results of this family with AGA. Specifically, according to prior work [18], none of the algorithms in *FAST* family always performs the best. Therefore, to show the superiority of our approach, we run all *FAST* algorithms and select the best one for each project. In other words, when comparing APFD, we keep the highest APFD, and when comparing time cost, we keep the lowest time cost. Moreover, due to the randomness in *FAST*, for each subject we applied each of these approaches 10 times and used their median effectiveness and efficiency results. Regarding the time cost, the same as

Section 5, we measure the efficiency of a TCP approach by including its preparation time, i.e., the preparation time used in *FAST*.<sup>4</sup>

### 7.1.1 FAST Results on Seeded Faults

The results of *FAST* are shown by the tenth and twelfth columns in Table 7 (Appendix C), available in the online supplemental material. Due to space limit, we do not present the results of all the five *FAST* algorithms, but the largest APFD value and smallest time cost among them for each subject. Note that usually a *FAST* algorithm cannot achieve both the largest APFD value and the smallest time cost. As the APFD results and time cost of AGA is already given by the eighth and ninth columns, we use column  $Win_{APFD}$  and column  $Win_{Time}$  to show whether  $APFD_{AGA} \geq APFD_{FAST}$  and  $Time_{AGA} < Time_{FAST}$ , respectively.

Regarding to APFD values, the AGA approach is much better than *FAST* in all subjects. More specifically, the differences between them are from 0.0456 to 0.3039, and 0.1702 on average. To statistically check their differences, we follow the similar procedure as above. We first use Shapiro-Wilk test to check the normality of residuals, and the  $p$ -value in AGA and *FAST* is 0.328 and 0.137, which cannot reject the hypothesis that they are normally distributed. Then, taken project size as a control variable, the Analysis of Covariance (ANCOVA) shows that  $p$ -value  $< 2 \times 10^{-16}$ , indicating the statistically significant difference between AGA and *FAST*. Moreover, the effect size (Cohen's  $d$ ) is 2.96 (huge effect) and Tukey's range test shows that the 95% confidence interval of their difference is [0.149, 0.192]. To sum up, AGA significantly outperforms *FAST* in terms of APFD because *FAST* algorithms are designed to sacrifice prioritization accuracy to achieve high efficiency by using hash signatures.

Regarding to the time cost, the time cost of AGA outperforms *FAST* on 52 out of 55 open-source subjects, and the speedup ratio of AGA over *FAST* is 4.29X. To statistically check their differences, we follow the similar procedure as above. We first use Shapiro-Wilk test to check the normality of residuals, and the  $p$ -value in AGA and *FAST* is  $4.92 \times 10^{-15}$  and  $4.05 \times 10^{-15}$ , which reject the hypothesis that they are normally distributed. Therefore, we use the proportional odds regression [40] and include project size as a control variable. The results show that the  $p$ -value of "group" is  $4.250 \times 10^{-4}$ , indicating significant difference between AGA and *FAST*, and the effect size (Cohen's  $d$ ) is 0.286 (medium effect). That is, the proposed AGA is more efficient to *FAST* (with 4.29X speedup ratio). This is a *surprising result* because AGA can even be faster than a technique that is designed to sacrifice effectiveness to reduce time cost. We also present the violin plot with included box plot in Fig. 3. On larger projects, the speedup ratios are smaller, which means *FAST* also scales up well on large-size projects, whereas it is less efficient than AGA.

4. The previous work *FAST* [18] separated their total running time into preparation time and prioritization time in their evaluation. However, preparation happens only once in BB approaches while not in WB approaches, because the input of BB approaches is test code. Given updated source code but out-of-date coverage information (from the previous version), we need not prioritize again and TCP results will not change. Otherwise, with updated coverage information, the whole process (including preparation) has to be repeated.

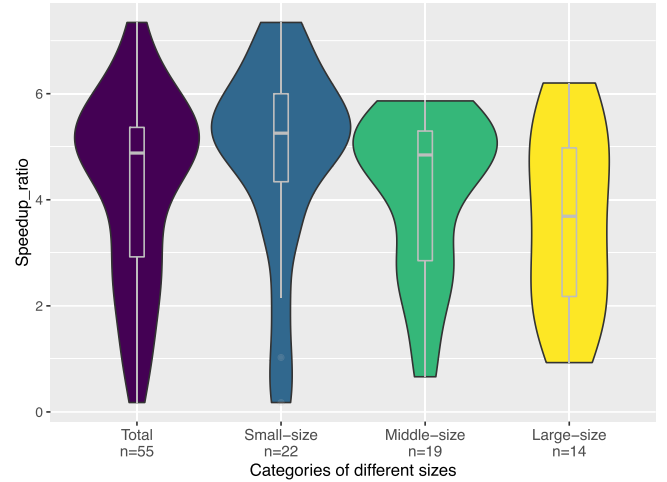


Fig. 3. Speedup ratios distribution of AGA over *FAST* on open-source projects.

### 7.1.2 FAST Results on Real Faults

Besides, as *FAST* is evaluated by some subjects of Defects4J [44] in the previous work [18], we apply the AGA approach to these subjects by reusing their artifact package (including subjects and code) for fair comparison. Moreover, we add the experiment on Mockito, which is also in Defects4J but does not appear in the experiment of *FAST*. Defects4J is the largest real-fault benchmark, so this experiment complements the previous experiments on seeded faults and can evaluate AGA on real faults. The comparison results are given by Table 4, where  $Win_{APFD}$  and  $Win_{Time}$  show whether the proposed AGA approach outperforms *FAST* in terms of APFD and time cost, respectively. From this table, AGA is more effective than *FAST* algorithms on 5 out of 6 projects and it achieves better time efficiency on all 6 projects (with 5.24X as average speedup ratio), which indicates the superiority of AGA. Also, from this experiment, we show that AGA is superior on real faults, too.

Actually, it is worth pointing out that as the authors of *FAST* [18] stated, no single *FAST* algorithm can be the best, which means the most effective algorithm in *FAST* may lead to somewhat higher time cost and the most efficient algorithm in *FAST* may lead to somewhat lower APFD value. That is, the results of *FAST* in Table 7 (Appendix C), available in the online supplemental material, and Table 4 are not results of one *FAST* algorithm, but the best results of all *FAST* algorithms. Moreover, even compared with these results, AGA is still promising considering both effectiveness and efficiency.

