

# Test Case Prioritization for Object-Oriented Software: An Adaptive Random Sequence Approach Based on Clustering\*

Jinfu Chen<sup>a</sup>, Lili Zhu<sup>a</sup>, Tsong Yueh Chen<sup>b</sup>, Dave Towey<sup>c</sup>, Fei-Ching Kuo<sup>b</sup>, Rubing Huang<sup>a</sup>, Yuchi Guo<sup>a</sup>

<sup>a</sup>(School of Computer Science and Communication Engineering, Jiangsu University, Zhenjiang, 202000, China)  
{jinfuchen, lilizhu, rbhuang, yuchiguo}@ujs.edu.cn

<sup>b</sup>(Department of Computer Science and Software Engineering, Swinburne University of Technology, Hawthorn, 3122, Australia)  
{tychen, dkuo}@swin.edu.au

<sup>c</sup>(School of Computer Science, The University of Nottingham Ningbo China, Ningbo, 315100, China)  
dave.towey@nottingham.edu.cn

---

## Abstract

Test case prioritization (TCP) attempts to improve fault detection effectiveness by scheduling the important test cases to be executed earlier, where the importance is determined by some criteria or strategies. Adaptive random sequences (ARSs) can be used to improve the effectiveness of TCP based on white-box information (such as code coverage information) or black-box information (such as test input information). To improve the testing effectiveness for object-oriented software in regression testing, in this paper, we present an ARS approach based on clustering techniques using black-box information. We use two clustering methods: (1) clustering test cases according to the number of objects and methods, using the K-means and K-medoids clustering algorithms; and (2) clustered based on an object and method invocation sequence similarity metric using the K-medoids clustering algorithm. Our approach can construct ARSs that attempt to make their neighboring test cases as diverse as possible. Experimental studies were also conducted to verify the proposed approach, with the results showing both enhanced probability of earlier fault detection, and higher effectiveness than random prioritization and method coverage TCP technique.

*Keywords:*

*Object-oriented software, Adaptive random sequence, Test cases prioritization, Cluster analysis, Test cases selection*

---

## 1. Introduction

Software testing is an important approach for ensuring the quality and reliability of software. Since the development of object-oriented (OO) technology, object-oriented software (OOS) has become widely used. However, testers may face challenges when attempting to apply traditional software testing approaches to OOS testing, due to some special characteristics of OO languages such as encapsulation, inheritance and polymorphism [1–3]. Many OOS testing approaches have been studied, including random testing (RT) [4], state-based testing [5], and sequence-based testing [6]. Among these approaches, RT has often been used in industry, partly due to its simplicity [7, 8]. Other testing approaches generally require more professional testing skills, and often focus on some specific kinds of software. A problem with the evolution of OOS is that test suites generated by these OOS testing approaches often include very large numbers of test cases, and hence execution of all of them can incur a very high cost [9–11].

In order to improve the testing efficiency of OOS in regression testing, we need to prioritize test cases to find faults as quickly as possible. Generally speaking, since only some test

inputs can detect faults, if these particular inputs could be prioritized for early execution, then the testing efficiency could be greatly improved. This kind of test case prioritization (TCP) should make it possible to detect faults earlier [12].

Current TCP techniques are developed based on white-box or black-box information [13]. The white-box information often includes program source code coverage, a program model and fault detection history; and black-box information usually includes test input information. In regression testing the white-box information is usually based on previous program versions, but the testing is done on the current version [14]. TCP techniques using black-box information do not have this problem.

Random sampling is a black-box prioritization technique, and is usually used as a benchmark for effectiveness evaluation of other prioritization techniques. In order to improve the effectiveness of random sequences and present a better prioritization benchmark in regression testing, research has resulted in a prioritization technique using Adaptive Random Sequences (ARSs) [13, 15]. ARSs can be regarded as an alternative random sequence, in which test cases are evenly spread in the input domain with the purpose of improving the performance of the random sequence. ARSs originated from the concept of Adaptive Random Testing (ART) [16–19], which is an enhanced version of RT that attempts to improve RT’s failure-detection effectiveness by evenly spreading test inputs throughout the entire in-

---

\*A preliminary version of this paper was presented at the 7th IEEE International Workshop on Program Debugging (IWPDP 2016) [21]

put domain. Adaptive random sampling can generate ARSs to make the selection of the ordered test cases as diverse across the input domain as possible [20].

ARSs have been applied to TCP for process-oriented software, based on ART techniques [15, 17, 19]. We first used the ARS technique on complex OO programs [21], and proposed an ARS approach for OOS test case prioritization. In this paper, we extend the previous work and use the notion of clustering to generate ARSs for OOS, with test cases of similar properties grouped into the same cluster, and test cases in the same cluster being different from those in other clusters. Intuitively speaking, test cases in the same cluster may have similar fault detection capability [22]. Thus test cases extracted from different clusters should have different properties, and hence should be able to detect different failures. Based on this intuition, we used cluster analysis technology to generate ARSs from different clusters, aiming to achieve an even spread of the prioritized adaptive sequence test cases across the input domain.

In this paper, we report on using method object clustering (MOClustering) and dissimilarity metric clustering (DMClustering) to generate ARSs. MOClustering forms clusters according to the number of objects and the length of method invocation sequences, using the K-means and K-medoids clustering algorithms. DMClustering uses K-medoids clustering algorithm and the structure information of test inputs to form clusters according to the Object and Method Invocation Sequence Similarity (OMISS) metric [23], which is a dissimilarity measurement for the test inputs of OO programs (based on calculation of the dissimilarity between two series of objects and between two sequences of method invocations). Additionally, a sampling strategy called MSampling (maximum sampling) is used to construct the ARSs within the MOClustering and the DMClustering frameworks. Because the proposed approach uses three clustering algorithms, three ARSs are constructed. We conducted empirical studies using seven open source subject programs, with the results showing that the proposed approaches can effectively prioritize the test cases and enhance the failure detection effectiveness. In particular, DMClustering outperforms other methods in testing large scale programs with complex structure.

The remainder of this paper is organized as follows. The research background is given in Section II. The three clustering algorithms are explained in Section III. The ARS generation algorithm is presented in Section IV. The results of our empirical studies and experimental analysis are reported in Section V. Some related work is discussed in Section VI. And the conclusion and future work are presented in Section VII.

## 2. Background

### 2.1. Regression testing

Regression testing is important for ensuring software quality and reliability. The purpose of regression testing is to ensure that the modified program still confirms to the software requirements [24]. Regression testing techniques usually involve test case reduction and test case prioritization [10, 25]. Test case reduction selects a subset of a given test suite, and aims to reduce

regression testing time by only re-running the test cases affected by code changes. Test case prioritization techniques aim to reorder test executions so as to maximize some objectives, such as detecting faults earlier or reducing the testing cost. Compared to test case reduction, test case prioritization may be a more conservative approach, because it does not discard test cases and only prioritizes them [10].

### 2.2. Cluster analysis

Cluster analysis can be used to improve software testing effectiveness, using the basic idea that test cases with similar properties be grouped into the same cluster: test cases in the same cluster are similar to each another but different from test cases in other clusters. In general, most clustering methods can be classified into one of the following five categories [26]: (1) partition methods; (2) hierarchical methods; (3) density-based methods; (4) grid-based methods; and (5) model-based methods.

### 2.3. Test Case Prioritization

The purpose of test case prioritization (TCP) is to increase the test suite’s rate of fault detection by scheduling test cases with higher priority to be executed earlier, according to some criteria. TCP can identify a permutation of a test suite, from the set of all possible permutations, that maximizes the value of a fitness function — where the function reflects a given testing goal, such as the number of detected faults. Rothermel et al. [12, 27] proposed the weighted average percentage of faults detected (APFD) as a metric to measure prioritization performance. If  $T$  represents an ordered test suite containing  $n$  test cases, and  $F$  represents a set of  $m$  failures detected by  $T$ , then  $TF_i$  represents the number of test cases executed in  $T$  before detecting fault  $i$ . The formula of APFD is defined as follows, with APFD values ranging from 0 to 1, and higher values indicating better fault detection rates.

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

Existing TCP techniques are classified as either white-box or black-box [24, 28]. Most white-box TCP techniques are based on the coverage information of the test suite for previous program versions. The white-box approaches use a selected test coverage criterion to prioritize the test suites. Test coverage criteria mainly include statement coverage, branch coverage, path coverage, method coverage and class coverage. Black-box TCP techniques usually prioritize the test suites using information associated with the test input and output information. Black-box TCP techniques mainly include combinatorial interaction testing, input model diversity and input (output) test set diameter.

### 2.4. Adaptive random sequence

Chen et al. proposed Adaptive Random Testing (ART) as an enhancement to RT [16, 17]. ART attempts to improve on RT’s failure-detection effectiveness by evenly spreading test inputs

1 throughout the entire input domain, using a similarity/dissimilarity 53  
2 metric [17, 18]. ART can be used not only to generate its 54  
3 own sequence of test cases, but also to order a given test suite  
4 to improve its chance of detecting failures earlier, with such an  
5 ordered sequence being called an Adaptive Random Sequence  
6 (ARS). Similar to ART, an ARS is also based on the idea of  
7 even spreading across the input domain — a concept that has  
8 been shown to effectively reveal failures faster. ARSs can be ap-  
9 plied to regression testing, and may be a simple, effective, and  
10 relatively low-overhead alternate to random sequences (RSs),  
11 which are commonly used in regression testing. Thus, we can  
12 use ARSs to prioritize test suites, and to enhance the perform-  
13 ance of regression testing for OOS.

## 14 2.5. Test Case Generation

15 In integration and system testing of OOS, a test case  $t$  can  
16 consist of two parts:  $t.OBJ$  and  $t.MINV$ , where  $t.OBJ$  is a list of  
17 objects and  $t.MINV$  is an ordered list of methods (representing  
18 a sequence of method invocations) in the test case. Before  
19 ordering the test cases, the test suites for regression testing must  
20 first be generated. The test suites are randomly generated in our  
21 approach. Since test cases are generated based on the class in-  
22 formation of the program under test, it is necessary to first ob-  
23 tain and analyze the class diagram. Visual Studio [29] was used  
24 to obtain the detailed class information of the subject programs,  
25 and the class diagrams.

26 The test suites were randomly generated, and the generation  
27 steps are as follows. First, the class diagram of the program  
28 under test is obtained. Based on this, the second step is to create  
29 a random number of objects, with random values assigned to each  
30 member object. Next, a random number of methods are generat-  
31 ed as the length of method sequence, and the method sequence  
32 is verified. Finally, values are assigned to the method param-  
33 eters by calling a random value generator for the corresponding  
34 data type. As a result, a test case is generated. The above steps  
35 were repeated until sufficiently many test cases were generated.

## 36 3. Clustering Algorithms

37 In this study, we used three methods to cluster test cas-  
38 es: MOClustering\_means (method object clustering with K-  
39 means), MOClustering\_medoids (method object clustering with  
40 K-medoids), and DMClustering (dissimilarity metric cluster-  
41 ing with K-medoids). MOClustering\_means and MOClustering\_ 79  
42 medoids used the Euclidean distance to calculate the dis-  
43 similarity between test cases, while DMClustering employed 80  
44 the OMISS metric to calculate the dissimilarity. In DMCluster-  
45 ing, because the OOS test inputs involved objects and methods 82  
46 rather than numerical data, the K-means could not be calculat-  
47 ed. Hence, only the K-medoids clustering algorithm was used 83  
48 in DMClustering.

### 49 3.1. Framework overview

50 Figure 1 shows the framework for our approaches. Before  
51 generating the test suites, the class diagram of the program  
52 under test is first obtained and analyzed. Then the test suites

are generated, with each test case consisting of objects and the  
methods called by these objects.

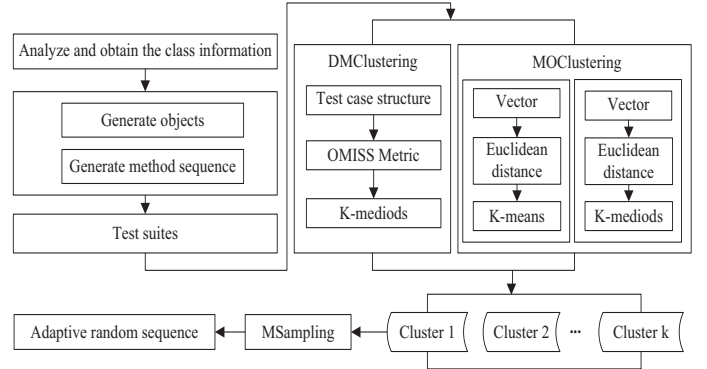


Figure 1: TCP Framework

Next, three methods are applied to cluster test cases of the  
constructed test suites. In MOClustering (method object cluster-  
ing), test cases are represented in the form of vectors, and K-  
means and K-medoids clustering algorithms are applied, using  
Euclidean distance, to cluster the test cases — MOClustering  
with K-medoids clustering algorithm is referred to as MOClus-  
tering\_medoids; and MOClustering algorithm with K-means is  
referred to as MOClustering\_means. In DMClustering, OMISS  
is used to calculate the dissimilarity between test cases, and the  
K-medoids clustering algorithm groups test cases into clusters.  
Finally, the adaptive random test sequences are generated using  
the MSampling (maximum sampling) strategy.

### 3.2. MOClustering

#### 3.2.1. Object method vector

When conducting OOS integration and system testing, typi-  
cally, a test case  $t$  will consist of a set of objects and an ordered  
list of methods. We therefore use an object method vector to  
represent a test case, defined as follows.

