

# AMT: An Efficient and Novel Method to Improve the Fault-detecting Efficiency of Metamorphic Testing

Chang-ai Sun, *Senior Member, IEEE*, Hepeng Dai, Huai Liu, *Member, IEEE*, and Tsong Yueh Chen, *Member, IEEE*,

**Abstract**—Metamorphic testing (MT) is a promising technique to alleviate the oracle problem, which first defines metamorphic relations (MRs) that are then used to generate new test cases (i.e. follow-up test cases) from the original test cases (i.e. source test cases), and verify the results of source and follow-up test cases. Many efforts have been reported to improve MT's efficiency by either generating better MRs that are more likely to be violated or selecting different test case selection strategies to generate source test cases. Unlike these efforts, we investigate how to improve the efficiency of MT in terms of test executions. Furthermore, traditional MT techniques often employ the random testing strategy (RT) to select source test cases for execution, which could be inefficient because the feedback information during the testing process is not leveraged. Consequently, we propose an adaptive metamorphic testing (AMT) technique to improve the efficiency of MT through controlling the execution process of MT. In the process of testing, AMT can quickly find source test cases and MRs with fault detected ability based on historical information of testing. We conducted an empirical study to evaluate the efficiency of the proposed technique with four laboratory programs, a GNU program, and an Alibaba program. Empirical results show that AMT outperforms traditional MT in terms of fault-detecting efficiency.

**Index Terms**—metamorphic testing, control test process, feedback, random testing, adaptive random testing, partition testing, adaptive partition testing

## 1 INTRODUCTION

TEST result verification is an essential step of software testing. A test oracle [1] is a mechanism that can exactly decide whether the output produced by a programs is correct. However, there are situations where it is difficult to decide whether the result of the software under test (SUT) agrees with the expected result. This situation is known as oracle problem [2], [3]. In order to alleviate the oracle problem, several techniques have been proposed, such as N-version testing [4], metamorphic testing (MT) [5], [6], assertions [7], and machine learning [8]. Among of them, MT obtains metamorphic relations (MRs) according to the properties of SUT. Then, MRs are used to generate new test cases called follow-up test cases from original test cases known as source test cases. Next, both source and follow-up test cases are executed and their result are verified against the corresponding MRs. MT checks whether a relevant MR holds among multiple executions. If an MR does not hold, then it can be concluded that the software system is faulty.

MT has been drawing increasing attention in the software testing community [6], [9], since this technique not only alleviates the oracle problem but also generates new

test cases from existing test cases. The technique has been successfully applied in a number of application domains and paradigms, including healthcare [10], bioinformatics [11], air traffic control [12], the web service [13], the RESTful web APIs [14], and testing of AI systems [15]–[17]. Observations from the existing evaluations indicate that MT is a new and promising strategy that complements the existing testing approaches [18].

There are two broad categories of methods to improve the effectiveness of MT: source test cases generating and MRs identifying. The former employs different test case selection strategies to generate source test cases for MT [19], [20]. The latter creates MRs by combining existing relations, or by generating them automatically [21]–[25].

The fault-detecting effectiveness of MT relies on the quality of MRs and the source test cases. Chen et al. [19] initially suggested all available MRs should be used as part of the test strategy. In the context of the resources for software development are always limited, Chen et al. [21] proposed “good” MRs should be those that can make the multiple executions of the program as different as possible. Furthermore, they concluded that good MRs should be preferably selected with regard to the algorithm under test following a white-box approach. However, this was later disputed by Mayer and Guderlei [26], who studied six subject programs for matrix determinant computation with seeded faults. They concluded that metamorphic relations in the form of equalities or linear equations as well as those close to the implementation strategy have limited effectiveness. Conversely, they reported that good metamorphic relations are usually strongly inspired by the semantics of

C.-A. Sun, and H. Dai are with the School of Computer and Communication Engineering, University of Science and Technology Beijing, Beijing 100083, China. E-mail: casun@ustb.edu.cn.

H. Liu is with the College of Engineering and Science, Victoria University, Melbourne VIC 8001, Australia. E-mail: Huai.Liu@vu.edu.au.