Considering the advantageous of AGA over *FAST*, it is interesting to analyze the secrets behind the observation. *FAST* approach achieves the efficiency improvement by using the algorithms used in big data domain to summarize the key information in coverage, but suffers from the effectiveness loss to some extent because some information is missing in the summarization. AGA consists of two parts, time complexity reduction and iteration number reduction. In particular, the former part is to use some extra data structures (e.g., indices) to summarize the coverage information of each test case, i.e., the statements covered by each test. With these data structures, AGA does not need to scan the coverage table whenever a test case is selected, and thus the



TABLE 4  
Results of Some Defects4J Projects

Projects	FAST		AGA				APFD_range	GAF	
	APFD	Time	APFD	Win <sub>APFD</sub>	Time	Win <sub>Time</sub>		APFD	Time
Closure	0.5219	11.7302	0.4347		2.0408	✓	0.0006	0.4354	9.4078
Math	0.5471	5.2600	0.6992	✓	0.9586	✓	0.0000	0.6992	7.4710
Lang	0.5627	0.4280	0.6094	✓	0.1513	✓	0.0000	0.6094	0.2150
Time	0.5463	2.8633	0.5469	✓	0.5042	✓	0.0034	0.5436	0.9870
Chart	0.5264	5.1456	0.7128	✓	0.9030	✓	NA	0.7128	5.0767
Mockito	0.5197	2.7311	0.5975	✓	0.4537	✓	0.0014	0.5961	2.7539
<b>Total</b>				5		6			

time cost of AGA reduces but its effectiveness maintains. To sum up, *FAST* suffers from effectiveness loss because it uses simplified information, while AGA does not because it uses the same information as before but in an easy-to-access way.

Moreover, similar to Table 2, we compute the gaps between the highest and lowest APFD among all iteration numbers for Defects4J subjects, which are shown in Column “APFD\_range” of Table 4. The range of “Chart” is marked as “NA” because it has only one iteration. As we can see, the gaps are extremely small, which also confirms the conclusion in Section 3.

Additionally, Column “GAF” of Table 4 shows the results of GA-first. AGA is much more efficient than GAF while achieves larger APFD, which is consistent with the conclusion in Section 6.3.

*Conclusion:* Surprisingly, AGA can achieve 4.29X speedup ratio compared to *FAST*, which targets improving time efficiency while sacrificing effectiveness. At the same time, the experimental results show that AGA is significantly better than *FAST* in terms of APFD values, and the average difference between them is 0.1702.

## 7.2 Comparison With Other TCP Techniques

Although only *FAST* has a close goal to ours, to better evaluate AGA, we also compare it with more representative TCP techniques. In particular, in this study we use the following TCP techniques whose input is only coverage information and which have been widely used in the literature [18], [31].

- *ART-D* [10] is a family of adaptive random-based TCP techniques guided by coverage information. At each iteration, a candidate set is dynamically created by randomly picking test cases from the set of not-yet-prioritized test cases as long as they can increase coverage. The test case in the candidate set that is the farthest away from the set of prioritized test cases is selected.
- *GA-S (Additional Spanning)* [54] is a variant of GA that at each iteration picks the test case that covers the largest number of uncovered elements among those in the “spanning set”. Here, an element subsumes another if covering the former guarantees covering the latter: The notion of a spanning set denotes the subset of non-subsumed elements.
- *GE* [8] is a genetic algorithm, which is a representative of search-based prioritization techniques and is

evaluated to be effective. In each iteration, it uses a fitness function to select individuals and then applies crossover and mutation operators to generate new individuals. Specifically, an individual (a sequence) is encoded as an array where each value indicates the position of a test case; The fitness function is defined by Baker’s linear ranking algorithm [55]; The crossover operator selects two parents and each of the two offspring is formed by combining the first several values in one parent and the remaining values in the other parent; The mutation operator randomly selects two values in an individual and exchanges their positions.

In this section, we reuse the implementation of *ART-D*, *GA-S*, and *GE* in [18], [22] and compare them with AGA on the 55 open-source projects. Considering the randomness of these techniques, each of them is run 10 times. The remaining setting of this experiment is the same as Section 5. Due to the space limit of Table 7 (Appendix C), available in the online supplemental material, we put the results in Table 5. In Table 5, each row represents one project, and the running time and APFD of AGA, *ART-D*, *GA-S*, and *GE* are shown separately.

The average speedup ratio of AGA over *ART-D* is 144.58X. Moreover, in all 55 projects, the APFD values of AGA are larger than *ART-D*, and the average APFD difference is 0.1384. That is, AGA always outperforms *ART-D* in terms of both effectiveness and efficiency. The average speedup ratio of AGA over *GA-S* is 182.27X. In 54 out of 55 projects, the APFD values of AGA is larger than *GA-S*, and the average APFD difference is 0.0708. The average speedup ratio of AGA over *GE* is 285.91X. In 50 out of 55 projects, the APFD values of AGA are larger than *GE*, and the average APFD difference is 0.0459. That is, compared with the three TCP techniques, our proposed AGA achieves both effectiveness and efficiency. Moreover, the time cost and APFD values of the compared TCP techniques distribute in a larger range than AGA, indicating that the latter can achieve stably promising performance.

*Conclusion:* As AGA aims to largely improve the TCP efficiency while preserving the high-effectiveness of GA, it outperforms *ART-D*, *GA-S*, and *GE* in terms of both efficiency and effectiveness.

## 8 INDUSTRIAL CASE STUDY

To show the practical usage of our approach, we then conducted an industrial case study as follows.