**Definition 1.** (object method vector,  $omv$ ): An object method  
vector of a test case is defined as an ordered pair of the number  
of its objects and the total number of methods called by all of its  
objects, denoted  $omv = \langle On, Mn \rangle$ , where  $On$  is the number of  
objects in the test input, and  $Mn$  is the total number of methods  
called by all objects that are in the test input.

For example, the object method vector for a test case  $t_1$  with  
three objects and five methods called by all objects is represent-  
ed as  $\langle 3, 5 \rangle$ .

#### 3.2.2. Distance measure

Because Euclidean distance is a natural measurement for  
distance between numerical data, it is used to measure the dis-  
tance between pairs of  $omv$ . If  $X$  is the  $omv$  of  $t_1$ , and  $Y$  is the  
 $omv$  of  $t_2$ , with  $X = \langle x_1, x_2 \rangle$  and  $Y = \langle y_1, y_2 \rangle$ , then the distance  
between  $X$  and  $Y$  is defined as:

$$d(X, Y) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2} \quad (2)$$

For example, if  $X$  is  $\langle 3, 5 \rangle$  and  $Y$  is  $\langle 3, 4 \rangle$ , then the Euclidean distance between  $X$  and  $Y$  is equal to 1, because  $d(X, Y) = \sqrt{(3-3)^2 + (5-4)^2} = \sqrt{1} = 1$ .

### 3.2.3. MOClustering\_means algorithm

In MOClustering\_means, test cases are clustered according to the numbers of objects and methods in each test case. The K-means clustering algorithm is efficient and scalable for large data sets, and was therefore used in MOClustering. The algorithm first selects  $K$  test cases as the initial data for each cluster. Each remaining test case is allocated to the closest cluster, defined by the lowest distance to the mean value of the cluster. The mean value of each cluster is then updated. This process is repeated until objects in each cluster no longer change or the sum of square error converges. After clustering, the test cases in the same cluster are expected to be similar each another, and different to those in other clusters.

MOClustering\_means is shown in Algorithm 1, and has three input parameters: *testcasepool* (the simulated input domain), *TNum* (the number of test cases to be selected from *testcasepool* to form a test suite) and  $K$  (the number of clusters to be generated). The algorithm will generate  $K$  clusters for *TNum* test cases selected from *testcasepool*. In MOClustering\_means, *TNum* test cases are first randomly selected to form a test suite that is to be prioritized; and the number of objects and methods is extracted from each chosen test case to construct the object method vectors set *OMV* for the *TNum* test cases, i.e., we construct the corresponding relationship between *OMV* and *TNum* test cases, and thus the test cases are grouped based on the corresponding clustering operation of the elements in *OMV*. Next, the first  $K$  test cases are selected as the initial cluster center of each cluster, and the mean value of each cluster updated according to Formula 3. Then, the Euclidean distance between each element of *OMV* and the mean value of each cluster are calculated, and the corresponding test case of each object method vector is assigned to the closest cluster. This is repeated until test cases in each cluster no longer change, or the sum of square error (Formula 4) converges. At this point,  $K$  clusters would have been generated and stored in the data set *clustering*.

Let  $OMV(c)$  be the set of object method vectors corresponding to cluster  $c$ . Suppose  $OMV(c) = \{omv_1, omv_2, \dots, omv_n\}$ , where  $omv_i = \langle On_i, Mn_i \rangle, i = 1, 2, \dots, n$ , where  $On_i$  is the number of objects of the test input  $t_i$  in  $c$ , and  $Mn_i$  is the sum of the number of methods called by each object of the test input  $t_i$  in  $c$ . Let  $avg(c)$  denote the mean of cluster  $c$  which is defined as a vector of two mean values shown below:

$$avg(c) = \left\langle \frac{\sum_{i=1}^n On_i}{n}, \frac{\sum_{i=1}^n Mn_i}{n} \right\rangle \quad (3)$$

Suppose that  $C$  is a cluster set, with  $C = \{c_1, c_2, \dots, c_K\}$ , and  $OMV(C)$  (or  $OMV(c_i)$ ) is the set of object method vectors corresponding to  $C$  (or  $c_i$ ), with  $OMV(C) = \bigcup_{i=1}^K OMV(c_i)$ , where  $OMV(c_i) = \{omv_{i1}, omv_{i2}, \dots, omv_{ih}\}$ . The mean value of cluster  $c_i$  —  $avg(c_i)$  — is calculated according to Formula

### Algorithm 1 MOClustering\_means (*testcasepool*, $K$ , *TNum*)

---

```

1: Construct OriginalTC = {} to store the selected test cases;
2: Construct OMV = {} to store the set of object method vectors;
3: Construct Clustering = {} to store the generated clusters;
4: Construct meanValue = {} to store the mean value of each cluster;
5: Choose TNum test cases from testcasepool randomly and add them to OriginalTC;
6: for ( $i=1$  to TNum)
7:    $On = |OriginalTC[i].Objects|$ ;  $||OriginalTC[i].Objects|$  is equal to the number of objects of OriginalTC[ $i$ ].
8:    $Mn = |OriginalTC[i].Methods|$ ;  $||OriginalTC[i].Methods|$  is equal to the number of methods of OriginalTC[ $i$ ].
9:    $OMV[i] = \langle On, Mn \rangle$ ;  $||$  The  $i$ th element of OMV is denoted by OMV[ $i$ ].
10: end for
11: Choose  $K$  elements from OMV and add the corresponding test cases to Clustering as the initial cluster center;
12: Set change = true;
13: while (change == true)
14:   Update meanValue for each cluster;
15:   for ( $i=1$  to TNum)
16:     for ( $j=1$  to  $K$ )
17:       calculate  $d(OMV[i], meanValue[j])$  according to Formula 2;
18:     end for
19:     Put the corresponding test case of OMV[ $i$ ] to the nearest cluster;
20:   end for
21:   if (each cluster keep invariant)
22:     Set change = false;
23:   else
24:     Set change = true;
25:   end if
26: end while
27: return Clustering

```

---

3. Let  $ES$  denote the sum of square error among cluster set  $C$ , which is defined as:

$$ES = \sum_{i=1}^K \sum_{j=1}^h d(omv_{ij}, avg(c_i))^2 \quad (4)$$

For example, suppose that the test suites have five test cases, and we extract the number of objects and methods from each test case to construct *OMV*:  $omv_1 = \langle 4, 3 \rangle$ ,  $omv_2 = \langle 2, 1 \rangle$ ,  $omv_3 = \langle 3, 4 \rangle$ ,  $omv_4 = \langle 1, 5 \rangle$  and  $omv_5 = \langle 3, 2 \rangle$ . Also assume that  $K$  is set to 2, and that test cases  $t_2$  and  $t_3$  are somehow chosen as the initial cluster centers. We first calculate the distance between  $omv_i (i = 1, 4, 5)$  and  $omv_2$ , and the distance between  $omv_i (i = 1, 4, 5)$  and  $omv_3$ , then put each test case into its nearest cluster. For example, since the distance between  $omv_1$  and  $omv_2$  is 2.83, and the distance between  $omv_1$  and  $omv_3$  is 1.41, then  $t_1$  should be put into cluster  $c_2$ . After the first round distribution, cluster  $c_1$  has two test cases  $t_2$  and  $t_5$ , and cluster  $c_2$  has three  $t_1, t_3$  and  $t_4$ . The mean value of the new clusters should next be updated. After the second round distribution, cluster  $c_1$  still has  $t_2$  and  $t_5$ , and cluster  $c_2$  still has three  $t_1, t_3$  and  $t_4$ . Because the clusters are the same as in the previous round, the process of clustering is completed. Figure 2 summarizes the three rounds of distribution for the above example.

### 3.2.4. MOClustering\_medoids algorithm

The K-medoids clustering algorithm randomly selects  $K$  test cases as the center points (also referred to as the representa-

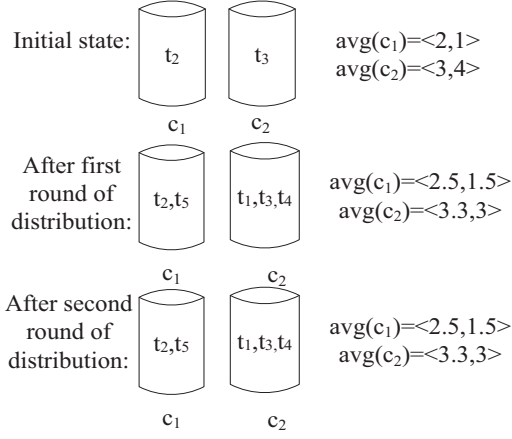


Figure 2: Illustration of MOClustering\_means clustering process

1 tive test cases) of  $K$  clusters, and whenever the clusters are  
 2 changed, the algorithm iteratively uses non-representative test  
 3 cases (non-center points) to replace the representative test case,  
 4 if necessary. The representative test case  $O$  is defined as follow  
 5 [30, 31].

6 **Definition 2.** A Representative Test Case,  $O$ , of a cluster is the  
 7 test case that has the minimum absolute error value in the cluster.  
 8

9 The absolute error value ( $E$ ) of the representative test case  
 10  $O$  is calculated by either Formula 5 or Formula 9 according to  
 11 which distance metric is being used. In MOClustering\_medoids,  
 12 test cases are clustered according to their *Object Method Vec-*  
 13 *tors*. Although K-means is efficient, it is also sensitive to outliers.  
 14 Thus, when a test case with extreme values appears, the  
 15 data distribution may be significantly distorted. The K-medoids  
 16 algorithm can reduce the sensitivity to outliers by selecting a  
 17 test case to represent the cluster without using the mean value.  
 18 Thus, K-medoids was used in the MOClustering method to  
 19 compare with MOClustering\_means. The algorithm first selects  
 20  $K$  test cases to set up the initial  $K$  clusters. Each remaining test  
 21 case is then allocated to the closest cluster, defined by the lowest  
 22 distance to the representative test case of the cluster. The  
 23 representative test case of each cluster is then updated. This  
 24 process is repeated until test cases in each cluster no longer  
 25 change. After clustering, the test cases in one cluster are close  
 26 to the representative test case of that cluster, and far away from  
 27 other clusters.

28 MOClustering\_medoids is shown in Algorithm 2, and has  
 29 three input parameters: *testcasepool* (the simulated input domain),  
 30  $TNum$  (the number of test cases to be selected from the  
 31 simulated input domain to form a test suite on which prioritization  
 32 is to be conducted) and  $K$  (the number of clusters to be  
 33 generated). That is, the algorithm will generate  $K$  clusters for  
 34  $TNum$  test cases selected from *testcasepool*. In MOClustering\_  
 35 medoids,  $TNum$  test cases are first randomly selected from  
 36 the simulated input domain as the initial data; and the number  
 37 of objects and methods is extracted from each chosen test case  
 38 to construct a data set  $OMV$  for these  $TNum$  test cases, i.e.,  
 39 we construct the corresponding relationship between  $OMV$  and

40  $TNum$  test cases, and the test cases are grouped based on the  
 41 corresponding clustering operation on the elements of  $OMV$ .  
 42 Next, the first  $K$  test cases corresponding to the first  $K$  elements  
 43 from  $OMV$  are selected as the initial representative test cases  
 44 for the  $K$  clusters and the selected representative test cases are  
 45 stored in *RepreTCASE*. Then, the Euclidean distance between  
 46 each element of  $OMV$  and the *omv* of the representative test  
 47 case  $O$  (of each cluster) is calculated, and the corresponding  
 48 test case of each object method vector is assigned to the closest  
 49 cluster. Finally, for every cluster, we consider each of its  
 50 non-representative test cases, denoted by  $O'$ , and calculate the  
 51 absolute error value  $E'$  of  $O'$  using Formula 5 – if  $E'$  is less than  
 52  $E$  which is the absolute value of  $O$ , then  $O$  is replaced with  $O'$ .  
 53 This is repeated until all clusters become steady, that is, there  
 54 are no changes in any clusters after an updating process. By  
 55 then,  $K$  clusters would have been generated and stored in the  
 56 data set *clustering*

---

#### Algorithm 2 MOClustering\_medoids (*testcasepool*, $K$ , $TNum$ )

---

```

1: Construct OriginalTCASE = {} to store the selected test cases;
2: Construct OMV = {} to store the set of object method vectors;
3: Construct Clustering = {} to store the generated clusters;
4: Construct RepreTCASE = {} to store representative test cases of each cluster;
5: Choose  $TNum$  test cases from testcasepool randomly and add them to OriginalTCASE;
6: for ( $i=1$  to  $TNum$ )
7:    $On = |OriginalTCASE[i].Objects|$ ; //  $|OriginalTCASE[i].Objects|$  is equal to the number of objects of OriginalTCASE[ $i$ ].
8:    $Mn = |OriginalTCASE[i].Methods|$ ; //  $|OriginalTCASE[i].Methods|$  is equal to the number of methods of OriginalTCASE[ $i$ ].
9:    $OMV[i] = \langle On, Mn \rangle$ ; // The element of OMV is denoted by  $OMV[i]$ .
10: end for
11: Choose  $K$  items from OMV and add the corresponding test cases to RepreTCASE as the initial representative test case;
12: Set change = true;
13: while (change == true)
14:   for ( $i=1$  to  $TNum$ )
15:     for ( $j=1$  to  $K$ )
16:       Calculate  $d(OMV[i], RepreTCASE[j])$  according to Formula 2 ;
17:     end for
18:     Put the corresponding test case of  $OMV[i]$  to the nearest cluster;
19:     Update the cluster that  $OMV[i]$  corresponds to in Clustering;
20:   end for
21:   for ( $i=1$  to  $K$ )
22:     for (each non-representative test case  $O'$  in the cluster )
23:       Compute its absolute error value  $E'$ ; // Formula 5
24:       if ( $E' < E$ )
25:          $RepreTCASE[i] = O'$ ;
26:       end if
27:     end for
28:   end for
29:   if (each RepreTCASE[ $i$ ] keep invariant )
30:     Set change = false;
31:   else
32:     Set change = true;
33:   end if
34: end while
35: return Clustering