T.Y. Chen is with the Department of Computer Science and Software Engineering, Swinburne University of Technology, Hawthorn VIC 3122, Australia. Email: tychen@swin.edu.au.

the program under test. Asrafi et al. [27] concluded that the higher the combined code coverage of the source and follow-up test cases, the more different are the executions, and the more effective is the MR through a case study. The above three methods discuss the properties of “good” MRs from the point of view of the black-box and the white-box. On the other hand, Just et al. [28] assessed the applicability of metamorphic testing for system and integration testing in the context of an image encoder, and concluded that the MRs derived from the components of a system are usually better at detecting faults than those MRs derived from the whole system. To select “good” MRs, all of these methods require test cases to be executed before testing SUT. Without doubt, testing overhead is increased.

As for selection source test cases, 57% of existing studies employed random testing (RT) to select test cases, and 34% of existing studies used existing test suites according to a survey report by Segura et al. [9]. RT randomly selects test cases from input domain (which refers to the set of all possible inputs of SUT). Although RT is simple to implement, RT does not make use of any execution information about SUT or the test history. Thus, traditional MT may be ineffectiveness in some situations. The existing test suite may not adequately test the new version SUT. Because of the above limitations, Barus et al. [20] employed adaptive random testing (ART) [29] that is a class of testing method aimed to improve the performance of RT by increasing the diversity across a program’s input domain of the test cases executed to generate source test cases for MT. However, the characteristics of ART itself (such as high time consumption, most algorithms only apply to numerical programs) will increase the time consumption of MT.

In previous and existing work, it is difficult for practitioners to choose the MRs based on the existing work or need to consume a lot of resources to choose those “good” MRs. On the other hand, researchers have mainly focused on the impact of MRs, with the impact of source test cases on MTs fault-detection effectiveness having somehow been neglected. In other words, investigation of the impact of source test cases on MR (and MT) effectiveness is an area yet to be explored [6]. Furthermore, researchers investigated the effectiveness of MT by either choosing different MRs alone, or choosing different test case selection strategies to generate source test cases alone.

Different from the previous methods, this paper proposed an innovative method that makes use of feedback information to select both source test case and corresponding MR based on software cybernetics [30], [31], aiming to maximize the fault-detection effectiveness of MT.

In this study, we investigate how to make use of feedback information in the previous tests to control the execution process of MT in terms of both selecting source test cases and MRs, and proposed adaptive metamorphic testing method (AMT). AMT is built on top of three insights: i) source test cases is one of factors that effects the fault-detecting efficiency of MT. However, random testing is commonly used to generate source test cases, and only a little researches related to use appropriate test cases selection strategies to generate source test cases [9]; ii) most of existing work help testers create MRs, aiming to improve the fault-detection efficiency of MT, while there seems to

be no investigation to study how to choose an MR based on the generated set of MRs; iii) an appropriate test cases selection strategy can generate “better” test cases (i.e. those test cases have a higher probability to detect a fault). In MT, if source test cases cannot detect a fault, and follow-up test case generated by an MR detect a fault, we still can conclude that this execution detect a fault. Therefore, based on a selected source test case, choosing a “better” MR can improve the fault-detection efficiency.

In short, AMT collects feedback information during the test, and makes use of those information to select source test case and MRs. As a result, a process-feedback metamorphic testing framework is proposed to improve the fault-detecting efficiency of MT. An empirical study was conducted to evaluate the efficiency of the proposed technique. Main contributions made in this paper are the following:

- 1 From a new perspective, we proposed a adaptive metamorphic testing method. This includes a general framework (AMT) that indicates how to make use of history testing information to select next source test case and dynamically select MR.
- 2 We particular developed two algorithms (MAPT\* and RAPT\*) to select source test cases, and two algorithms (PBMR and RMRS) to select MRs for MT. Then the test case selection algorithms and MR selection algorithms were combined to obtain four AMT strategies: i) MR-AMT and RR-AMT (MAPT\* selects test cases and RMRS selects MRs, and RAPT\* selects test cases and RMRS selects MRs, respectively), which first select a partition according to the test profile, and randomly select a test case from the selected partition and randomly selects a MR from the set of MRs whose source test cases belong to selected partition, then update test profile according to the result of test execution; ii) MP-AMT and RP-AMT (MAPT\* selects test cases and PBMR selects MRs, and RAPT\* selects test cases and PBMRs selects MRs, respectively) which first select a partition according to the test profile, and randomly select a test case from the selected partition and selects a MR from the set of MRs whose source test cases belong to selected partition according to PBMR algorithms, then update test profile according to the result of test execution.
- 3 We evaluated the performance of AMT through a series of empirical studies on 6 programs, including four laboratory programs, GNU `grep`, and Alibaba `FastJson`. These studies show that AMT has significantly higher fault-detection efficiency than traditional MT that randomly selects source test cases and MRs.

The rest of this paper is organized as follows. Section 2 introduces the underlying concepts for MT, APT and CPM. Section 3 presents the AMT framework, the strategies of test case selection, and the strategies of MRs selection. Section 4 describes an empirical study where the proposed AMT is used to test four laboratory programs, a Gun program, and an Alibaba program, the results of which are summarized in Section 5. Section 6 discusses related work.

## 2 BACKGROUND

In this section, we present some of the underlying concepts for MT, and APT.

## 2.1 Metamorphic Testing (MT)

MT is a novel technique to alleviate the oracle problem: Instead of applying an oracle, MT uses a set of MRs (corresponding to some specific properties of the SUT) to verify the test results [5]. MT is normally conducted according to the following steps:

- Step1. Identify an MR from the specification of the SUT.
- Step2. Generate the source test case  $stc$  using the traditional test cases generation techniques.
- Step3. Derive the follow-up test case  $ftc$  from the  $stc$  based on the MR.
- Step3. execute  $stc$  and  $ftc$  and get their outputs  $O_s$  and  $O_f$ .
- Step4. Verify  $stc$ ,  $ftc$ ,  $O_s$ , and  $O_f$  against the MR: If the MR does not hold, a fault is said to be detected.

The above steps could be repeated for a set of MRs.

Let us use a simple example to illustrate how MT works. For instance, consider a program  $P(G, a, b)$  purportedly calculating the length of the shortest path between nodes  $a$  and  $b$  in an undirected graph  $G$ . When  $G$  is nontrivial, the program is difficult to test because no oracle can be practically applied. Nevertheless, we can perform MT. Let  $(G_1, a_1, b_1)$  and  $(G_2, a_2, b_2)$  be two inputs, where  $G_2$  is a permutation of  $G_1$  (that is,  $G_2$  and  $G_1$  are isomorphic), and  $(a_1, b_1)$  in  $G_1$  correspond to  $(a_2, b_2)$  in  $G_2$ . Then an MR can be identified as follows:  $P(G_1, a_1, b_1) = P(G_2, a_2, b_2)$ . A metamorphic test using this MR will run  $P$  twice, namely, a source execution on the source test case  $(G_1, a_1, b_1)$  and a follow-up execution on the follow-up test case  $(G_2, a_2, b_2)$ .

## 2.2 Category Partition Method (CPM)

Category partition method (CPM) is widely used specification-based testing technique [33]. It helps software testers create test cases by refining the functional specification of a program into test specifications. The method consists of the following steps.

- Step1. Decompose the functional specification into functional units that can be tested independently, then Identify the parameters (the explicit inputs to a functional unit) and environment conditions (the state of the system at the time of execution) that are known as *categories*.
- Step2. Partition each *category* into *choices*, which include all the different kinds of values that are possible for that *category*.
- Step3. Determine the constraints among the *choices* of different *categories*, and Write the test specification (which is a list of categories, choices, and constraints in a predefined format) using the test specification language TSL.
- Step4. Use a generator to produce *test frames* from the test specification. Each generated *test frame* is a set of choices. Then, create a test case by selecting a single element from each *choice* in each generated *test frame*.

In this study, we made use of CPM to create test cases. During the process of testing, the choice information can be used as the basis for choosing the MRs. The proposed

AMT selects both source test cases and corresponding MRs by making use of information which is collected during the test process that means, we need to establish a connection between the source test cases and MRs. The connection provides a guideline for us to obtain the set of candidate MRs according to a source test case or select a source test case for an MR.