TABLE 5  
Comparison With Other TCP Techniques on Open-Source Subjects

Project	AGA		ART-D		GA-S		GE	
	Time	APFD	Time	APFD	Time	APFD	Time	APFD
$S_1$	0.0157	0.9070	0.0832	0.8440	0.5878	0.8812	0.0650	0.8747
$S_2$	0.0058	0.8380	0.0224	0.7508	0.3312	0.7870	0.0660	0.8373
$S_3$	0.0222	0.8848	0.0196	0.7681	0.1593	0.8108	0.0900	0.8705
$S_4$	0.0789	0.8509	0.0478	0.5933	0.1079	0.6555	0.4440	0.8061
$S_5$	0.0089	0.8527	0.1361	0.7405	0.6358	0.7662	0.9430	0.8310
$S_6$	0.0046	0.8101	0.0957	0.6909	0.2560	0.6800	0.9140	0.7935
$S_7$	0.0076	0.9059	0.0150	0.7711	0.1116	0.8255	0.0970	0.8855
$S_8$	0.0180	0.8898	0.0752	0.7459	2.7672	0.8153	0.1160	0.8985
$S_9$	0.3363	0.9144	53.9033	0.8821	33.7447	0.8995	38.4580	0.8949
$S_{10}$	0.0247	0.9518	0.5895	0.8009	0.9681	0.9017	0.8060	0.9172
$S_{11}$	0.0553	0.8766	0.4818	0.7640	10.6250	0.7950	0.4070	0.8501
$S_{12}$	0.0106	0.8864	0.0708	0.7589	0.3043	0.8375	0.2970	0.8667
$S_{13}$	0.0385	0.8615	0.1763	0.7816	2.0406	0.7667	0.2000	0.7783
$S_{14}$	0.0385	0.9188	0.6426	0.7800	0.8556	0.8779	0.6880	0.9011
$S_{15}$	0.0055	0.8582	0.1225	0.7285	0.2513	0.7117	0.7360	0.8470
$S_{16}$	0.0152	0.8031	0.3687	0.6578	0.9745	0.6509	3.0660	0.7760
$S_{17}$	0.1317	0.9183	7.8831	0.8316	7.8195	0.8893	2.7440	0.8785
$S_{18}$	0.0226	0.9028	0.1431	0.7783	1.5687	0.8039	0.3120	0.8950
$S_{19}$	0.6995	0.9033	34.4428	0.7553	374.3919	0.8381	4.1710	0.8528
$S_{20}$	0.0469	0.8013	0.2180	0.7021	36.2190	0.7293	0.3410	0.8185
$S_{21}$	0.0248	0.8328	0.4520	0.7156	1.5459	0.7676	1.5660	0.7899
$S_{22}$	0.0215	0.8642	0.2999	0.7660	0.8994	0.8151	0.5080	0.8570
$S_{23}$	0.0962	0.8198	4.8919	0.6979	4.2931	0.7782	3.6190	0.7866
$S_{24}$	0.0187	0.9858	0.0420	0.9213	12.3507	0.9857	0.1230	0.9860
$S_{25}$	1.0579	0.8401	70.1140	0.8099	387.4319	0.8426	3.5190	0.8283
$S_{26}$	0.0294	0.8339	0.4066	0.6026	0.3047	0.6927	2.6380	0.7833
$S_{27}$	0.0303	0.9614	0.0254	0.7695	0.2441	0.8579	0.2200	0.9501
$S_{28}$	0.1132	0.9164	11.6980	0.7642	4.6810	0.8734	8.5760	0.8534
$S_{29}$	2.3900	0.9490	130.3359	0.8254	7955.7540	0.9180	3.7060	0.9131
$S_{30}$	0.1804	0.9617	6.9994	0.8572	28.4835	0.9202	1.6830	0.9285
$S_{31}$	0.3170	0.9426	46.5734	0.8342	15.9675	0.9191	11.7620	0.9065
$S_{32}$	0.1270	0.8911	22.8351	0.7165	3.1369	0.8231	53.4990	0.7696
$S_{33}$	0.0921	0.8662	11.0712	0.6612	5.5717	0.7356	23.6020	0.7568
$S_{34}$	0.7903	0.9328	120.8907	0.8271	128.1646	0.8861	21.5730	0.8780
$S_{35}$	0.3631	0.9467	32.4553	0.7277	20.0292	0.8694	5.1690	0.9105
$S_{36}$	23.7849	0.9371	2,397.2400	0.8207	1,976.4397	0.8597	11.6880	0.8570
$S_{37}$	0.3582	0.8507	33.6262	0.6876	23.6568	0.7621	1,979.6000	0.7165
$S_{38}$	0.2794	0.8657	115.8419	0.7753	35.4853	0.7747	337.0780	0.8072
$S_{39}$	2.2078	0.9545	1,114.6899	0.7931	265.5048	0.9289	168.0300	0.7733
$S_{40}$	0.5942	0.9244	118.3193	0.7621	17.4111	0.8672	112.3900	0.8156
$S_{41}$	0.1562	0.9106	19.2647	0.7195	4.2101	0.8454	40.8410	0.8116
$S_{42}$	0.0287	0.8569	0.0123	0.7409	0.2401	0.8176	0.1300	0.8649
$S_{43}$	0.2558	0.8924	32.1075	0.7331	38.5073	0.7915	24.0450	0.8208
$S_{44}$	0.8465	0.9240	162.8027	0.8437	147.3191	0.9090	34.0820	0.8864
$S_{45}$	0.0597	0.8464	0.8793	0.6822	7.4841	0.7811	4.6150	0.7918
$S_{46}$	0.0858	0.8656	3.7066	0.6834	17.5147	0.7494	4.0400	0.8003
$S_{47}$	0.0629	0.8750	0.7825	0.6884	1.0181	0.7588	3.6900	0.8225
$S_{48}$	0.0025	0.7939	0.0095	0.6455	0.0386	0.7089	0.0790	0.7873
$S_{49}$	0.1120	0.8009	0.8459	0.6463	5.6348	0.7600	4.3070	0.7884
$S_{50}$	0.0047	0.8517	0.0164	0.5108	0.0585	0.6855	0.1830	0.8559
$S_{51}$	1.1036	0.8671	203.2114	0.6595	505.1426	0.8080	61.9450	0.7261
$S_{52}$	0.4052	0.9542	26.1336	0.8454	200.9616	0.8889	20.1360	0.9233
$S_{53}$	0.2298	0.8710	13.8199	0.6956	32.9077	0.8101	33.7620	0.7768
$S_{54}$	2.6648	0.9089	3,373.5142	0.7578	252.3093	0.8478	13,384.3880	0.7902
$S_{55}$	18.1040	0.9544	55,120.6523	0.8613	4536.0047	0.9292	13,516.8050	0.8747

Baidu is a famous Internet service provider with over 600M monthly active users. In their regression testing infrastructure, test case prioritization is frequently needed and they have been adopting Greedy Additional (GA) strategy for a long time because of its simple idea and relatively high effectiveness. However, they often complain about the long

running time of GA, which deviates from the original intention of test case prioritization, that is to accelerate the process of detecting faults.