```

---

Suppose that  $OMV(c)$  is the set of object method vectors corresponding to  $c$ , and  $OMV(c) = \{omv_1, omv_2, \dots, omv_n\}$ . Let  $E$  denote the absolute error value of a test case  $O$  in cluster  $c$ , and  $omv(O)$  be the element of  $OMV(c)$  corresponding to  $O$ . In MOClustering\_medoids, the absolute error value of the test

1 case  $O$  is defined as:

$$E = \sum_{i=1}^n d(omv_i, omv(O)) \quad (5)$$

2 For example, suppose that the constructed test suite has  
 3 five test cases (that is,  $TNum$  is 5) and their respective  $OMV$ :  
 4  $omv_1 = \langle 4, 3 \rangle$ ,  $omv_2 = \langle 2, 1 \rangle$ ,  $omv_3 = \langle 3, 4 \rangle$ ,  $omv_4 = \langle$   
 5  $1, 5 \rangle$  and  $omv_5 = \langle 3, 2 \rangle$ . Also assume that  $K$  is set to 2,  
 6 i.e., there are two clusters,  $c_1$  and  $c_2$ . The calculation process  
 7 of the earlier stage is the same as in Algorithm 2. Suppose we  
 8 somehow choose two test cases as the initial representative test  
 9 cases:  $t_2$  for  $c_1$  and  $t_3$  for  $c_2$ . We first calculate the distance  
 10 between  $omv_i$  ( $i = 1, 4, 5$ ) and  $omv_2$ , and the distance between  
 11  $omv_i$  ( $i = 1, 4, 5$ ) and  $omv_3$ , then put each test case into the  $n$ -  
 12 nearest cluster. For example, as the distance between  $omv_1$  and  
 13  $omv_2$  is 2.83, and the distance between  $omv_1$  and  $omv_3$  is 1.41,  
 14 then  $omv_1$  should be put into cluster  $c_2$ . After the end of the  
 15 first round distribution based on the similar operations, cluster  
 16  $c_1$  has two test cases ( $t_2$  and  $t_5$ ), and cluster  $c_2$  has three test  
 17 cases ( $t_1$ ,  $t_3$  and  $t_4$ ). Then we need to see whether the represen-  
 18 tative test case of each cluster needs to be updated or not. For  
 19 example, in  $c_2$ , consider  $t_1$ . Calculate its absolute error value  
 20  $E_1$  according to Formula 5. If  $E_1$  is smaller than  $E_3$  (which is  
 21  $t_3$ 's  $E$ ), then  $t_1$  replaces  $t_3$  to become the new representative test  
 22 case. Other test cases in  $c_2$  are also examined. If no change  
 23 is observed for representative test cases of any cluster, then the  
 24 clustering process is completed. Figure 3 summarizes the three  
 25 rounds of distribution for the above example.

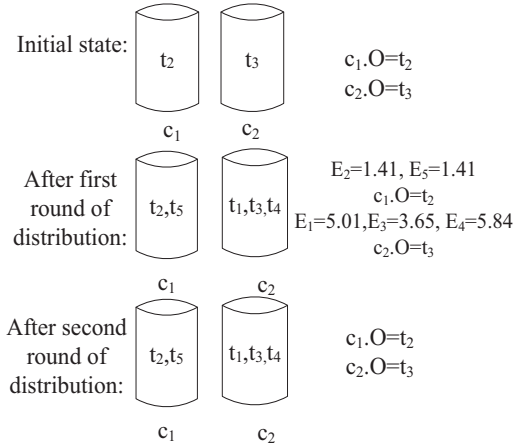


Figure 3: Illustration of MOClustering\_medoids clustering process

### 26 3.3. DMClustering

#### 27 3.3.1. OMISS metric

28 The OOS test input structure may be very complex because  
 29 it may include different combinations of objects and methods,  
 30 including multiple classes, multiple objects, inherited elements,  
 31 reference objects, self-defined methods, and method invocation  
 32 sequences. To investigate the impact of using different distance  
 33 metrics on test case prioritization, we use our recently devel-  
 34 oped OMISS metric to calculate the distance between test cases  
 35 in the clustering process.

36 According to the OMISS metric [23], a test input  $t$  con-  
 37 sists of an object set ( $OBJ$ ) and a method invocation sequence  
 38 ( $MINV$ ), i.e.,  $t = \{t.OBJ, t.MINV\}$ . The distance between test  
 39 inputs ( $TestcaseDistance$ ) is defined as the sum of the distance  
 40 of object sets ( $TCobjectDistance$ ) and the distance of method  
 41 invocation sequences ( $TCmSeqDist$ ), as shown in Formula 6. In  
 42 Formula 6,  $t_1.OBJ$  and  $t_2.OBJ$  refer to the objects sets in  $test-$   
 43  $case1$  and  $testcase2$ , respectively; and  $t_1.MINV$  and  $t_2.MINV$   
 44 represent the method invocation sets of  $testcase1$  and  $testcase2$ ,  
 45 respectively.

$$TestcaseDistance(t_1, t_2) = TCobjectDistance(t_1.OBJ, t_2.OBJ) + TCmSeqDist(t_1.MINV, t_2.MINV) \quad (6)$$

The distance between two object sets ( $TCobjectDistance$ )  
 is calculated by comparing each pair of objects in the two sets,  
 and is defined as the minimum sum of distances amongst all  
 possible objects pairing between  $t_1.OBJ$  and  $t_2.OBJ$ . An object  
 can be divided into two parts: the attribute section and behavior  
 section. The attribute section includes self-defined attributes  
 (the attributes are defined by the current class), inherited at-  
 tributes, and reference attributes. The behavior section includes  
 self-defined methods and inherited methods. Hence, the distance  
 between objects ( $ObjectDistance$ ) is determined by the  
 attribute section ( $AttributeDistance$ ) and the behavior section  
 ( $BehaviorDistance$ ) of the object. The distance between ob-  
 jects is defined in Formula 7, where  $p.A$  refers to the attribute  
 section of object  $p$ ,  $q.A$  refers to the attribute section of object  
 $q$ ,  $p.B$  means the behavior section of object  $p$ , and  $q.B$  means  
 the behavior section of object  $q$ .

$$ObjectDistance(p, q) = AttributeDistance(p.A, q.A) + BehaviorDistance(p.B, q.B) \quad (7)$$

36 The distance between the two method invocation sequences,  
 37 which is defined in Formula 8, includes the length difference,  
 38 the set difference and the sequence difference. The sequence  
 39 difference is calculated by  $SequenceDissimilarity(t_1.MINV,$   
 40  $t_2.MINV)$  in Formula 8 based on the ordered lists, and is equal  
 41 to the number of common methods in the same position divided  
 42 by the number of methods in the shorter sequence. For exam-  
 43 ple, if there are two method invocation sequences,  $t_1.MINV =$   
 44  $\{m_3, m_2, m_1\}$ , which has three methods, and  $t_2.MINV = \{m_4, m_2,$   
 45  $m_1, m_3, m_5\}$ , which has five methods, then the length differ-  
 46 ence is 2; the set difference is 0.4 ( $1-3/5$ ), because  $t_1.MINV$  and  
 47  $t_2.MINV$  have three common methods ( $m_1, m_2$  and  $m_3$ ) and  
 48 five different methods ( $m_1, m_2, m_3, m_4,$  and  $m_5$ ); the sequence  
 49 difference is 0.667 ( $=2/3$ ), because the second and third meth-  
 50 ods of  $t_1.MINV$  are equal to the second and third methods of  
 51  $t_2.MINV$ ; and  $t_1.MINV$  is the shorter sequence, with a total of  
 52 three methods. Therefore the distance between  $t_1.MINV$  and  
 53  $t_2.MINV$  is 3.067 ( $2+0.4+0.667$ ).

$$\begin{aligned}
& TCMSeqDist(t_1.MINV, t_2.MINV) \\
&= |length(t_1.MINV) - length(t_2.MINV)| \\
&\quad + \left(1 - \frac{|t_1.MINV \cap t_2.MINV|}{|t_1.MINV \cup t_2.MINV|}\right) \\
&+ SequenceDissimilarity(t_1.MINV, t_2.MINV)
\end{aligned} \tag{8}$$

The detailed explanations and examples with regard to Formulas 6, 7, and 8 can be found in [23].

### 3.3.2. DMClustering algorithm

Because DMClustering applies to objects and methods, which are not numerical data, the K-means algorithm could not be used. Hence, only the K-medoids clustering algorithm was used in DMClustering.

---

#### Algorithm 3 DMClustering (*testcasepool*, *K*, *TNum*)

---

```

1: Construct OriginalTCASE = {} to store the selected test cases;
2: Construct Clustering = {} to store the generated clusters;
3: Construct RepreTCASE = {} to store representative test cases of each cluster;
4: Choose TNum test cases from testcasepool randomly and add them to OriginalTCASE;
5: Choose K items from OriginalTCASE and add them to RepreTCASE as the
   initial representative test case;
6: Set change = true;
7: while (change == true)
8:   for (i=1 to TNum)
9:     for (j=1 to K)
10:      Calculate TestcaseDistance(OriginalTCASE[i], RepreTCASE[j]);
11:      // Formula 6
12:     end for
13:     Put OriginalTCASE[i] to the nearest cluster;
14:     Update the cluster of OriginalTCASE[i] in Clustering;
15:   end for
16:   for (i=1 to K)
17:     for (each non-representative test case O' in the cluster )
18:       Compute its absolute error value E'; // Formula 9
19:       if (E' < E)
20:         RepreTCASE[i] = O';
21:       end if
22:     end for
23:   end for
24:   if (each RepreTCASE[i] keep invariant )
25:     Set change = false;
26:   else
27:     Set change = true;
28:   end if
29: end while
30: return Clustering

```

---

DMClustering is shown in Algorithm 3, and has three input parameters: *testcasepool* (the simulated input domain), *TNum* (the number of test cases to be selected from *testcasepool* to form a test suite) and *K* (the number of clusters to be generated). The algorithm generates *K* clusters for *TNum* test cases selected from *testcasepool*. In DMClustering, *TNum* test cases are first randomly selected from the *testcasepool* as the initial data, and are then added to *OriginalTCASE*. Next, *K* items from *OriginalTCASE* are selected as the initial representative test case *O* of each cluster, and the generated representative test case is stored in *RepreTCASE*. Then, the difference between each remaining test case and each representative test case *O* (of each

cluster) are calculated with the OMISS metric (Formulas 6, 7, and 8), and each test case is assigned to the nearest cluster. Finally, for each non-representative test case *O'*, its absolute error value (*E'*) is calculated using Formula 9 – if *E'* is less than *E* (the absolute value of *O*), then *O* is replaced with *O'* and the clusters are updated. This is repeated until items in each cluster no longer change, at which point, *K* clusters will have been generated and stored in *clustering*.

Suppose that *T* is the set of test cases for cluster *c*,  $T = \{t_1, t_2, \dots, t_n\}$ . Let *E* denote the absolute error value of the representative test case *O* in cluster *c*. In DMClustering, the absolute error value of the test case *O* is defined as:

$$E = \sum_{i=1}^n OMIS S(t_i, O) \tag{9}$$

For example, assume that a cluster *c*<sub>1</sub> has three test cases, *t*<sub>1</sub>, *t*<sub>2</sub>, and *t*<sub>3</sub>. First, calculate the sum of the distances *OMIS S*(*t*<sub>2</sub>, *t*<sub>1</sub>) and *OMIS S*(*t*<sub>3</sub>, *t*<sub>1</sub>), denoted *E*<sub>1</sub> (the *E* for *t*<sub>1</sub>). Then, calculate the sum of *OMIS S*(*t*<sub>1</sub>, *t*<sub>2</sub>) and *OMIS S*(*t*<sub>3</sub>, *t*<sub>2</sub>), denoted *E*<sub>2</sub> (the *E* for *t*<sub>2</sub>), and the sum of *OMIS S*(*t*<sub>1</sub>, *t*<sub>3</sub>) and *OMIS S*(*t*<sub>2</sub>, *t*<sub>3</sub>), denoted *E*<sub>3</sub> (the *E* for *t*<sub>3</sub>). If *E*<sub>3</sub> is smaller than *E*<sub>1</sub> and *E*<sub>2</sub>, then *t*<sub>3</sub> is the representative test case of cluster *c*<sub>1</sub>.

## 4. Adaptive Random Sequence Generation

After all test cases have been clustered, a sampling strategy is needed to choose test cases from the clusters. The traditional random sampling strategy selects *n* test cases randomly from the entire pool of test cases. Some of these *n* test cases may be from the same cluster, which may have similar properties, including the ability to detect the same fault. Such a test case sequence may lead to a poor fault detection rate. The same problem occurs if random sampling is applied to choose a cluster, from which a test case is then selected. To maintain the diversity in test cases, we use a new MSampling (maximum) sampling mechanism.