## 2.3 Adaptive Partition Testing (APT)

Based on software cybernetics, which aims to explore the interplay between software engineering and control theory, Sun [32] propose a new testing approach, adaptive partition testing (APT), where test cases are randomly selected from some partition  $c_i$  whose probability  $p_i$  of being selected is adaptively adjusted along the testing process. The main concept of APT is to increase the selection probabilities of the partitions which have higher defect detected rates. If a defect is detected by some test cases from partition  $c_i$ , then  $c_i$  is considered to have a higher defect detected rate and an increment is added to  $p_i$ ; Otherwise,  $p_i$  is decreased.

Furthermore, they particularly develop two algorithms, Markov-chain based adaptive partition testing (MAPT) and reward-punishment based adaptive partition testing (RAPT), to implement the proposed approach. MAPT and RAPT algorithms assume that testers can obtain test oracles, therefore, MAPT and RAPT are difficult to apply when there are no test oracles or it is difficult to obtain the test oracles.

### 2.3.1 Markov-Chain Based Adaptive Partition Testing (MAPT)

According to the concept of Markov chain, given two states  $i$  and  $j$ , the probability of transitioning from  $i$  to  $j$  is represented by  $p_{i,j} = Pr\{j|i\}$ . For  $i = 1, 2, \dots, m$  we can then construct a Markov matrix as follows:

$$\mathcal{P} = \begin{pmatrix} p_{1,1} & p_{1,2} & \cdots & p_{1,m} \\ p_{2,1} & p_{2,2} & \cdots & p_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ p_{m,1} & p_{m,2} & \cdots & p_{m,m} \end{pmatrix} \quad (1)$$

In MAPT, each partition is considered as a state in the Markov matrix. If a partition  $s_i$  is selected for conducting a test, the probability of selecting  $s_i$  for conducting the next test will be  $p_{i,j}$ . MAPT will adaptively adjust the value of each  $p_{i,j}$  according to the testing result. The algorithm of MAPT is available in [32].

### 2.3.2 reward-punishment based adaptive partition testing (RAPT)

Based on the reward and punishment mechanism, RAPT attempts to be quicker at selecting the fault-revealing test cases. Two parameters  $Rew_i$  and  $Pun_i$  are used in RAPT to determine to what extent a partition  $c_i$  can be rewarded and punished, respectively. If a test case in  $c_i$  reveals a fault,  $Rew_i$  will be incremented by 1 and  $Pun_i$  will become 0, and test cases will be repeatedly selected from  $c_i$  until a non-fault-revealing test case is selected from  $c_i$ . If a test case selected from  $c_i$  does not reveal a fault,  $Rew_i$  will become 0 and  $Pun_i$  will be incremented by 1. If  $Pun_i$  reaches a preset bound value  $Bou_i$ ,  $c_i$  will be regarded to have a very low failure rate, and its corresponding  $p_i$  will become 0. The algorithm of RAPT is available in [32].

### 3 ADAPTIVE METAMORPHIC TESTING

In this section, we present the motivation of this paper, describe a framework for process-feedback MT, and algorithms about selection source test cases and MRs.

#### 3.1 Motivation

Since MT was first published, a considerable number of studies have been reported from various aspects [9]. To improve the efficiency of MT, most of studies have paid their attention to identify the better MRs, which are more likely to be violated. For the efficiency of MRs, several factors such as the difference between the source and follow-up test cases [21], [22] and the detecting-faults capacity of MRs compared to existing test oracles [34], have been investigated.

Since the follow-up test cases are generated based on source test cases and MRs, in addition to the so-called good MRs, source test cases also have an impact on the efficiency of MT. However, 57% of existing studies employed RT to select test cases, and 34% of existing studies used existing test suites according to a survey report by Segura et al. [9]. In this study, we investigate the strategies of selection test cases and MRs, and its impacts on the fault-detecting efficiency of MT.