To check the performance of our AGA approach in real-world scenarios, we collected 22 versions of five industrial projects from Baidu, each of which is taken as a subject in



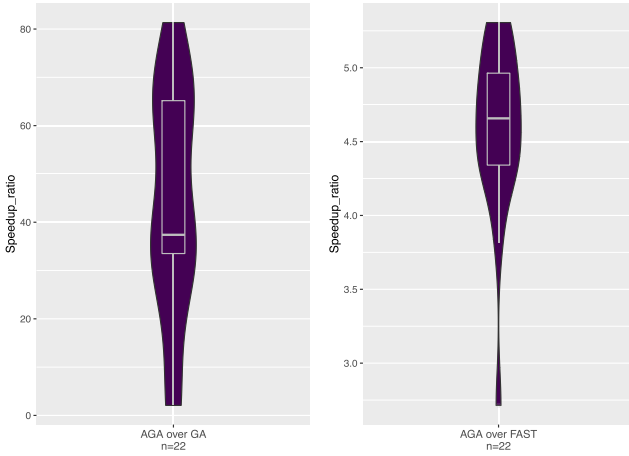


Fig. 4. Speedup ratios distribution of AGA over GA and *FAST* on industrial projects.

this study. More specifically, these subjects are collected from December 2017 to February 2018 and October 2018 to November 2018, and all of them are written in C. As shown in the first three columns in Table 8 (Appendix D), available in the online supplemental material, we summarize the SLOC and number of test cases of each subject. The SLOCs range from 20K to 500K while the numbers of test cases range from 202 to 4,246. Besides, we used C-Cover [56] to collect statement coverage for each industrial subject.

In Table 8 (Appendix D), available in the online supplemental material, we report the time cost of GA, AGA, and *FAST*, respectively. When the time cost of AGA is less than GA, we mark it with ✓. As we can see, in all 22 subjects, the time cost of AGA is much lower than that of GA, and the speedup ratio is 44.27X on average. In general, our AGA approach is demonstrated to be efficient on industrial subjects from Baidu. For example, for the subject  $\mathcal{I}_1$ , its original prioritization time is larger than 29,000 seconds, which may be unbearable in practice. However, through AGA, the prioritization time is reduced to less than 360 seconds. On the other hand, the surprisingly high efficiency of AGA also indicates the ubiquitous existence of many redundant accesses of data in industrial projects. We also present the violin plot with included box plot in Fig. 4 to show the distribution and variation. As we can see, on most projects, AGA has a large improvement compared to GA.

Provided adjacency matrix as input, the average speedup ratio of AGA over GA is 61.43X.

After we report the results, developers in Baidu verified (1) the time cost of our implementation of GA is close to their inner implementations, and (2) the speedup ratio is significant and our technique improves their prioritization efficiency, because their implementation only works on small projects, not large projects.

In addition, we also compared our approach with *FAST*, and the experimental setup is the same with Section 7. Surprisingly, AGA outperforms *FAST* again. Specifically, on all 22 subjects, AGA is faster than *FAST* and the average speedup ratio is 4.58X. To statistically check their differences, we follow the similar procedure as above. We first use Shapiro-Wilk test to check the normality of residuals, and the *p-value* in AGA and *FAST* is  $5 * 10^{-4}$  and  $2.5 * 10^{-3}$ , which reject the hypothesis that they are normally distributed. Similarly to

the above, proportional odds regression [40] is used and we introduce a variable “group” representing AGA and *FAST* and take project size as a control variable. The results show that the *p-value* of “group” is  $1.35 * 10^{-6}$ , indicating significant difference between AGA and *FAST*, and the effect size (Cohen’s *d*) is 1.28 (very large effect). We also present the violin plot with included box plot in Fig. 4. Again, we find that on most projects, AGA has a large improvement than *FAST*. In [18], the authors proposed *FAST* to solve the scalability problem of TCP techniques with the decrease of effectiveness. Their approach is evaluated to be efficient when the project size grows up rapidly. However, our AGA approach is even more efficient than *FAST*, and this means AGA may scale up better and is practical in real-world scenarios. Also, recall that when we compare AGA with *FAST* on open-source projects, the *p-value* is larger and the effect size is smaller than here, and we conjecture that this is due to the relatively small sizes of open-source projects.

Additionally, besides *FAST*, which targets the TCP efficiency problem, we also compare AGA with other more general TCP techniques as we have done in Section 7. Specifically, we run ART-D, GA-S, and GE on the 22 subjects. Considering the randomness of these techniques, each of them is run 10 times. The results are shown in Table 8 (Appendix D), available in the online supplemental material. As we can see, these techniques are much slower than AGA, even GA. The average speedup ratio of AGA over ART-D, GA-S, and GE is 993.37X, 4230.53X, and 123.25X, respectively.

It is worth noting that on one hand, Baidu is sensitive to the positions of detected faults in history, thus these positions are not available to us. On the other hand, they only provide coverage data after desensitization and we do not have access to the source code of these subjects (to create mutants) due to the confidential policy. As a result, we cannot compare the effectiveness of the three approaches in terms of APFD in industrial subjects.

*Conclusion:* Our AGA approach achieves 44.27X speedup ratio compared to GA. AGA even outperforms *FAST* in terms of time efficiency (4.58X), and the difference is statistically significant. This indicates that AGA is practical in real-world scenarios.

## 9 DISCUSSION

*Space Comparison.* From the space complexity analysis in Section 2, AGA consumes at most twice more space than GA, which is acceptable in practice. Moreover, AGA does not require high performance servers, e.g., the time cost of AGA on the two largest open-source projects (i.e., commons-math & camel-core) is only 153.42s and 187.82s (on a personal computer whose Intel Core-i5 with 8GB memory), almost the same as Table 7 (Appendix C), available in the online supplemental material.

*Impact of Seeded Faults and Real Faults.* Previous work [24], [57], [58] has explored the relationship between seeded faults and real faults and they may have different characteristics, which has potential influence on the evaluation on test case prioritization, fault localization, etc. In this paper,

we evaluate our approach on both 55 open-source subjects with seeded faults and Defects4J dataset with real faults. The high performance of AGA on both of them can illustrate its superiority well.

*Discussion on Other TCP Approaches.* Researchers have put dedicated efforts in TCP and have proposed a large number of TCP techniques since then. Many approaches take other information rather than coverage information (e.g., test inputs, test outputs, mutants) as input, so they are in different dimensions. However, even taken all kinds of approaches into consideration, the GA approach remains one of the most effective strategies in terms of fault-detection rate [7], [8], [10], [18]. So, we target GA in this paper and AGA can be better than other approaches.