MSampling is explained in Algorithm 4. It has three input parameters: *K* (the number of generated clusters), *n* (the specified number of test cases to be prioritized), and *clustering* (the *K* clusters generated by MOClustering and DMClustering), where *n* is less than or equal to the number of test cases in all *K* clusters. The specific steps of MSampling are: (1) Randomly choose an initial cluster. (2) When these clusters are generated by MOClustering, the distances between the selected cluster and the unselected clusters are calculated using Formulas 10 and 11. When the clusters are generated using DMClustering, the distance is calculated with Formula 12. (3) The most distant cluster is selected next. (4) Steps 2 and 3 are repeated until all clusters are selected, and an ordered sequence of clusters is generated. (5) According to the order of clusters in the sequence, randomly choose a unique test case from each cluster, in sequence. (6) Repeat Step 5 until the specified number (*n*) of prioritized test cases has been obtained. If the number of test cases selected in the current cluster is equal to the length of this cluster, we should jump to the next cluster. The prioritized test case sequence is stored in the data set *GTCases*.

**Algorithm 4** MSampling( $K, n, clustering$ )

---

```

1: Construct a set to store  $K$  clusters  $OC = \{c_1, c_2, \dots, c_i, \dots, c_K\}$ ;
2: Construct  $C = ()$  to store the chosen cluster;
3: Construct  $GTCases = ()$  to store the prioritized test case sequence;
4: Randomly choose a cluster  $c$ ;
5: Add  $c$  to  $C$ ;
6: while !(all clusters are added to  $C$ )
7:   for ( $i = 1$  to  $K$ )
8:     if ( $OC[i]$  is not added to  $C$ )
9:       Calculate the distance between  $C$  and  $OC[i]$ ;
10:    end if
11:  end for
12:  Update  $c =$  the cluster that has the farthest distance with  $C$ ;
13:  Add  $c$  to  $C$ ;
14: end while
15: while !(the number of test cases in  $GTCases$  is up to  $n$ )
16:  for (each  $c$  in  $C$  (in their order in  $C$  and assume  $C$  is circular))
17:    if (the number of test cases selected in  $c <$  the length of  $c$ )
18:      Take a test case  $t$  from each cluster in turn;
19:      Append  $t$  to  $GTCases$ ;
20:    else
21:      Jump to the next cluster;
22:    end if
23:  end for
24: end while
25: return  $GTCases$ ;

```

---

When these clusters are generated by MOClustering\_means, then, if  $C$  is a cluster set, and  $C = \{c_1, c_2, \dots, c_K\}$ ,  $AVG$  is the mean value set, and  $AVG = \{avg_1, avg_2, \dots, avg_K\}$ , where  $avg_i$  is the mean value of  $c_i$ . Let  $DMO\_M(c_i, c_j)$  be the distance between clusters  $c_i$  and  $c_j$  ( $i, j = 1, 2, \dots, K$ ). The distance between any two clusters is defined as:

$$DMO\_M(c_i, c_j) = d(avg_i, avg_j) \quad (10)$$

Similarly, when the clusters are generated by MOClustering\_medoids, if  $C$  is a cluster set, and  $C = \{c_1, c_2, \dots, c_K\}$ ,  $OS$  is a representative test cases set, and  $OS = \{o_1, o_2, \dots, o_K\}$ , with  $o_i$  being the representative test case of the corresponding  $c_i$ , ( $i = 1, 2, \dots, K$ ). Let  $DMO\_K(c_i, c_j)$  be the distance between clusters  $c_i$  and  $c_j$ , ( $i, j = 1, 2, \dots, K$ ). The distance between clusters is defined as:

$$DMO\_K(c_i, c_j) = d(omv(o_i), omv(o_j)) \quad (11)$$

With DMClustering, if  $C$  is a cluster set, and  $C = \{c_1, c_2, \dots, c_K\}$ ,  $OS$  is a representative test cases set, and  $OS = \{o_1, o_2, \dots, o_K\}$ , with  $o_i$  being the representative test case of the corresponding  $c_i$ . Let  $DDM(c_i, c_j)$  stand for the distance between clusters  $c_i$  and  $c_j$ , which is defined as:

$$DDM(c_i, c_j) = OMIS(S(o_i, o_j)) \quad (12)$$

For example, if we have three clusters,  $c_1 = \{t_{11}, t_{12}\}$ ,  $c_2 = \{t_{21}, t_{22}\}$ , and  $c_3 = \{t_{31}, t_{32}\}$ , then suppose  $c_2$  is chosen as the first cluster, and the distances between  $c_1$  and  $c_2$  and between  $c_3$  and  $c_2$  are calculated. If the distance between  $c_1$  and  $c_2$  is less than that between  $c_3$  and  $c_2$ , then the order of clusters is  $c_2, c_3$  and  $c_1$ . Based on Algorithm 4, test cases are selected from  $c_2, c_3$  and  $c_1$  in sequence. Suppose  $GTCases = (t_{21}, t_{31}, t_{11})$  after the first round of selection, then the next round of selection

is conducted. At the end of the sampling strategy, a final adaptive random sequence for the prioritized test cases is generated:  $GTCases = (t_{21}, t_{31}, t_{11}, t_{22}, t_{32}, t_{12})$ . Test cases in the sequence are expected to be evenly spread in the input domain and are executed in this order in the testing framework.

## 5. Empirical Studies And Analysis

### 5.1. Setup of the empirical studies

Mutant programs are often used in empirical studies to investigate the fault detection effectiveness of different testing methods. Given the same test inputs, if the outputs produced by a mutant version are different from the outputs produced by the original program, then these test inputs can be regarded as failure-causing inputs. Our study also used mutation programs to evaluate how quickly a test case prioritization method could find failure-causing inputs.

Table 1 presents the seven subject programs investigated in the experiment. The programs were all written in the C++ or C# language and are from some open sources websites [32–35]. Faults were manually seeded into the subject program methods based on common mutation operators. In this study, we used the following 13 operators [36], which generate some typical program faults.

- 1) arithmetic operators replacement (AOR);
- 2) logical operators replacement (LOR);
- 3) relational operators replacement (ROR);
- 4) constant for scalar variable replacement (CSR);
- 5) scalar variable for scalar variable replacement (SVR);
- 6) scalar variable for constant replacement (SCR);
- 7) array reference for constant replacement (ACR);
- 8) new method invocation with child class type (NMI);
- 9) argument order change (AOC);
- 10) accessor method change (AMeC);
- 11) access modifier change (AMoC);
- 12) hiding variable deletion (HVD);
- 13) property replacement with member field (PRM).

Of these 13 mutation operators, the last six are OO-specific, and are used to generate OO-specific faults. Table 2 shows the type of mutation operators and the number of faults seeded for each program. The machine used to conduct the testing has an Intel dual core i3-2120 3.3 GHz processor, 4 GB of RAM, and runs under the Windows 7 operating system.

### 5.2. Effectiveness measure criteria

In our study, we used three measures to compare the TCP approaches:  $F_m$  (F-measure) – the number of the test cases executed before finding the first fault;  $E$  – the total number of distinct faults detected by a specific number of test cases; and  $APFD$  – the weighted average percentage of faults detected. A testing approach is considered effective if it has a low F-measure, a high  $E$ , and a high  $APFD$  value [37]. In this study, we compared MOClustering\_means, MOClustering\_medoids, DM-Clustering, and RT-ms (RT with method sequence — a random sequence generation approach for OOS test cases with method



Table 1: SUBJECT PROGRAMS

| ID | name                 | Lines of code | Num. of public classes | Num. of public methods | Num. of faults | Description  |
|----|----------------------|---------------|------------------------|------------------------|----------------|--|
| 1  | CCoinBox [32]        | 120           | 1                      | 7                      | 4              | C++ library that simulates a vending machine                   |
| 2  | WindShieldWiper [32] | 233           | 1                      | 13                     | 4              | C++ library that simulates a windshield wiper                  |
| 3  | SATM [32]            | 197           | 1                      | 9                      | 4              | C++ library that simulates an Automatic Teller Machine         |
| 4  | RabbitsAndFoxes [33] | 770           | 6                      | 33                     | 9              | C# program that simulates a predator-prey model                |
| 5  | WaveletLibrary [34]  | 2406          | 12                     | 84                     | 15             | C# library for wavelet algorithms                              |
| 6  | IceChat [34]         | 571000        | 101                    | 271                    | 24             | C# program that implements an IRC (Internet Relay Chat) Client |
| 7  | CSPspEmu [35]        | 406808        | 443                    | 1433                   | 26             | C# program for a PSP (PlayStation Portable) emulator           |

Table 2: MUTATION OPERATORS AND THE NUMBER OF FAULTS SEEDED

| ID | Num. of faults | Mutation operators (number)   |
|----|----------------|---|
| 1  | 4              | AOR(1), LOR(2), ROR(1)  |
| 2  | 4              | AOR(1), LOR(1), ROR(1), ACR(1)  |
| 3  | 4              | AOR(1), LOR(1), ROR(1), SCR(1)  |
| 4  | 9              | AOR(1), LOR(1), SVR(1), NMI(1), AOC(1), AMeC(1), AMoC(1), HVD(1), PRM(1), LOR(1), SVR(1), CSR(2), SCR(1), ACR(2), |
| 5  | 15             | NMI(1), AOC(1), AMeC(1), AMoC(2), HVD(2), PRM(1)  |
| 6  | 24             | AOR(2), LOR(1), ROR(1), SVR(2), CSR(2), SCR(1), ACR(2), NMI(2), AOC(3), AMeC(2), AMoC(2), HVD(2), PRM(2)          |
| 7  | 26             | AOR(2), LOR(1), ROR(1), SVR(2), CSR(1), SCR(1), ACR(2), NMI(2), AOC(3), AMeC(3), AMoC(3), HVD(3), PRM(2)          |

performance.

In this study, in order to find its most suitable value,  $K$  was set to 2%, 5%, 10%, 15%, 20%, 25%, and 30% of the total number of test cases (5000 test cases). Based on the overall experimental results, appropriate values of  $K$  for each subject program were determined, as shown in Table 3.

Table 3: THE VALUE OF  $K$  FOR EACH SUBJECT PROGRAMS

| ID | Name            | $K$ | Percentage of the total number of test cases |
|----|-----------------|-----|--|
| 1  | CCoinBox        | 500 | 10%  |
| 2  | WindShieldWiper | 500 | 10%  |
| 3  | SATM            | 500 | 10%  |
| 4  | RabbitsAndFoxes | 750 | 15%  |
| 5  | WaveletLibrary  | 750 | 15%  |
| 6  | IceChat         | 750 | 15%  |
| 7  | CSPspEmu        | 750 | 15%  |

invocation sequence), and Method\_Coverage (a method coverage TCP technique).

In order to properly assess the statistical significance of the differences between our methods and other methods, we conducted the effective statistical analysis based on the p-values and effect size (set at a 5% level of significance) using the unpaired two-tailed Wilcoxon-Mann-Whitney test and the non-parametric Vargha and Delaney effect size measure [38–40].

The p-value is used to show the statistical significance of difference. If the p-value (probability value) is less than 0.05, which means that there is significant difference between the two compared methods, otherwise not [38]. Additionally, we used the non-parametric effect size (ES) measure to show the probability that one method is better than another [39]. That is, when we get the ES for any two methods A and B, a higher ES value indicates higher probability showing A is better than B. In this study, we used R language [41] to obtain the p-value and ES value for the pair-wise TCP techniques.

### 5.3. Experimental parameters

For both the  $K$ -means and the  $K$ -medoids clustering algorithms,  $K$  is the main input parameter. If the value of  $K$  is not suitable, low quality clusters may be generated: if test cases are clustered into too many clusters, then some similar test cases may be put into different clusters; if they are clustered into too few, then dissimilar test cases may be put into the same cluster. Both of these situations may lead to poor failure detection

In addition, in all experiments ( $F_m$ ,  $E$  and  $APFD$ ),  $testcase\_pool$  in Algorithms 1 to 3 simulated the input domain, and  $T\_Num$  in Algorithms 1 to 3 was the total number of test cases (5000). The value of  $n$  in Algorithm 4 was set to 100, 500, 1000, 1500, 2000, 2500, 3000, 3500, 4000 and 5000.

### 5.4. Experiments

To evaluate the effectiveness of our approaches, we attempted to answer the following three research questions:

**RQ1:** Do cluster TCP techniques perform better than prioritization with random sequences or method coverage, in terms of  $F_m$ ?

**RQ2:** Do cluster TCP techniques perform better than prioritization with random sequences or method coverage, in terms of  $E$ ?

**RQ3:** Do cluster TCP techniques perform better than prioritization with random sequences or method coverage, in terms of  $APFD$ ?

#### 5.4.1. Results and discussion

1) Do MOClustering means, MOClustering medoids and DM-Clustering perform better than prioritization with random sequences and Method\_Coverage, in terms of  $F_m$ ?

Table 4 summarizes the  $F_m$  results for the five different methods. All results in the table were averaged over 100 runs of tests for each subject program, each time with a different seed.