It has been pointed out that fault-detecting inputs tend to cluster into “continuous regions” [35], [36], that is, the test cases in some partitions are more likely to detect faults than the test cases in other partitions. Inspired by the observation, AMT takes full advantage of feedback information to update the test profile, aiming at increasing the selection probabilities of partitions with larger failure rates. Accordingly, the MRs whose sources test cases belonging to the partitions with larger failure rates, are more likely to be selected and violated. Therefore, AMT is expected to detect faults more efficient than traditional MT.

#### 3.2 Framework

Considering the principles of software cybernetics and the features of MT, we propose an adaptive metamorphic testing framework, as illustrated in Figure 1. As shown in Figure 1, there are two feedback loops in the AMT framework. The first feedback loop constitutes the software under test, the database (history of testing data) and the testing strategy, where the history of testing data is used to generate the next test case and MR. The second feedback loop constitutes the software under test, the testing history, the parameter adjustment and the testing strategy, where the history of testing data is used to improve the underlying testing strategy. The improvement may lead the partition with higher failure rate has greater probability of being selected and the partition with lower failure rate has a lower probability of being selected. Interactions between AMT components are depicted in the framework. We next discuss the individual framework components.

1 *Controller*. The controller is the testing strategy, which is responsible for guiding the selection of test case and MR. In this paper, our testing strategy is based on partition testing (PT), which refers to a class of testing techniques that classify the input domain into a number of partitions [37]. After constructing partitions, Testers need to

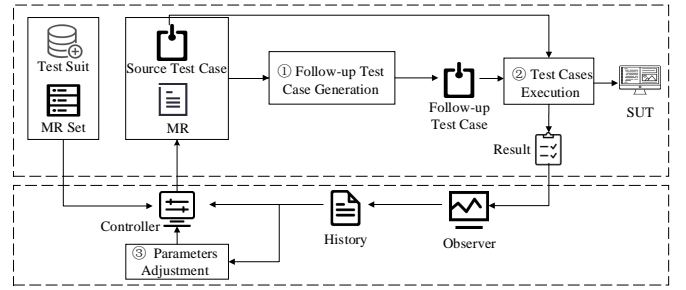


Fig. 1. The Framework of AMT

initialize the test profile, a simple way of doing which would be the use of a uniform probability distribution  $P_1 = P_2 = \dots = P_k$ , where  $k$  denotes the number of partitions, and  $p_i (i = 1, 2, \dots, k)$  denotes the probability of selecting the  $i^{th}$  partition. Besides, the tester needs to create a subset  $R_i (i = 1, 2, \dots, k)$  of MRs for each partition  $s_i$  so that any test case belonging to partition  $p_i$  can generate follow-up test case driven by any MR belonging to  $R_i$ . Then *Controller* selects a partition according to the test profile, and selects source test case and MR based on given source test case selection strategies and MR selection strategies, respectively. The test profile is updated during the testing process according to the testing results.

- 2 *Follow-up Test Case Generation*. The follow-up test case is derived from the selected source test case based on the selected MR.
- 3 *Test Cases Execution*. The relevant AMT component receives the generated source test case and follow-up test case, and executes them on SUT.
- 4 *Parameters Adjustment*. Upon completion of each test, its pass or fail status is determined by comparing the actual and expected results (a pass status if both are the same). The pass or fail status is then used to adjust the values of the probability adjusting factor.

#### 3.3 Source Test Case Selection

The previous MAPT and RAPT algorithms assume that testers can obtain test oracles, therefore, MAPT and RAPT are difficult to apply when there are no test oracles or it is difficult to obtain the test oracles. In this paper, we proposed two test case selection strategies MAPT\* and RAPT\* based on the MAPT and RAPT to select test cases in the context of MT.