## 10 RELATED WORK

Test case prioritization attracts much attention since this problem was raised at the end of the 20th century, and the work on test case prioritization can be classified into prioritization algorithms [8], [10], [59], [60], [61], [62], coverage criteria used in prioritization [3], [5], [28], [63], [64], [65], [66], [67], [68], measurement used to estimate prioritization effectiveness [2], [5], [69], and empirical studies [1], [3], [5], [12], [31], [70], [71], [72], [73], [74]. Moreover, a number of surveys on test case prioritization are also given in the literature [75], [76], [77]. For example, Catal *et al.* [76] conducted a systematic study of TCP techniques in 2001-2011 including 120 papers published in that time period. Due to the space limit, we do not list all the prioritization work here, but introduce some very recently published work. Di *et al.* [78] proposed Hyper-volume-based Genetic Algorithm to prioritize test cases using multiple test coverage criteria. Azizi *et al.* [79] proposed a graph-based framework to map the prioritization problem to a graph traversal algorithm. Chen *et al.* [80] gave an adaptive random sequences approach based on clustering techniques using black-box information. Different from them, our work targets the effective GA algorithm and attempts to solve its efficiency problem.

Moreover, some researchers noticed the efficiency problem of TCP and began to work on it. Henard *et al.* [12] said, “if prioritization takes too long, then it eats into the time available to run the prioritized test suite.” That is, for large software, it is necessary to take the scalability of TCP into consideration. Marijan *et al.* [81] proposed ROCKET to prioritize test cases based on historical failure data, test execution time and domain-specific heuristics to improve the efficiency in the scenario of continuous integration. Knauss *et al.* [82] proposed to analyze the correlation between test failures and source code changes to rapidly prioritize test cases. Elbaum *et al.* [13] introduced two techniques that use readily available test execution history data to determine what test cases are worth executing and execute them with higher priority. Recently, Miranda *et al.* [18] introduced the FAST techniques to provide similarity-based test case prioritization techniques with scalable improvements. Our work is related to the above work because all of them target TCP efficiency problem. However, the above work either does not take advantage of the coverage information which results in lower effectiveness or addresses the efficiency problem alone without balancing or even

sacrificing the effectiveness. That is, to our best knowledge, none of the existing work can improve the efficiency of GA while maintaining its widely-recognized effectiveness. Our work achieves this goal and AGA is particularly advantageous for large-scale industrial projects.

## 11 CONCLUSION

In this paper, we make a deep analysis of the Greedy Additional algorithm (GA) for test case prioritization (TCP) problem and propose AGA to improve its efficiency while preserving effectiveness. On one hand, we find the redundant data accesses in GA and take the use of extra data structures to cut down them, which leads to an optimized time complexity from  $\mathcal{O}(m^2n)$  to  $\mathcal{O}(kmn)$  given  $n > m$ , where  $m$  is the number of test cases,  $n$  is the number of program elements, and  $k$  is the iteration number. On the other hand, we notice the impacts of iteration numbers on the effectiveness and efficiency of GA and propose to reduce it to a relatively small value to improve efficiency while preserving effectiveness. Overall, we achieve an  $\mathcal{O}(mn)$  algorithm for prioritization.

We performed comprehensive experiments on 55 open-source projects to show the effectiveness and efficiency of AGA. On one hand, AGA can achieve the same average effectiveness as the GA approach, whose performance is considered to be high, and at the same time, the efficiency of AGA is much higher than GA. Specifically, our AGA approach can achieve 5.95X/27.72X speedup ratio over GA on average on two input formats. On the other hand, compared with FAST, which was recently proposed to solve the TCP efficiency problem while sacrificing effectiveness to some extent, AGA achieves 0.1702 higher APFD values on average and surprisingly the average speedup ratio of AGA over FAST is 4.29X.

Additionally, we conducted an industrial case study on 22 industrial subjects, collected from Baidu, which is a famous Internet service provider with over 600M monthly active users. The experimental results show that the average speedup ratios of AGA over GA and FAST are 44.27X/61.43X and 4.58X (with significant difference and very large effect), respectively.

To the best of our knowledge, this is the first attempt to alleviating the efficiency problem of the Greedy Additional TCP approach while maintaining its effectiveness. It is worth noting that the efficiency of TCP algorithm is especially important when software becomes larger, that is to say, in real-world scenarios. Our empirical evidence indicates that AGA is particularly more advantageous for large-scale industrial projects.

## ACKNOWLEDGMENTS

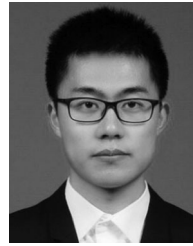
The authors would like to thank all the reviewers for their valuable comments and suggestions.

## REFERENCES

- [1] G. Rothermel, R. H. Untch, C. Chu, and M. J. Harrold, “Prioritizing test cases for regression testing,” *IEEE Trans. Softw. Eng.*, vol. 27, no. 10, pp. 929–948, Oct. 2001.
- [2] S. Elbaum, A. Malishevsky, and G. Rothermel, “Incorporating varying test costs and fault severities into test case prioritization,” in *Proc. Int. Conf. Softw. Eng.*, 2001, pp. 329–338.
- [3] S. Elbaum, A. G. Malishevsky, and G. Rothermel, “Test case prioritization: A family of empirical studies,” *IEEE Trans. Softw. Eng.*, vol. 28, no. 2, pp. 159–182, Feb. 2002.