Table 4:  $F_M$  OF VARIOUS TCP METHODS

| ID   | $F_m$                 |                         |                  |                     |       |
|------|-----------------------|-------------------------|------------------|---------------------|-------|
|      | MOClustering<br>means | MOClustering<br>medoids | DMClus<br>tering | Method<br>_Coverage | RT-ms |
| 1    | 53.94                 | 70.37                   | 72.15            | 71.92               | 74.87 |
| 2    | 63.88                 | 58.54                   | 58.85            | 68.84               | 75.54 |
| 3    | 49.93                 | 46.67                   | 47.17            | 49.05               | 52.90 |
| 4    | 21.88                 | 27.08                   | 22.91            | 24.02               | 28.45 |
| 5    | 6.96                  | 8.58                    | 6.75             | 8.44                | 8.66  |
| 6    | 37.58                 | 36.54                   | 33.14            | 40.19               | 55.00 |
| 7    | 85.14                 | 77.14                   | 71.98            | 83.53               | 89.67 |
| mean | 45.62                 | 46.42                   | 44.71            | 49.43               | 55.01 |
| sDev | 48.30                 | 51.47                   | 49.82            | 41.29               | 58.20 |

1 Table 4 shows that, for the CcoinBox program, MOClustering  
2 \_means used the least number of test cases to detect the first  
3 failure, followed by MOClustering\_medoids, Method\_Coverage,  
4 DMClustering and RT-ms. For programs WindShieldWiper,  
5 MOClustering\_medoids found the first fault with the least num-  
6 ber of test cases, followed by DMClustering, MOClustering\_me-  
7 ans, Method\_Coverage and RT-ms. For programs SATM, MO-  
8 Clustering\_medoids found the first fault with the least number  
9 of test cases, followed by DMClustering, Method\_Coverage,  
10 MOClustering\_means and RT-ms. For the RabbitsAndFoxes  
11 program, the number of test cases used by MOClustering\_means  
12 and DMClustering to detect the first failure was similar, and less  
13 than that for Method\_Coverage, MOClustering\_medoids and RT-  
14 ms. For the WaveletLibrary, IceChat, and CSPspEmu program-  
15 s, DMClustering used the least number of test cases to find  
16 the first failure, and RT-ms used the most. For the program-  
17 s IceChat and CSPspEmu, MOClustering\_medoids performed  
18 better than MOClustering\_means and RT-ms, but for the pro-  
19 gram WaveletLibrary, MOClustering\_means performed better  
20 than Method\_Coverage, MOClustering\_medoids and RT-ms. Th-  
21 erefore, in terms of the  $F_m$ , on average, DMClustering per-  
22 formed best, especially for the large-scale programs, followed  
23 by MOClustering\_means, MOClustering\_medoids, Method\_Cov-  
24 erage and RT-ms. Compared with RT-ms, DMClustering achie-  
25 ved an average of 18.72% improvement; MOClustering\_means  
26 achieved an average of 17.07%; and MOClustering\_medoids  
27 achieved an average of 15.62% improvement. Compared with  
28 Method\_Coverage, DMClustering achieved an average of 9.55%  
29 improvement; MOClustering\_means achieved an average of 7.71%  
30 and MOClustering\_medoids achieved an average of 6.09% im-  
31 provement. Hence, the proposed cluster TCP techniques always  
32 performed better than prioritization with random sequences and  
33 method coverage prioritization, in terms of  $F_m$ .

34 In order to further analyze the  $F_m$  of each testing method for  
35 each subject program, Tables 4 and 5 also summarize the main  
36 statistical measures including  $sDev$  (standard deviation) for the  
37 7 subject programs. The standard deviation for RT-ms is the  
38 biggest (58.20), which indicates that its data points are spread  
39 out over a wider range than other TCP techniques.

40 Figure 4 to 10 are box-plots diagrams showing the  $F_m$  re-  
41 sults for the seven subject programs, with the data in each box-  
42 plot being the  $F_m$  results over 100 runs for each subject program  
43 with different seeds.

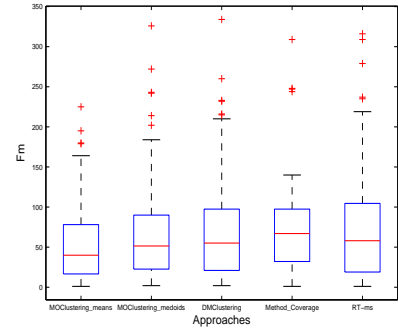
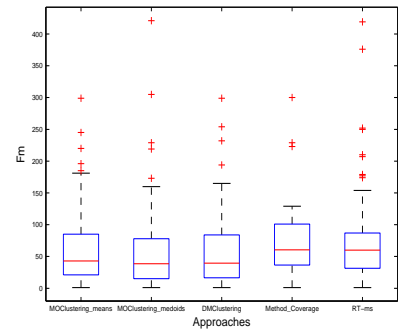
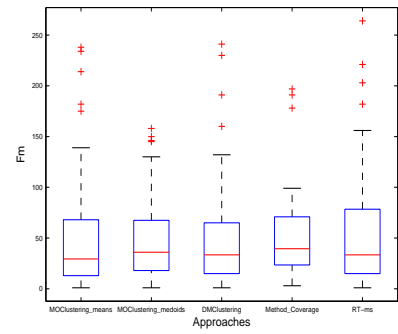
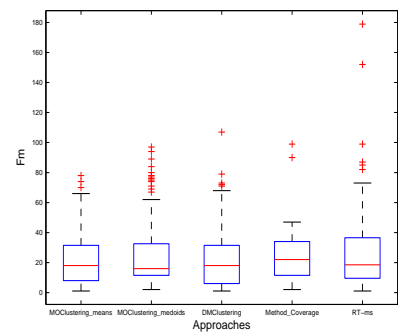
Figure 4:  $F_m$  experimental results for CcoinBoxFigure 5:  $F_m$  experimental results for WindShieldWiperFigure 6:  $F_m$  experimental results for SATMFigure 7:  $F_m$  experimental results for RabbitsAndFoxes

Table 5: STATISTICAL RESULT OF  $F_M$  FOR 7 SUBJECT PROGRAMS

| ID |      | MOClustering<br>_means | MOClustering<br>_medoids | DMClustering | Method<br>_Coverage | RT-ms |
|----|------|------------------------|--------------------------|--------------|---------------------|-------|
| 1  | mean | 53.94                  | 70.37                    | 72.15        | 71.92               | 74.87 |
|    | sDev | 48.93                  | 64.43                    | 64.39        | 53.93               | 70.01 |
| 2  | mean | 63.88                  | 58.54                    | 58.85        | 68.84               | 75.54 |
|    | sDev | 61.27                  | 66.24                    | 57.75        | 48.59               | 71.40 |
| 3  | mean | 49.93                  | 46.67                    | 47.17        | 49.05               | 52.90 |
|    | sDev | 43.13                  | 38.56                    | 40.43        | 36.79               | 51.90 |
| 4  | mean | 21.88                  | 27.08                    | 22.91        | 24.02               | 28.45 |
|    | sDev | 17.37                  | 24.50                    | 20.80        | 16.42               | 29.32 |
| 5  | mean | 6.96                   | 8.58                     | 6.75         | 8.44                | 8.66  |
|    | sDev | 6.22                   | 6.88                     | 5.39         | 5.50                | 7.67  |
| 6  | mean | 37.58                  | 36.54                    | 33.14        | 40.19               | 55.00 |
|    | sDev | 34.09                  | 35.88                    | 34.57        | 34.13               | 39.36 |
| 7  | mean | 85.14                  | 77.14                    | 71.98        | 83.53               | 89.67 |
|    | sDev | 53.06                  | 56.03                    | 54.40        | 60.42               | 61.91 |

Table 6: COMPARISON BETWEEN VARIOUS PAIRS OF METHODS USING P-VALUE AND EFFECTIVE SIZE METHODS ON  $F_M$

| Pair of methods      | MOClustering<br>_means and<br>RT-ms | MOClustering<br>_medoids and<br>RT-ms | DMClustering<br>and RT-ms | MOClustering<br>_means and<br>Method<br>_Coverage | MOClustering<br>_medoids and<br>Method<br>_Coverage | DMClustering<br>and Method<br>_Coverage | MOClustering<br>_means and<br>DMClustering | MOClustering<br>_medoids and<br>DMClustering | MOClustering<br>_means and<br>MOClustering<br>_medoids |
|----------------------|-------------------------------------|---------------------------------------|---------------------------|---|---|---|--|--|--|
| <b>P-value</b>       | 0.007193                            | 0.006183                              | 0.000192                  | 0.000943  | 0.000692  | 1.64E-05                                | 0.261891                                   | 0.360047                                     | 0.844637   |
| <b>ES</b>            | 0.5434847                           | 0.5422582                             | 0.557556                  | 0.551042  | 0.5523612   | 0.566514                                | 0.4826837                                  | 0.485873                                     | 0.503026   |
| <b>Better Method</b> | MOClustering<br>_means              | MOClustering<br>_medoids              | DMClustering              | MOClustering<br>_means                            | MOClustering<br>_medoids                            | DMClustering                            | DMClustering                               | DMClustering                                 | MOClustering<br>_means                                 |

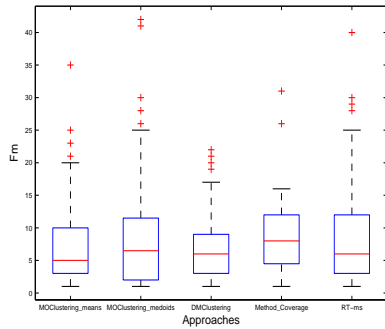


Figure 8:  $F_m$  experimental results for WaveletLibrary

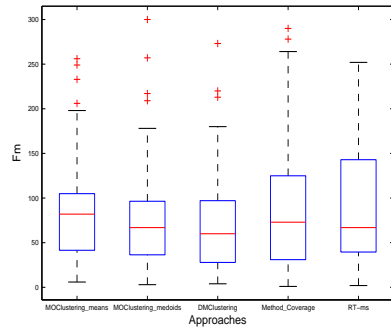


Figure 10:  $F_m$  experimental results for CSPspEmu

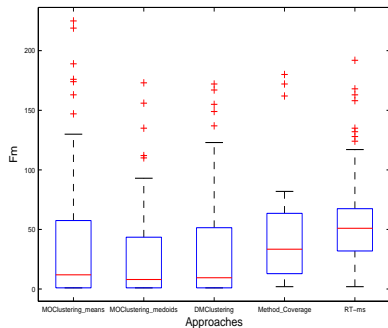


Figure 9:  $F_m$  experimental results for IceChat

1 As can be observed from Figure 4, both the outlying and the  
2 maximum observed values of MOClustering\_means are far s-  
3 maller than corresponding values of RT-ms. Figure 5 shows that  
4 the performances of the five methods are similar, but the medi-  
5 ans of MOClustering\_means, MOClustering\_medoids and DM-  
6 Clustering are much smaller than the medians of Method.Coverage  
7 and RT-ms. This implies that in most cases, Method.Coverage  
8 and RT-ms required more test cases to find the first fault. As  
9 shown by Figure 6, three cluster TCP techniques outperform  
10 other methods with smaller medians. As observed from Fig-  
11 ure 7, the performances of MOClustering\_medoids, MOClus-  
12 tering\_means, DMClustering, and Method.Coverage are simi-  
13 lar, while they outperform RT-ms with larger outlying point  
14 values and maximum values. As seen from Figures 8 and 10, DM-  
15 Clustering has the shorter IQR (interquartile range) and smaller  
16 medians than Method.Coverage and RT-ms, which means that  
17 its performance is more stable than that of these two methods

for the two larger programs. Figure 9 shows that all three cluster TCP techniques have smaller medians than Method\_Coverage and RT-ms, which implies that their performances are better than those of these two methods for the program on average. From Figures 4 to 10, we can find that the cluster TCP techniques have shorter IQRs and smaller medians than the other methods. Hence, they have more stable performance, especially for the larger programs.

In order to further study the significance of the differences in  $F_m$ , we report in Table 6 the p-value and effective size (ES) [38] for pairwise comparisons between the representative techniques from two different groups, from which we can find that the difference of  $F_m$  between our methods and RT-ms is significant (because the p-value is less than 0.05), and the difference in  $F_m$  between DMClustering and Method\_Coverage is also significant. But the difference among our methods is not significant (because the p-value is larger than 0.05). Through a further analysis on the ES values for different pairwise comparisons, we can find that these values between our methods and other methods including RT-ms and Method\_Coverage are larger than 0.5, which indicates that our methods perform better than RT-ms and Method\_Coverage. Column "Better Method" of Table 6 presents the better method of the relevant pair. In three cluster TCP techniques, DMClustering performs best, followed by MOClustering\_means and MOClustering\_medoids on average.

We also analyzed the time taken to detect the first failure ( $F_m$ -time) for the different methods for the seven subject programs. Table 7 shows the  $F_m$ -time results for the five different methods. RT-ms required the least amount of time to detect the first failure. The testing time depends on the specific program under test, and the testing time generally includes both test case generation and execution time, which is usually the main cost in real testing activities. The testing times for MOClustering\_means, MOClustering\_medoids and DMClustering were not more than twice that of RT-ms on average.

Table 7:  $F_m$ -TIME OF VARIOUS TCP METHODS

| ID   | $F_m$ -time (Seconds) |                      |              |                 |       |
|------|-----------------------|----------------------|--------------|-----------------|-------|
|      | MOClustering_means    | MOClustering_medoids | DMClustering | Method_Coverage | RT-ms |
| 1    | 0.71                  | 0.86                 | 1.13         | 0.67            | 0.64  |
| 2    | 1.15                  | 1.37                 | 1.85         | 1.05            | 0.96  |
| 3    | 0.79                  | 0.83                 | 1.17         | 0.71            | 0.68  |
| 4    | 0.83                  | 0.92                 | 1.22         | 0.73            | 0.67  |
| 5    | 0.74                  | 0.82                 | 1.06         | 0.69            | 0.62  |
| 6    | 1.28                  | 1.64                 | 2.13         | 1.14            | 0.97  |
| 7    | 1.84                  | 2.36                 | 3.04         | 1.46            | 1.35  |
| Mean | 1.05                  | 1.31                 | 1.66         | 0.92            | 0.84  |

RT-ms performs better than MOClustering\_means, MOClustering\_medoids, DMClustering and Method\_Coverage in terms of  $F_m$ -time, but it has low effectiveness in terms of  $F_m$ . DMClustering outperforms RT-ms in terms of  $F_m$ . Due to the complex structure of OOS test inputs in the subject programs under test, OMISS requires more time to calculate the distance between test inputs. Hence, DMClustering improves the fault detection effectiveness, but at the expense of more time for computing the OMISS metric.