##### 3.3.1 MAPT\*

Suppose that source test case  $tc$  and corresponding follow-up test case  $tc'$  are belonging to partition  $s_s$  and  $s_f$ , respectively. if their results violate the related MR,  $\forall i = 1, 2, \dots, m$  and  $s, f \neq i$ , we update test profile using the following equations. When source test case and follow-up test case belong to same partition ( $s = f$ ), we set

$$p'_{s,i} = \begin{cases} p_{s,i} - \frac{\gamma \times p_{s,s}}{m-1} & \text{if } p_{s,i} > \frac{\gamma \times p_{s,s}}{m-1} \\ p_{s,i} & \text{if } p_{s,i} \leq \frac{\gamma \times p_{s,s}}{m-1} \end{cases} \quad (2)$$

and then,

$$p'_{s,s} = 1 - \sum_{i=0, i \neq s}^m p'_{s,i}. \quad (3)$$

Alternatively, if  $s \neq f$ , we set

$$p'_{s,i} = \begin{cases} p_{s,i} - \frac{\gamma \times p_{s,s}}{m-1} & \text{if } p_{s,i} > \frac{\gamma \times p_{s,s}}{m-1} \\ p_{s,i} & \text{if } p_{s,i} \leq \frac{\gamma \times p_{s,s}}{m-1} \end{cases} \quad (4)$$

$$p'_{f,i} = \begin{cases} p_{f,i} - \frac{\gamma \times p_{f,f}}{m-1} & \text{if } p_{f,i} > \frac{\gamma \times p_{f,f}}{m-1} \\ p_{f,i} & \text{if } p_{f,i} \leq \frac{\gamma \times p_{f,f}}{m-1} \end{cases} \quad (5)$$

and then,

$$p'_{s,s} = 1 - \sum_{i=0, i \neq f}^m p_{s,i} \quad (6)$$

$$p'_{f,f} = 1 - \sum_{i=0, i \neq f}^m p_{f,i} \quad (7)$$

When source test case and corresponding follow-up test case donot detect a fault, we employ the following equations to update test profile. If source test case and correspnding follow-up test case belong to same partition ( $s = f$ ), we set

$$p'_{s,i} = \begin{cases} p_{s,i} + \frac{\tau \times p_{s,i}}{m-1} & \text{if } p_{s,s} > \frac{\tau \times (1 - p_{s,s})}{m-1} \\ p_{s,i} & \text{if } p_{s,s} \leq \frac{\tau \times (1 - p_{s,s})}{m-1} \end{cases} \quad (8)$$

and then,

$$p'_{s,s} = \begin{cases} p_{s,s} - \frac{\tau \times (1 - p_{s,s})}{m-1} & \text{if } p_{s,s} > \frac{\tau \times (1 - p_{s,s})}{m-1} \\ p_{s,s} & \text{if } p_{s,s} \leq \frac{\tau \times (1 - p_{s,s})}{m-1} \end{cases} \quad (9)$$

Alternatively, if  $s \neq f$ , we set

$$p'_{s,i} = \begin{cases} p_{s,i} + \frac{\tau \times p_{s,s}}{m-1} & \text{if } p_{s,s} > \frac{\tau \times (1 - p_{s,s})}{m-1} \\ p_{s,i} & \text{if } p_{s,s} \leq \frac{\tau \times (1 - p_{s,s})}{m-1} \end{cases} \quad (10)$$

$$p'_{s,s} = \begin{cases} p_{s,s} - \frac{\tau \times (1 - p_{s,s})}{m-1} & \text{if } p_{s,s} > \frac{\tau \times (1 - p_{s,s})}{m-1} \\ p_{s,s} & \text{if } p_{s,s} \leq \frac{\tau \times (1 - p_{s,s})}{m-1} \end{cases} \quad (11)$$

$$p'_{f,i} = \begin{cases} p_{f,i} + \frac{\tau \times p_{f,f}}{m-1} & \text{if } p_{f,f} > \frac{\tau \times (1 - p_{f,f})}{m-1} \\ p_{f,i} & \text{if } p_{f,f} \leq \frac{\tau \times (1 - p_{f,f})}{m-1} \end{cases} \quad (12)$$

$$p'_{f,f} = \begin{cases} p_{f,f} - \frac{\tau \times (1 - p_{f,f})}{m-1} & \text{if } p_{f,f} > \frac{\tau \times (1 - p_{f,f})}{m-1} \\ p_{f,f} & \text{if } p_{f,f} \leq \frac{\tau \times (1 - p_{f,f})}{m-1} \end{cases} \quad (13)$$