- [4] X. Qu, M. B. Cohen, and G. Rothermel, "Configuration-aware regression testing: An empirical study of sampling and prioritization," in *Proc. Int. Symp. Softw. Testing Anal.*, 2008, pp. 75–86.
- [5] G. Rothermel, R. H. Untch, C. Chu, and M. J. Harrold, "Test case prioritization: An empirical study," in *Proc. IEEE Int. Conf. Softw. Maintenance*, 1999, pp. 179–188.
- [6] W. E. Wong, J. R. Horgan, S. London, and H. Agrawal, "A study of effective regression testing in practice," in *Proc. Int. Symp. Softw. Rel. Eng.*, 1997, pp. 264–274.
- [7] L. Zhang, D. Hao, L. Zhang, G. Rothermel, and H. Mei, "Bridging the gap between the total and additional test-case prioritization strategies," in *Proc. Int. Conf. Softw. Eng.*, 2013, pp. 192–201.
- [8] Z. Li, M. Harman, and R. M. Hierons, "Search algorithms for regression test case prioritization," *IEEE Trans. Softw. Eng.*, vol. 33, no. 4, pp. 225–237, Apr. 2007.
- [9] S. Lin, "Computer solutions of the traveling salesman problem," *Bell Syst. Tech. J.*, vol. 44, no. 10, pp. 2245–2269, 1965.
- [10] B. Jiang, Z. Zhang, W. K. Chan, and T. H. Tse, "Adaptive random test case prioritization," in *Proc. IEEE/ACM Int. Conf. Automated Softw. Eng.*, 2009, pp. 233–244.
- [11] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction to Algorithms*. Cambridge, MA, USA: MIT Press, 2009.
- [12] C. Henard, M. Papadakis, M. Harman, Y. Jia, and Y. L. Traon, "Comparing white-box and black-box test prioritization," in *Proc. IEEE/ACM Int. Conf. Softw. Eng.*, 2016, pp. 523–534.
- [13] S. Elbaum, G. Rothermel, and J. Penix, "Techniques for improving regression testing in continuous integration development environments," in *Proc. ACM SIGSOFT Int. Symp. Found. Softw. Eng.*, 2014, pp. 235–245.
- [14] M. V. Mäntylä, B. Adams, F. Khomh, E. Engström, and K. Petersen, "On rapid releases and software testing: A case study and a semi-systematic literature review," *Empirical Softw. Eng.*, vol. 20, no. 5, pp. 1384–1425, 2015.
- [15] A. Memon, Z. Gao, B. Nguyen, S. Dhanda, E. Nickell, R. Siemborski, and J. Micco, "Taming google-scale continuous testing," in *Proc. Int. Conf. Softw. Eng. Softw. Eng. Pract. Track*, 2017, pp. 233–242.
- [16] A. Kumar, "Development at the speed and scale of google," in *Proc. Int. Softw. Dev. Conf.*, 2010.
- [17] Y. Lu *et al.*, "How does regression test prioritization perform in real-world software evolution?," in *Proc. IEEE/ACM Int. Conf. Softw. Eng.*, 2016, pp. 535–546.
- [18] B. Miranda, E. Cruciani, R. Verdecchia, and A. Bertolino, "Fast approaches to scalable similarity-based test case prioritization," in *Proc. Int. Conf. Softw. Eng.*, 2018, pp. 222–232.
- [19] R. A. Fisher, "Statistical methods for research workers," in *Breakthroughs in Statistics*. Berlin, Germany: Springer, 1992, pp. 66–70.
- [20] Q. Luo, K. Moran, L. Zhang, and D. Poshyvanyk, "How do static and dynamic test case prioritization techniques perform on modern software systems? An extensive study on GitHub projects," *IEEE Trans. Softw. Eng.*, vol. 45, no. 11, pp. 1054–1080, Nov. 2018.
- [21] J. Zhou, J. Chen, and D. Hao, "Parallel test prioritization," *ACM Trans. Softw. Eng. Methodol.*, vol. 31, no. 1, pp. 1–50, 2021.
- [22] J. Chen *et al.*, "Optimizing test prioritization via test distribution analysis," in *Proc. ACM Joint Meeting Eur. Softw. Eng. Conf. Symp. Found. Softw. Eng.*, 2018, pp. 656–667.
- [23] S. Wang, J. Nam, and L. Tan, "QTEP: Quality-aware test case prioritization," in *Proc. Joint Meeting Found. Softw. Eng.*, 2017, pp. 523–534.
- [24] R. Just, D. Jalali, L. Inozemtseva, M. D. Ernst, R. Holmes, and G. Fraser, "Are mutants a valid substitute for real faults in software testing?," in *Proc. ACM SIGSOFT Int. Symp. Found. Softw. Eng.*, 2014, pp. 654–665.
- [25] J. H. Andrews, L. C. Briand, and Y. Labiche, "Is mutation an appropriate tool for testing experiments?," in *Proc. Int. Conf. Softw. Eng.*, 2005, pp. 402–411.
- [26] H. Do and G. Rothermel, "On the use of mutation faults in empirical assessments of test case prioritization techniques," *IEEE Trans. Softw. Eng.*, vol. 32, no. 9, pp. 733–752, Sep. 2006.
- [27] Y. Lou, D. Hao, and L. Zhang, "Mutation-based test-case prioritization in software evolution," in *Proc. IEEE Int. Symp. Softw. Rel. Eng.*, 2015, pp. 46–57.
- [28] H. Mei, D. Hao, L. Zhang, L. Zhang, J. Zhou, and G. Rothermel, "A static approach to prioritizing JUnit test cases," *IEEE Trans. Softw. Eng.*, vol. 38, no. 6, pp. 1258–1275, Nov./Dec. 2012.
- [29] M. J. Arafeen and H. Do, "Test case prioritization using requirements-based clustering," in *Proc. IEEE Int. Conf. Softw. Testing Verification Validation*, 2013, pp. 312–321.
- [30] H. Do, S. Mirarab, L. Tahvildari, and G. Rothermel, "The effects of time constraints on test case prioritization: A series of controlled experiments," *IEEE Trans. Softw. Eng.*, vol. 36, no. 5, pp. 593–617, Sep./Oct. 2010.
- [31] Q. Luo, K. Moran, and D. Poshyvanyk, "A large-scale empirical comparison of static and dynamic test case prioritization techniques," in *Proc. ACM SIGSOFT Int. Symp. Found. Softw. Eng.*, 2016, pp. 559–570.
- [32] Pit mutation testing. 2021. [Online]. Available: <http://pitest.org/atlassian/clover-bitbucket>.
- [33] atlassian / clover – bitbucket. 2021. [Online]. Available: <https://bitbucket.org/atlassian/clover/src/default/>
- [34] R. Gopinath, C. Jensen, and A. Groce, "Mutations: How close are they to real faults?," in *Proc. IEEE Int. Symp. Softw. Rel. Eng.*, 2014, pp. 189–200.
- [35] Q. Luo, K. Moran, D. Poshyvanyk, and M. D. Penta, "Assessing test case prioritization on real faults and mutants," in *Proc. IEEE Int. Conf. Softw. Maintenance Evol.*, 2018, pp. 240–251.
- [36] M. Papadakis, C. Henard, and Y. L. Traon, "Sampling program inputs with mutation analysis: Going beyond combinatorial interaction testing," in *Proc. IEEE Int. Conf. Softw. Testing, Verification Validation*, 2014, pp. 1–10.
- [37] J. Petke, S. Yoo, M. B. Cohen, and M. Harman, "Efficiency and early fault detection with lower and higher strength combinatorial interaction testing," in *Proc. Joint Meeting Found. Softw. Eng.*, 2013, pp. 26–36.
- [38] S. S. Shapiro and M. B. Wilk, "An analysis of variance test for normality (complete samples)," *Biometrika*, vol. 52, no. 3/4, pp. 591–611, 1965.
- [39] H. B. Mann and D. R. Whitney, "On a test of whether one of two random variables is stochastically larger than the other," *Ann. Math. Statist.*, vol. 18, pp. 50–60, 1947.
- [40] P. McCullagh, "Regression models for ordinal data," *J. Roy. Statist. Soc. Ser. B (Methodological)*, vol. 42, no. 2, pp. 109–127, 1980.
- [41] J. Cohen, *Statistical Power Analysis For The Behavioral Sciences*. New York, NY, USA: Academic Press, 2013.
- [42] M. S. Bartlett, "Properties of sufficiency and statistical tests," *Proc. Roy. Soc. London. Ser. A-Math. Phys. Sci.*, vol. 160, no. 901, pp. 268–282, 1937.
- [43] J. W. Tukey, "Comparing individual means in the analysis of variance," *Biometrics*, vol. 5, pp. 99–114, 1949.
- [44] R. Just, D. Jalali, and M. D. Ernst, "Defects4J: A database of existing faults to enable controlled testing studies for java programs," in *Proc. Int. Symp. Softw. Testing Anal.*, 2014, pp. 437–440.
- [45] D. Paterson, G. M. Kapfhammer, G. Fraser, and P. McMin, "Using controlled numbers of real faults and mutants to empirically evaluate coverage-based test case prioritization," in *Proc. Int. Workshop Automat. Softw. Test*, 2018, pp. 57–63.
- [46] M. A. Hasan, M. A. Rahman, and M. S. Siddik, "Test case prioritization based on dissimilarity clustering using historical data analysis," in *Proc. Int. Conf. Inf., Commun. Comput. Technol.*, 2017, pp. 269–281.
- [47] T. B. Noor and H. Hemmati, "A similarity-based approach for test case prioritization using historical failure data," in *Proc. IEEE Int. Symp. Softw. Rel. Eng.*, 2015, pp. 58–68.
- [48] A. Haghighatkah, M. Mäntylä, M. Oivo, and P. Kuvaja, "Test case prioritization using test similarities," in *Proc. Int. Conf. Product-Focused Softw. Process Improvement*, 2018, pp. 243–259.
- [49] X. Li, W. Li, Y. Zhang, and L. Zhang, "DeepFLI: Integrating multiple fault diagnosis dimensions for deep fault localization," in *Proc. ACM SIGSOFT Int. Symp. Softw. Testing Anal.*, 2019, pp. 169–180.
- [50] X. Li and L. Zhang, "Transforming programs and tests in tandem for fault localization," *Proc. ACM Program. Lang.*, vol. 1, pp. 1–30, 2017.
- [51] S. Pearson *et al.*, "Evaluating and improving fault localization," in *Proc. IEEE/ACM Int. Conf. Softw. Eng.*, 2017, pp. 609–620.
- [52] J. Sohn and S. Yoo, "FLUCCS: Using code and change metrics to improve fault localization," in *Proc. ACM SIGSOFT Int. Symp. Softw. Testing Anal.*, 2017, pp. 273–283.
- [53] M. Zhang, X. Li, L. Zhang, and S. Khurshid, "Boosting spectrum-based fault localization using pagerank," in *Proc. ACM SIGSOFT Int. Symp. Softw. Testing Anal.*, 2017, pp. 261–272.
- [54] M. Marré and A. Bertolino, "Using spanning sets for coverage testing," *IEEE Trans. Softw. Eng.*, vol. 29, no. 11, pp. 974–984, Nov. 2003.
- [55] J. E. Baker, "Adaptive selection methods for genetic algorithms," in *Proc. Int. Conf. Genet. Algorithms Appl.*, 1985, pp. 101–111.
- [56] Bullseye testing technology. 2021. [Online]. Available: <http://www.bullseye.com/>
- [57] M. Papadakis, D. Shin, S. Yoo, and D.-H. Bae, "Are mutation scores correlated with real fault detection? A large scale empirical study on the relationship between mutants and real faults," in *Proc. IEEE/ACM Int. Conf. Softw. Eng.*, 2018, pp. 537–548.