On the other hand, MOClustering (especially MOClustering\_means)

outperforms RT-ms in terms of  $F_m$ , and outperforms DMClustering in terms of  $F_m$ -time. Since the distance between test inputs in MOClustering is calculated using the Euclidean distance which is much simpler than OMISS metric used in DMClustering, the  $F_m$ -time of MOClustering was less than that of DMClustering on average. Therefore, according to the different testing requirements, we have a trade-off for employing different methods. In other words, when we know the approximate execution time for specific subject programs, we may be able to determine which method should be used based on  $F_m$  performance. For example, if the test case execution time is less than or equal to the test case generation time, then we may consider the influence of the generation time; but, if the generation time is much less than the execution time, then we may ignore its impact.

2) Do MOClustering\_means, MOClustering\_medoids and DMClustering perform better than prioritization with random sequences and Method\_Coverage, in terms of E?

Table 8 shows the total number of distinct faults detected for seven subject programs using ten different test suite sizes – 100, 500, 1000, 1500, 2000, 2500, 3000, 3500, 4000 and 5000. All results were again obtained over 100 runs, with different seeds for each run.

Table 8: THE SUM OF FAULT DETECTED FOR ALL SEVEN SUBJECT PROGRAMS WITH DIFFERENT NUMBERS OF TEST CASES

| Number of Test Cases | E                  |                      |              |                 |       |
|----------------------|--------------------|----------------------|--------------|-----------------|-------|
|                      | MOClustering_means | MOClustering_medoids | DMClustering | Method_Coverage | RT-ms |
| 100                  | 16.10              | 16.10                | 17.01        | 15.82           | 13.79 |
| 500                  | 38.92              | 37.80                | 39.20        | 38.50           | 36.12 |
| 1000                 | 46.55              | 45.71                | 46.97        | 44.10           | 42.77 |
| 1500                 | 50.47              | 49.70                | 50.89        | 48.58           | 46.55 |
| 2000                 | 53.41              | 52.64                | 54.11        | 51.52           | 49.14 |
| 2500                 | 55.86              | 54.74                | 56.42        | 52.92           | 51.52 |
| 3000                 | 57.75              | 56.77                | 58.52        | 56.14           | 53.55 |
| 3500                 | 59.29              | 58.38                | 59.78        | 57.40           | 55.30 |
| 4000                 | 60.76              | 60.06                | 61.25        | 58.80           | 57.19 |
| 5000                 | 62.97              | 62.91                | 63.33        | 62.83           | 62.56 |

Table 8 shows, as expected, that as the number of test cases used increases, the sum of distinct faults detected also increases. Furthermore, DMClustering has the best performance among the testing methods, followed by MOClustering\_means, MOClustering\_medoids, Method\_Coverage and RT-ms.

Figure 11 shows the total number of distinct faults detected by a number ( $n$ ) of test inputs generated by each testing method, across all subject programs. We found that DMClustering outperformed all other methods, followed by MOClustering\_means, MOClustering\_medoids, Method\_Coverage and RT-ms, regardless of the value of  $n$ .

In order to further analyze the difference between different methods for each program as the number of test cases increases, Figure 12 to 18 show the number of detected faults in ten stages – 100, 500, 1000, 1500, 2000, 2500, 3000, 3500, 4000 and 5000. Since different test suites may detect different numbers of faults in the 100 runs, the results were averaged over 100 runs, each time with a different seed and different test suites.

Looking at Figure 12, 13, and 15, it appears that MOClustering\_means has the best performance; MOClustering\_medoids performs best in Figure 14; and in Figure 16 to 18, it is DM-

1 Clustering that finds the most faults, regardless of the number of  
 2 test cases used. This is because the distance metric used in DM-  
 3 Clustering is more effective when applied to large-scale programs,  
 4 but MOClustering\_means and MOClustering\_medoids  
 5 are more effective in relatively small-scale programs. In Figure  
 6 14 and 16, we can observe that the lines of MOClustering\_  
 7 means, MOClustering\_medoid, DMClustering, Method\_C-  
 8 overage and RT-ms are almost coincident when the number of  
 9 test cases reaches 2000. The most appropriate explanation is  
 10 that the rates of fault detection for SATM and WaveletLibrary  
 11 are very high, and the faults are easily found. In Figure 17 and  
 12 18, all methods display a trend of finding more faults as the  
 13 number of test cases use increases. Through a further analys-  
 14 is, we can observe that some seeded faults in program IceChat  
 15 and CSPspEmu are very difficult to be detected by random test  
 16 cases, because they are associated with very lower failure rates.  
 17 Thus, 5000 test suite is not large enough to detect all faults for  
 18 these two programs, but large enough to detect all faults for the  
 19 other programs.

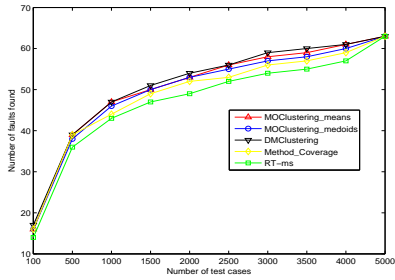


Figure 11: Relationship between average number of distinct faults found and number of test cases used for all seven subject programs

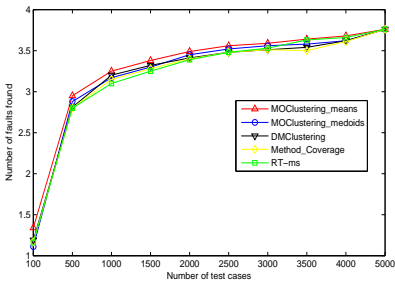


Figure 12: Relationship between the average number of faults found and the number of test cases used for CcoinBox

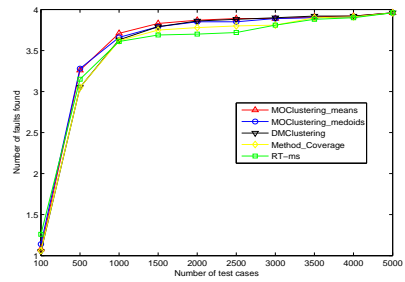


Figure 13: Relationship between the average number of faults found and the number of test cases used for WindShieldWiper

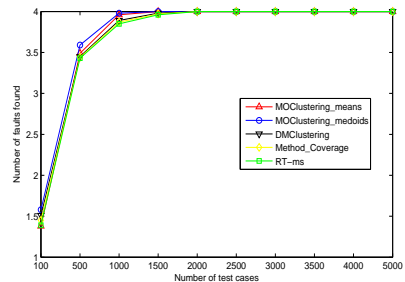


Figure 14: Relationship between the average number of faults found and the number of test cases used for SATM

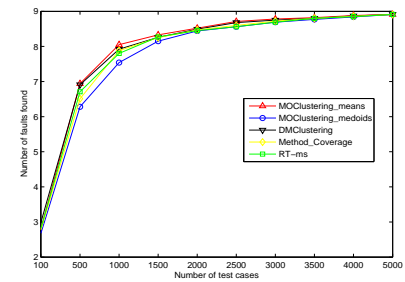


Figure 15: Relationship between the average number of faults found and the number of test cases used for RabbitsAndFoxes

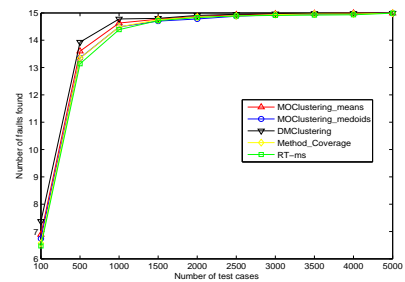


Figure 16: Relationship between the average number of faults found and the number of test cases used for WaveletLibrary

Table 9: COMPARISON BETWEEN VARIOUS PAIRS OF METHODS USING P-VALUE AND EFFECTIVE SIZE METHODS ON E WITH DIFFERENT NUMBERS OF TEST CASES

| Number of Test Cases | DMClustering and RT-ms |          | MOClustering _means and RT-ms |          | MOClustering _medoids and RT-ms |          | DMClustering and Method _Coverage |          | MOClustering _means and Method _Coverage |          | MOClustering _medoids and Method _Coverage |          |
|----------------------|------------------------|----------|-------------------------------|----------|---------------------------------|----------|-----------------------------------|----------|--|----------|--|----------|
|                      | P-value                | ES       | P-value                       | ES       | P-value                         | ES       | P-value                           | ES       | P-value                                  | ES       | P-value                                    | ES       |
| 100                  | 0.000017               | 0.642921 | 0.008701                      | 0.574365 | 0.005234                        | 0.566491 | 0.003482                          | 0.531962 | 0.008672                                 | 0.523585 | 0.009823                                   | 0.523491 |
| 500                  | 0.000236               | 0.623832 | 0.007124                      | 0.553474 | 0.004352                        | 0.537582 | 0.007573                          | 0.542852 | 0.010583                                 | 0.524774 | 0.198732                                   | 0.512827 |
| 1000                 | 0.000648               | 0.620743 | 0.009833                      | 0.544585 | 0.006763                        | 0.528643 | 0.006462                          | 0.533946 | 0.008472                                 | 0.535668 | 0.007410                                   | 0.526918 |
| 1500                 | 0.000092               | 0.615612 | 0.007721                      | 0.555694 | 0.005542                        | 0.539532 | 0.007573                          | 0.524739 | 0.006384                                 | 0.526754 | 0.009321                                   | 0.527826 |
| 2000                 | 0.000025               | 0.609758 | 0.005643                      | 0.576783 | 0.003631                        | 0.558443 | 0.005682                          | 0.535648 | 0.005493                                 | 0.537863 | 0.008432                                   | 0.528935 |
| 2500                 | 0.000138               | 0.607869 | 0.004532                      | 0.567892 | 0.005742                        | 0.549556 | 0.006894                          | 0.536757 | 0.007502                                 | 0.528974 | 0.009543                                   | 0.520624 |
| 3000                 | 0.000246               | 0.598750 | 0.008621                      | 0.558763 | 0.006853                        | 0.530447 | 0.008013                          | 0.547868 | 0.009611                                 | 0.538083 | 0.004432                                   | 0.528530 |
| 3500                 | 0.000571               | 0.579862 | 0.006732                      | 0.559854 | 0.003962                        | 0.551536 | 0.005124                          | 0.538757 | 0.006520                                 | 0.529195 | 0.007541                                   | 0.539423 |
| 4000                 | 0.000304               | 0.558751 | 0.003510                      | 0.548743 | 0.004851                        | 0.532425 | 0.006251                          | 0.529847 | 0.008651                                 | 0.525206 | 0.008657                                   | 0.528934 |
| 5000                 | 0.846479               | 0.509167 | 0.899845                      | 0.508654 | 0.913564                        | 0.501065 | 0.935468                          | 0.500957 | 0.957851                                 | 0.500672 | 0.965874                                   | 0.500478 |

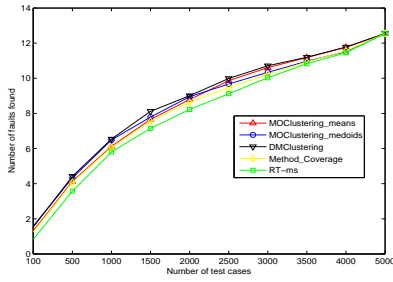


Figure 17: Relationship between the average number of faults found and the number of test cases used for IceChat

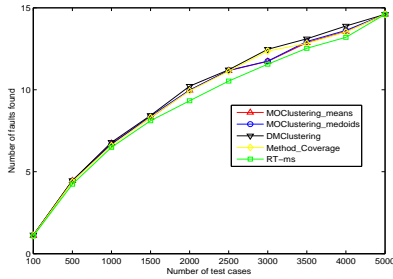


Figure 18: Relationship between the average number of faults found and the number of test cases used for CSPspEmu

From Table 8 and Figures 12 to 18, we have the following observations: the number of faults detected increases as  $n$  increases; based on the 5000 test inputs ( $TNum$  in Algorithms 1 to 3), DMClustering outperforms other methods, followed by MOClustering\_means, MOClustering\_medoids, Method\_Coverage and RT-ms (regardless of the value of  $n$ ).

In order to further analyze the significance of the difference in  $E$  with different test cases, we report in Table 9 the p-value and effective size (ES) for pairwise comparisons between the representative techniques from two different groups. We find that the difference between our methods and RT-ms is significant (because the p-value is less than 0.05), and the difference between our methods and Method\_Coverage is also significant, in most cases with different number of test cases. Through a further analysis of the ES values for different pairwise comparisons, we found that the ES values between our methods and RT-ms and Method\_Coverage are larger than 0.5,

which indicates that our methods perform better than RT-ms and Method\_Coverage. Amongst the three cluster TCP techniques, DMClustering performs best, followed by MOClustering\_means and MOClustering\_medoids. In addition, when the number of test cases is 5000, all methods have similar results, which can be seen based on the values of p-value and ES. The reason for this is that 5000 test cases can find most of the faults in most of the subject programs.

3) Do MOClustering\_means, MOClustering\_medoids and DMClustering perform better than prioritization with random sequences and Method\_Coverage, in terms of APFD?

Table 10 shows the average APFD values of the seven subject programs. All results were averaged over 100 runs, each time with a different seed.

Table 10: APFD OF VARIOUS TCP METHODS

| ID   | APFD                |                       |              |                  |       |
|------|---------------------|-----------------------|--------------|------------------|-------|
|      | MOClustering _means | MOClustering _medoids | DMClustering | Method _Coverage | RT-ms |
| 1    | 0.90                | 0.89                  | 0.88         | 0.88             | 0.87  |
| 2    | 0.93                | 0.93                  | 0.94         | 0.93             | 0.92  |
| 3    | 0.96                | 0.96                  | 0.96         | 0.95             | 0.94  |
| 4    | 0.92                | 0.90                  | 0.91         | 0.90             | 0.90  |
| 5    | 0.96                | 0.95                  | 0.96         | 0.95             | 0.94  |
| 6    | 0.70                | 0.70                  | 0.72         | 0.70             | 0.69  |
| 7    | 0.69                | 0.68                  | 0.76         | 0.68             | 0.67  |
| Mean | 0.87                | 0.86                  | 0.88         | 0.86             | 0.85  |
| sDev | 0.11                | 0.11                  | 0.10         | 0.12             | 0.12  |

As Table 10 shows, for program CCoinBox, MOClustering\_means performs best, followed by MOClustering\_medoids, DMClustering, Method\_Coverage and RT-ms. For program RabbitsAndFoxes, MOClustering\_means also performs best, and DMClustering outperforms MOClustering\_medoids, Method\_Coverage and RT-ms. For program SATM, the APFD values of three proposed methods are the same, and are much better than Method\_Coverage and RT-ms. In programs WindShieldWiper, WaveletLibrary, IceChat and CSPspEmu, DMClustering performs best, followed by MOClustering\_means, MOClustering\_medoids, Method\_Coverage and RT-ms in programs WindShieldWiper and WaveletLibrary, and followed by Method\_Coverage, MOClustering\_means, MOClustering\_medoids and RT-ms in programs IceChat and CSPspEmu. On average, DMClustering performs best, followed by MOClustering\_means, MOClustering\_medoids, Method\_Coverage and RT-ms.

In order to further analyze the APFD of each testing method

1 for each subject program, Table 11 summarizes the major s-  
 2 tistical measures including  $sDev$  (standard deviation) for the  
 3 7 subject programs. The standard deviation for RT-ms is the  
 4 biggest (0.12), which indicates that the data points are spread  
 5 out over a wider range of values than other TCP techniques.

6 Figure 19 to 25 show the APFD box-plots for the seven sub-  
 7 ject programs. In the figures, the x-axis has the prioritization  
 8 methods and the y-axis gives the APFD values for each method.

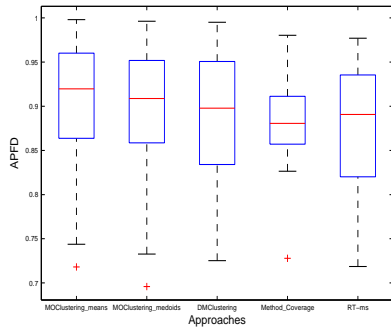


Figure 19: APFD values for CcoinBox

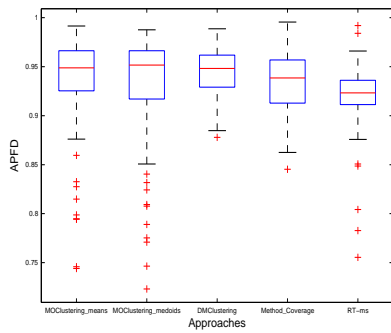


Figure 20: APFD values for WindShieldWiper

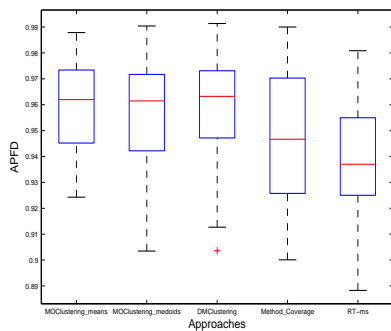


Figure 21: APFD values for SATM

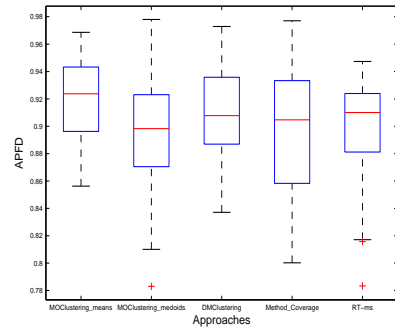


Figure 22: APFD values for RabbitsAndFoxes

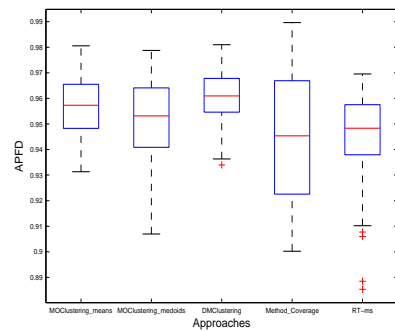


Figure 23: APFD values for WaveletLibrary

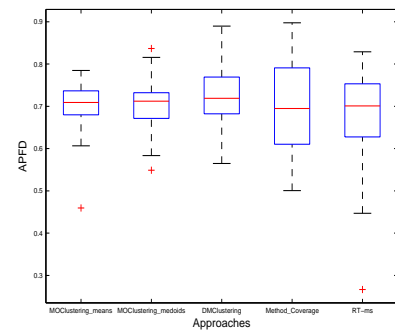


Figure 24: APFD values for IceChat

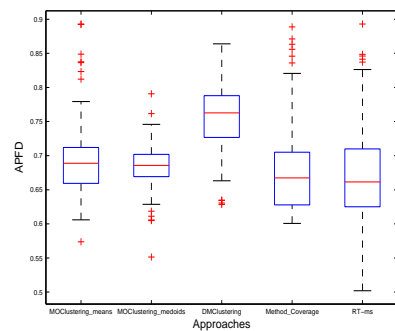


Figure 25: APFD values for CSPpEmu

Table 11: STATISTICAL RESULT OF APFD FOR 7 SUBJECT PROGRAMS

| ID |      | MOClustering<br>_means | MOClustering<br>_medoids | DMClustering | Method<br>_Coverage | RT-ms |
|----|------|------------------------|--------------------------|--------------|---------------------|-------|
| 1  | mean | 0.90                   | 0.89                     | 0.88         | 0.88                | 0.87  |
|    | sDev | 0.06                   | 0.05                     | 0.06         | 0.04                | 0.07  |
| 2  | mean | 0.93                   | 0.93                     | 0.94         | 0.93                | 0.92  |
|    | sDev | 0.03                   | 0.03                     | 0.02         | 0.03                | 0.03  |
| 3  | mean | 0.96                   | 0.96                     | 0.96         | 0.95                | 0.94  |
|    | sDev | 0.02                   | 0.02                     | 0.02         | 0.03                | 0.02  |
| 4  | mean | 0.92                   | 0.90                     | 0.91         | 0.90                | 0.90  |
|    | sDev | 0.03                   | 0.02                     | 0.03         | 0.05                | 0.03  |
| 5  | mean | 0.96                   | 0.95                     | 0.96         | 0.95                | 0.94  |
|    | sDev | 0.01                   | 0.02                     | 0.01         | 0.03                | 0.02  |
| 6  | mean | 0.70                   | 0.70                     | 0.72         | 0.70                | 0.69  |
|    | sDev | 0.05                   | 0.05                     | 0.07         | 0.11                | 0.09  |
| 7  | mean | 0.69                   | 0.68                     | 0.76         | 0.68                | 0.67  |
|    | sDev | 0.06                   | 0.03                     | 0.05         | 0.07                | 0.07  |

Table 12: COMPARISON BETWEEN VARIOUS PAIRS OF METHODS USING P-VALUE AND EFFECTIVE SIZE METHODS ON APFD

| Pair of methods      | MOClustering<br>_means and<br>RT-ms | MOClustering<br>_medoids and<br>RT-ms | DMClustering<br>and RT-ms | MOClustering<br>_means and<br>Method<br>_Coverage | MOClustering<br>_medoids and<br>Method<br>_Coverage | DMClustering<br>and Method<br>_Coverage | MOClustering<br>_means and<br>DMClustering | MOClustering<br>_medoids and<br>DMClustering | MOClustering<br>_means and<br>MOClustering<br>_medoids |
|----------------------|-------------------------------------|---------------------------------------|---------------------------|---|---|---|--|--|--|
| <b>P-value</b>       | 8.58E-08                            | 0.002017                              | 2.59E-09                  | 0.003581  | 0.004118  | 0.001594                                | 0.202008                                   | 0.197287                                     | 0.604694   |
| <b>ES</b>            | 0.582653                            | 0.547663                              | 0.591930                  | 0.544963  | 0.535095  | 0.548731                                | 0.480305                                   | 0.480291                                     | 0.503989   |
| <b>Better Method</b> | MOClustering<br>_means              | MOClustering<br>_medoids              | DMClustering              | MOClustering<br>_means                            | MOClustering<br>_medoids                            | DMClustering                            | DMClustering                               | DMClustering                                 | MOClustering<br>_means                                 |

In Figures 19, 20 and 21, the upper quartile and median values of three cluster TCP techniques are all higher than those of Method.Coverage and RT-ms, implying a better performance with respect to APFD values. In Figure 22, the lower quartile, median and upper quartile values for MOClustering\_means are all higher than those of the other methods. In Figures 23, 24, and 25, the lower quartile and median values for DMClustering are all higher than those of the other methods.

In order to further analyze the significance of the difference in APFD, we report in Table 12 the p-value and effective size (ES) for pairwise comparisons between the representative techniques from two different groups. We find that the difference between our approaches and RT-ms and Method.Coverage is significant (because the p-value is less than 0.05). But the difference among our proposed clustering methods is not significant (because the p-value is larger than 0.05). Through a further analysis on the ES values for different pairwise comparisons, we can find that these values between our methods and RT-ms and Method.Coverage are larger than 0.5, which indicates that our methods perform better than RT-ms and Method.Coverage. Column "Result" of Table 12 gives the better method of the relevant pair. Amongst three cluster TCP techniques, DMClustering performs best, followed by MOClustering\_means and MOClustering\_medoids on average.

In summary, based on Table 4 to 12, and Figure 4 to 25, we have the following observations: (1) DMClustering performs better than other methods for larger programs (WaveletLibrary, IceChat and CSPspEmu); (2) MOClustering\_means and MOClustering\_medoids have good performance with smaller programs, and MOClustering\_means is relatively more effective than MOClustering\_medoids in terms of  $F_m$ ,  $E$  and APFD; (3)

generally speaking, DMClustering performs best in terms of  $F_m$ ,  $E$  and APFD on average; (4) the three cluster TCP techniques (DMClustering, MOClustering\_means and MOClustering\_medoids) outperform Method.Coverage (method coverage prioritization) in terms of  $F_m$ ,  $E$  and APFD; and (5) the four TCP techniques (DMClustering, MOClustering\_means, MOClustering\_medoids and Method.Coverage) perform better than RT-ms (prioritization with random sequences) in terms of  $F_m$ ,  $E$  and APFD. All the above outperformance is statistically significant.

#### 5.4.2. Threats to validity

Although we believe that the experiment was well-designed and implemented, the study may still face some threats to its validity, as explained in the following. In the clustering algorithms, the number of clusters  $K$  is generally required to be known in advance. Obviously, the value of  $K$  has a significant influence on the clustering quality. Hence, if  $K$  is not correctly chosen, then the clustering analysis algorithm may produce low quality results. In this study, the value for  $K$  was determined experimentally, but in some other studies, it was determined according to the gap statistic algorithm [42] and the distribution characteristics of the test cases. In addition, the subject programs were downloaded from some open source websites, but these subject programs may not be associated with any test cases. Although we tried our best to find some OO programs (in C# or C++), with real test cases for OO integration testing in some famous software repositories such as the Software Infrastructure Repository (SIR) [43], unfortunately, we did not find suitable ones. Hence, we developed a tool which randomly generates test cases for these subject programs. In the absence of real test suites, we believe that random test suites are fair and



1 reasonable solutions.

2 In this study, the mutants in the seven subject program-  
3 s were generated by hand, due to a lack of good automatic  
4 mutation tools for both C++ and C# programs. However, the  
5 location and type of seeded faults were selected using a random  
6 number generator, thus making the process semirandom  
7 and semiautomatic. Additionally, in order to reduce the threats,  
8 we manually filter as many subsumed mutants [44] as possible.

## 9 6. Related Work

10 Chen et al. [17] first suggested how to use ART in test case  
11 prioritization, calling such an approach an adaptive random sequence  
12 (ARS), and explaining how it could be a cost-effective  
13 alternative to random sequences. Rothermel et al. [27] proposed  
14 several code coverage based TCP approaches, including total  
15 statement coverage prioritization, additional statement  
16 coverage prioritization, total branch coverage prioritization,  
17 and total fault-exposing potential prioritization. Their experimental  
18 results show that these methods can improve the fault  
19 detection rates of test suites.

20 Cluster analysis has drawn a lot of attention in the TCP  
21 community. Dickinson et al. [45] proposed a clustering based  
22 test case filtering technique that improves on the efficiency of  
23 random sampling by using an agglomerative hierarchical clustering  
24 algorithm, which is a bottom-up approach, where each  
25 test case is used as a cluster, and the clusters with minimal  
26 dissimilarity are merged into larger clusters until a predefined  
27 number of clusters remain. They studied several dissimilarity  
28 metrics, including binary metric, proportional metric, SD (standard  
29 deviation) metric, histogram metric, linear regression metric,  
30 count-binary metric, and proportional-binary metric. The  
31 inputs to the cluster analysis are function call profiles. In the  
32 profile, each pair of methods is represented as an entry showing  
33 the frequency of the executed methods. Although this approach  
34 can reflect the dynamic behavior of test cases, only the methods'  
35 execution time (including whether or not the method is  
36 executed) is used.