- [58] M. Daran and P. Thévenod-Fosse, "Software error analysis: A real case study involving real faults and mutations," *ACM SIGSOFT Softw. Eng. Notes*, vol. 21, no. 3, pp. 158–171, 1996.
- [59] G. Fraser and F. Wotawa, "Test-case prioritization with model-checkers," in *Proc. IASTED Int. Conf. Web-Based Education*, 2007, pp. 267–272.
- [60] S. Yoo, M. Harman, P. Tonella, and A. Susi, "Clustering test cases to achieve effective and scalable prioritisation incorporating expert knowledge," in *Proc. Int. Symp. Softw. Testing Anal.*, 2009, pp. 201–212.
- [61] R. K. Saha, L. Zhang, S. Khurshid, and D. E. Perry, "An information retrieval approach for regression test prioritization based on program changes," in *Proc. IEEE/ACM IEEE Int. Conf. Softw. Eng.*, 2015, pp. 268–279.
- [62] Z. Ma and J. Zhao, "Test case prioritization based on analysis of program structure," in *Proc. Asia-Pacific Softw. Eng. Conf.*, 2008, pp. 471–478.
- [63] S. Elbaum, A. G. Malishevsky, and G. Rothermel, "Prioritizing Test Cases for regression testing," in *Proc. ACM SIGSOFT Int. Symp. Softw. Testing Anal.*, 2000, pp. 102–112.
- [64] H. Do, G. Rothermel, and A. Kinneer, "Empirical studies of test case prioritization in a JUnit testing environment," in *Proc. Int. Symp. Softw. Rel. Eng.*, 2004, pp. 113–124.
- [65] J. A. Jones and M. J. Harrold, "Test-suite reduction and prioritization for modified condition/decision coverage," *IEEE Trans. Softw. Eng.*, vol. 29, no. 3, pp. 195–209, Mar. 2003.
- [66] L. Zhang, J. Zhou, D. Hao, L. Zhang, and H. Mei, "Prioritizing JUnit test cases in absence of coverage information," in *Proc. IEEE Int. Conf. Softw. Maintenance*, 2009, pp. 19–28.
- [67] B. Korel, L. H. Tahat, and M. Harman, "Test prioritization using system models," in *Proc. IEEE Int. Conf. Softw. Maintenance*, 2005, pp. 559–568.
- [68] L. Mei, Z. Zhang, W. Chan, and T. Tse, "Test case prioritization for regression testing of service-oriented business applications," in *Proc. Int. Conf. World Wide Web*, 2009, pp. 901–910.
- [69] G. M. Kapfhammer and M. L. Soffa, "Using coverage effectiveness to evaluate test suite prioritizations," in *Proc. ACM Int. Workshop Empirical Assessment Softw. Eng. Lang. Technol., Held Conjunction IEEE/ACM Int. Conf. Automated Softw. Eng.*, 2007, pp. 19–20.
- [70] D. Shin, S. Yoo, M. Papadakis, and D.-H. Bae, "Empirical evaluation of mutation-based test case prioritization techniques," *Softw. Testing, Verification Rel.*, vol. 29, no. 1–2, 2019, Art. no. e1695.
- [71] H. Do, S. Mirarab, L. Tahvildari, and G. Rothermel, "An empirical study of the effect of time constraints on the cost-benefits of regression testing," in *Proc. ACM SIGSOFT Int. Symp. Found. Softw. Eng.*, 2008, pp. 71–82.
- [72] D. Hao, L. Zhang, and H. Mei, "Test-case prioritization: Achievements and challenges," *Front. Comput. Sci.*, vol. 10, no. 5, pp. 769–777, 2016.
- [73] M. G. Epitropakis, S. Yoo, M. Harman, and E. K. Burke, "Empirical evaluation of pareto efficient multi-objective regression test case prioritisation," in *Proc. Int. Sym. Softw. Testing Anal.*, 2015, pp. 234–245.
- [74] D. Hao, L. Zhang, L. Zang, Y. Wang, X. Wu, and T. Xie, "To be optimal or not in test-case prioritization," *IEEE Trans. Softw. Eng.*, vol. 42, no. 5, pp. 490–505, May 2015.
- [75] S. Yoo and M. Harman, "Regression testing minimization, selection and prioritization: A survey," *Softw. Testing, Verification Rel.*, vol. 22, no. 2, pp. 67–120, 2012.
- [76] C. Catal and D. Mishra, "Test case prioritization: A systematic mapping study," *Softw. Qual. J.*, vol. 21, no. 3, pp. 445–478, 2013.
- [77] S. Mohanty, A. A. Acharya, and D. P. Mohapatra, "A survey on model based test case prioritization," *Int. J. Comput. Sci. Inf. Technol.*, vol. 2, no. 3, pp. 1042–1047, 2011.
- [78] D. D. Nucci, A. Panichella, A. Zaidman, and A. D. Lucia, "A test case prioritization genetic algorithm guided by the hypervolume indicator," *IEEE Trans. Softw. Eng.*, vol. 46, no. 6, pp. 674–696, Jun. 2018.
- [79] M. Azizi and H. Do, "Graphite: A greedy graph-based technique for regression test case prioritization," in *Proc. IEEE Int. Symp. Softw. Rel. Eng. Workshops*, 2018, pp. 245–251.
- [80] J. Chen *et al.*, "Test case prioritization for object-oriented software: An adaptive random sequence approach based on clustering," *J. Syst. Softw.*, vol. 135, pp. 107–125, 2018.
- [81] D. Marijan, A. Gotlieb, and S. Sen, "Test case prioritization for continuous regression testing: An industrial case study," in *Proc. IEEE Int. Conf. Softw. Maintenance*, 2013, pp. 540–543.
- [82] E. Knauss, M. Staron, W. Meding, O. Söder, A. Nilsson, and M. Castell, "Supporting continuous integration by code-churn based test selection," in *Proc. Int. Workshop Rapid Continuous Softw. Eng.*, 2015, pp. 19–25.