37 Yoo et al. [46] proposed a cluster-based TCP technique that  
38 significantly reduces the required number of pair-wise comparisons.  
39 Their clustering method partitions test cases into different  
40 subsets based on their dynamic runtime behavior, with  
41 test cases in each group having common properties. The clustering  
42 approach uses binary strings to represent test inputs and  
43 whether or not a statement is executed: If the source code  
44 statement has been executed, the digit of the corresponding bit  
45 in the binary string is set to 1; otherwise, it is set to 0. Zhang  
46 et al. [15] proposed online and offline ARS-based TCP techniques  
47 using black-box information based on the string distances of the  
48 input data, without referring to the execution history and code  
49 coverage information. The offline TCP algorithm selects new  
50 test cases farthest from all prioritized ones; and the online algo-  
51 rithm uses feedback information such that the next prioritized  
52 test case depends on the existing execution results.

## 53 7. Conclusion And Future Work

54 Software testing is an important aspect of examining the  
55 quality and reliability of object-oriented software (OOS). Be-  
56 cause OOS test cases may be very complex, traditional software  
57 testing approaches may not be appropriate for testing OOS. Al-  
58 though studies have been carried out to enhance OOS testing,  
59 OOS test case prioritization (TCP) has not yet been fully ex-  
60 plored. TCP can increase fault detection rates by optimizing  
61 test case execution sequences such that more important test cas-  
62 es are executed earlier - based on some criteria. Cluster analysis  
63 has recently been applied to improving TCP effectiveness.

64 In this paper, in order to improve the effectiveness of TCP  
65 for OOS, we have proposed an ARS approach based on cluster-  
66 ing techniques. We used three clustering methods to define our  
67 ARS methods: MOClustering\_means, MOClustering\_medoids  
68 and DMClustering. In MOClustering\_means and MOClustering\_  
69 medoids, test cases are clustered according to the number  
70 of objects and methods, using K-means and K-medoids cluster-  
71 ing algorithms. In these two methods, the Object Method  
72 Vector is constructed to calculate the distance between test cas-  
73 es using the Euclidean distance formula. In DMClustering, test  
74 cases are clustered based on an object and method invocation  
75 sequence similarity (OMISS) metric with the K-medoids cluster-  
76 ing algorithm. Furthermore, a sampling strategy MSampling  
77 is used to construct the ARSs. The final prioritized test case  
78 sequence is generated from the  $K$  clusters. The experimen-  
79 tal results show that the three proposed cluster methods out-  
80 perform Method.Coverage and RT-ms in terms of  $F_m$ ,  $E$  and  
81 APFD; and DMClustering performs best overall, and is there-  
82 fore a good choice for test case prioritization, especially for  
83 large scale OOS testing. Furthermore, all the better perfor-  
84 mances are statistically significant.

85 Based on the observations from our experimentations, we  
86 recommend that for large programs, it is better to set  $K$  to around  
87 15% of the total number of test cases, and for small programs,  
88 it is better to set it to around 10%.

89 In future, we will conduct further investigations into OOS  
90 test case features, and add other important information to the  
91 Object Method Vector and OMISS metric to enhance the prob-  
92 ability of the selected test cases for OOS to be more evenly  
93 spread across the input domain. We also will improve the sam-  
94 pling strategy to better optimize the TCP test cases.

## 95 Acknowledgement

96 This work is partly supported by National Natural Science Foun-  
97 dation of China (NSFC grant numbers: 61202110 and 61502205),  
98 and the Postdoctoral Science Foundation of China (Grant num-  
99 bers: 2015M571687 and 2015M581739).

## 100 References

- 101 [1] R. V. Binder, "Testing object-oriented software: a survey," *Software Test-*  
102 *ing Verification and Reliability*, vol. 6, no. 3, pp. 125–252, 1996.
- 103 [2] M. Pezze and M. Young, "Testing object-oriented software," in *26th*  
104 *International Conference on Software Engineering Proceedings (ICSE*  
105 *2004)*, United Kingdom, pp. 739–740, IEEE, 2004.

- [3] H. Y. Chen, T. H. Tse, F. T. Chan, and T. Y. Chen, "In black and white: an integrated approach to class-level testing of object-oriented programs," *Acm Transactions on Software Engineering and Methodology*, vol. 7, no. 3, pp. 250–295, 1998.
- [4] S. Anand, E. K. Burke, T. Y. Chen, J. Clark, M. B. Cohen, W. Grieskamp, M. Harman, M. J. Harrold, P. McMinn, and A. Bertolino, "An orchestrated survey of methodologies for automated software test case generation," *Journal of Systems and Software*, vol. 86, no. 8, pp. 1978–2001, 2013.
- [5] N. E. Holt, L. C. Briand, and R. Torkar, "Empirical evaluations on the cost-effectiveness of state-based testing: An industrial case study," *Information and Software Technology*, vol. 56, no. 8, pp. 890–910, 2014.
- [6] S. U. Hui, Y. Zhang, H. Yao, and R. Fei, "Object-oriented software cluster-level testing based on uml sequence diagram," *Computer Engineering*, vol. 31, no. 24, pp. 78–80, 2005.
- [7] R. Hamlet, "Random testing," *Encyclopedia of Software Engineering*, John Wiley and Sons, 2002.
- [8] T. Y. Chen, K. Fei-Ching, T. Dave, and Z. Z. Quan, "A revisit of three studies related to random testing," *Science China Information Sciences*, vol. 58, no. 5, pp. 1–9, 2015.
- [9] C. Nie, H. Wu, X. Niu, F. C. Kuo, H. Leung, and C. J. Colbourn, "Combinatorial testing, random testing, and adaptive random testing for detecting interaction triggered failures," *Information and Software Technology*, vol. 62, pp. 198–213, 2015.
- [10] Y. Ledru, A. Petrenko, S. Boroday, and N. Mandran, "Prioritizing test cases with string distances," *Automated Software Engineering*, vol. 19, no. 1, pp. 65–95, 2012.
- [11] I. Ciupa, A. Leitner, M. Oriol, and B. Meyer, "Artoo:adaptive random testing for object-oriented software," in *ACM/IEEE 30th International Conference on Software Engineering (ICSE 2008)*, IEEE, New York, USA, pp. 71–80, 2008.
- [12] S. Elbaum, A. G. Malishevsky, and G. Rothermel, "Test case prioritization: A family of empirical studies," *IEEE Transactions on Software Engineering*, vol. 28, no. 2, pp. 159–182, 2002.
- [13] X. Zhang, X. Xie, and T. Y. Chen, "Test case prioritization using adaptive random sequence with category-partition-based distance," in *IEEE International Conference on Software Quality, Reliability and Security (QRS 2016)*, IEEE, Washington, USA, pp. 374–385, 2016.
- [14] Q. Luo, K. Moran, and D. Poshvanyk, "A large-scale empirical comparison of static and dynamic test case prioritization techniques," in *ACM SIGSOFT Symposium on the Foundations of Software Engineering (FSE 2016)*, ACM, Washington, USA, pp. 559–570, 2016.
- [15] X. Zhang, T. Y. Chen, and H. Liu, "An application of adaptive random sequence in test case prioritization," in *International Conference on Software Engineering and Knowledge Engineering (SEKE 2014)*, IEEE, Canada, pp. 126–131, 2014.
- [16] T. Y. Chen, H. Leung, and I. K. Mak, "Adaptive random testing," in *Proceedings of the 9th Asian Computing Science Conference (ASIAN 2004)*, Thailand, pp. 320–329, Springer LNCS, 2004.
- [17] T. Y. Chen, F. C. Kuo, R. G. Merkel, and T. H. Tse, "Adaptive random testing: The art of test case diversity," *Journal of Systems and Software*, vol. 83, no. 1, pp. 60–66, 2010.
- [18] T. Y. Chen, F. C. Kuo, H. Liu, and W. E. Wong, "Code coverage of adaptive random testing," *IEEE Transactions on Reliability*, vol. 62, no. 1, pp. 226–237, 2013.
- [19] A. C. Barus, T. Y. Chen, F. C. Kuo, H. Liu, R. Merkel, and G. Rothermel, "A cost-effective random testing method for programs with non-numeric inputs," *IEEE Transactions on Computers*, vol. 65, no. 12, pp. 3509–3523, 2016.
- [20] H. Liu, F. C. Kuo, D. Towey, and T. Y. Chen, "How effectively does metamorphic testing alleviate the oracle problem?," *IEEE Transactions on Software Engineering*, vol. 40, no. 1, pp. 4–22, 2014.
- [21] J. Chen, L. Zhu, T. Y. Chen, R. Huang, D. Towey, F. C. Kuo, and Y. Guo, "An adaptive sequence approach for oos test case prioritization," in *IEEE International Symposium on Software Reliability Engineering Workshops (ISSRE-IWPD 2016)*, IEEE, Canada, pp. 205–212, 2016.
- [22] S. M. Aqilburney and H. Tariq, "K-means cluster analysis for image segmentation," *International Journal of Computer Applications*, vol. 96, no. 4, pp. 1–8, 2014.
- [23] J. Chen, F. C. Kuo, T. Y. Chen, D. Towey, C. Su, and R. Huang, "A similarity metric for the inputs of oo programs and its application in adaptive random testing," *IEEE Transactions on Reliability*, vol. PP, no. 99, pp. 1–30.
- [24] O. Legunsen, F. Hariri, A. Shi, Y. Lu, L. Zhang, and D. Marinov, "An extensive study of static regression test selection in modern software evolution," in *ACM Sigsoft International Symposium on Foundations of Software Engineering (FSE 2016)*, ACM, Washington, USA, pp. 583–594, 2016.
- [25] A. Gonzalez-Sanchez, E. Piel, R. Abreu, H. G. Gross, and A. J. C. Van Gemund, "Prioritizing tests for software fault diagnosis," *Software Practice and Experience*, vol. 41, no. 10, pp. 1105–1129, 2011.
- [26] S. R. Suganya and G. S. Devi, "Data mining: concepts and techniques," *Data Mining and Knowledge Engineering*, pp. 929–948, 2010.
- [27] G. Rothermel, R. H. Untch, C. Chu, and M. J. Harrold, "Prioritizing test cases for regression testing," *IEEE Transactions on Software Engineering*, vol. 27, no. 10, pp. 929–948, 2002.
- [28] C. Henard, M. Papadakis, M. Harman, Y. Jia, and Y. L. Traon, "Comparing white-box and black-box test prioritization," in *International Conference on Software Engineering*, pp. 523–534, 2016.
- [29] "Microsoft visual studio." <https://www.visualstudio.com>. 2013.
- [30] L. Kaufmann and P. J. Rousseeuw, "Clustering by means of medoids," in *Statistical Data Analysis Based on the L1-norm and Related Methods*, North-Holland, pp. 405–416, 1987.
- [31] H. S. Park and C. H. Jun, "A simple and fast algorithm for k-medoids clustering," *Expert Systems with Applications*, vol. 36, no. 2, pp. 3336–3341, 2009.
- [32] "Codeforce-free open source codes forge and sharing." <http://www.codeforce.com>. 2013.
- [33] "Sourceforge-download, develop and publish free open source software." <http://sourceforge.net>. 2013.
- [34] "Codeplex-open source project hosting." <http://www.codeplex.com>. 2013.
- [35] "Github, where software is built." <https://github.com>. 2015.
- [36] Y. Jia and M. Harman, "An analysis and survey of the development of mutation testing," *IEEE Transactions on Software Engineering*, vol. 37, no. 5, pp. 649–678, 2010.
- [37] T. Y. Chen and R. Merkel, "An upper bound on software testing effectiveness," *Acm Transactions on Software Engineering and Methodology*, vol. 17, no. 3, pp. 1–27, 2008.
- [38] A. Arcuri and L. Briand, "A hitchhiker's guide to statistical tests for assessing randomized algorithms in software engineering," *Software Testing Verification and Reliability*, vol. 24, no. 3, pp. 219–250, 2014.
- [39] A. Vargha and H. D. Delaney, "A critique and improvement of the CL common language effect size statistics of mcgraw and wong," *Journal of Educational and Behavioral Statistics*, vol. 25, no. 2, pp. 101–132, 2000.
- [40] M. Harman, P. McMinn, J. Souza, and S. Yoo, "Search based software engineering: techniques, taxonomy, tutorial," in *Empirical Software Engineering and Verification*, pp. 1–59, 2012.
- [41] R. D. C. Team, "Development core team, r: A language and environment for statistical computing," 2009.
- [42] R. Tibshirani, G. Walther, and T. Hastie, "Estimating the number of clusters in a data set via the gap statistic," *Journal of the Royal Statistical Society*, vol. 63, no. 2, pp. 411–423, 2001.
- [43] "Software-artifact infrastructure repository." <http://sir.unl.edu/portal/index.php>. 2016.
- [44] M. Papadakis, C. Henard, M. Harman, Y. Jia, and Y. L. Traon, "Threats to the validity of mutation-based test assessment," in *International Symposium on Software Testing and Analysis*, pp. 354–365, 2016.
- [45] W. Dickinson, D. Leon, and A. Fodgurski, "Finding failures by cluster analysis of execution profiles," in *23rd IEEE International Conference on Software Engineering Proceedings (ICSE 2001)*, IEEE, Ontario, Canada, pp. 339–348, 2001.
- [46] S. Yoo, M. Harman, P. Tonella, and A. Susi, "Clustering test cases to achieve effective and scalable prioritisation incorporating expert knowledge," in *Eighteenth International Symposium on Software Testing and Analysis, ISSTA 2009, Chicago, Il, Usa, July*, pp. 201–212, 2009.