**Feng Li** received the BS degree from Peking University in 2018. He is currently working toward the PhD degree with the School of Computer Science, Peking University. His research interests include software testing and analysis.



**Jianyi Zhou** received the BS and MS degrees from Beihang University in 2014 and 2017, respectively. He is currently working toward the PhD degree with the School of Computer Science, Peking University. His research interests include software testing and analysis.



**Yin Zhu Li** received the MS degree in computer science and technology in 2012 from Tianjin Normal University. She is currently an employee with Baidu Online Network Technology, Beijing, Company, Ltd., mainly working on automation testing. Her research interests include intelligent testing, including test case selection, test case generation, and fault localization.



**Dan Hao** received the BS degree in computer science from the Harbin Institute of Technology in 2002 and the PhD degree in computer science from Peking University in 2008. She is currently an associate professor with the School of Computer Science, Peking University, China. She was a program co-chair of ASE 2021 and SANER 2022, the general co-chair of SPLC 2018, program committees of many prestigious conferences, including ICSE, FSE, ASE, and ISSTA. Her research interests include software testing and debugging.



**Lu Zhang** received the BSc and PhD degrees in computer science from Peking University in 1995 and 2000, respectively. He is currently a professor with the School of Computer Science, Peking University, China. He was a postdoctoral researcher with Oxford Brookes University and the University of Liverpool, U.K. He was on the program committees of many prestigious conferences, including FSE, OOPSLA, ISSTA, and ASE. He was a program co-chair of SCAM 2008 and a program co-chair of ICSME 2017. His research interests

include software testing and analysis, program comprehension, software maintenance and evolution, software reuse and component-based software development, and service computing. He has been on the editorial boards of the *Journal of Software Maintenance and Evolution: Research and Practice* and *Software Testing, Verification and Reliability*.

► For more information on this or any other computing topic, please visit our Digital Library at [www.computer.org/csdl](http://www.computer.org/csdl).