### 3.3.2 RAPT\*

When source test case and follow-up test case belong to same partition ( $s = f$ ), we set,

$$p'_i = \begin{cases} p_i - \frac{(1 + \ln Re_{w_i}) \times \epsilon}{m-1} & \text{if } p_i > \frac{(1 + \ln Re_{w_i}) \times \epsilon}{m-1} \\ 0 & \text{if } p_i \leq \frac{(1 + \ln Re_{w_i}) \times \epsilon}{m-1} \end{cases} \quad (14)$$

and then we have

$$p'_s = 1 - \sum_{i=0, i \neq s}^m p'_i \quad (15)$$

Alternativelt, if  $s \neq f$ , we set

$$p'_i = \begin{cases} p_i - \frac{(1 + \ln Re_{w_i}) \times \epsilon}{m-2} & \text{if } p_i > \frac{(1 + \ln Re_{w_i}) \times \epsilon}{m-2} \\ 0 & \text{if } p_i \leq \frac{(1 + \ln Re_{w_i}) \times \epsilon}{m-2} \end{cases} \quad (16)$$

and then we have

$$p'_s = p_s + \frac{1 - \sum_{i=0, i \neq s, i \neq f}^m p'_i - p_s - p_f}{2} \quad (17)$$

$$p'_f = p_f + \frac{1 - \sum_{i=0, i \neq s, i \neq f}^m p'_i - p_s - p_f}{2} \quad (18)$$

When source test case and corresponding follow-up test case donot detect a fault, we employ the following equations to update test profile. If source test case and correspnding follow-up test case belong to same partition ( $s = f$ ), we set

$$p'_s = \begin{cases} p_s - \delta & \text{if } p_s > \delta \\ 0 & \text{if } p_s \leq \delta \text{ or } Pun_i = Bou_i, \end{cases} \quad (19)$$

and then

$$p'_i = \begin{cases} p_i + \frac{\delta}{m-1} & \text{if } p_s > \delta \\ p_i + \frac{p_s}{m-1} & \text{if } p_s \leq \delta \text{ or } Pun_i = Bou_i, \end{cases} \quad (20)$$

Alternatively, if  $s \neq f$ , we have

$$p'_s = \begin{cases} p_s - \delta & \text{if } p_s > \delta \\ 0 & \text{if } p_s \leq \delta \text{ or } Pun_i = Bou_i, \end{cases} \quad (21)$$

$$p'_f = \begin{cases} p_f - \delta & \text{if } p_f > \delta \\ 0 & \text{if } p_f \leq \delta \text{ or } Pun_i = Bou_i, \end{cases} \quad (22)$$

and then

$$p'_i = p_i + \frac{(p_s - p'_s) + (p_f - p'_f)}{m-2} \quad (23)$$

In AMT, once a source test case is selected, the corresponding candidate MRs are available. We first proposed a simple strategy based on random mechanism to select an MR from the set of candidate MRs, then we proposed an innovative strategy to select an MR that can make the execution of follow-up test case as different as possible from the source test case.



### 3.4 Dynamic MR Selection Strategies

#### 3.4.1 Randomly MR-Selection Strategy (RMRS)

RSMR randomly selects an from the set of candidate MRs, which is a straightforward method without considering extra information.

#### 3.4.2 Properties-Based strategy of MR selection (PBMR)

The fault-detection efficiency of MT highly dependent on the specific MRs that are used, and selecting effective MRs is thus a critical step when applying MT. Chen et al. reported that good metamorphic relations are those that can make the execution of the source-test case as different as possible to its follow-up test case. Moreover, They defined the “difference among execution” as any aspects of program runs (e.g., paths traversed). Before presenting the properties-based strategy of MR selection (PBMR), we first introduce a new metric called category-partition-based metric (CP-distance) that measures the distance between the source test cases and corresponding follow-up test cases. CP-distance makes use of the concepts of categories and choices from the CPM method described in Section 2.2. CP-distance have the capacity to reflect the difference among the source test cases and follow-up test cases, since the categories and choices are defined based on the software functionalities.

More formally, let us denote the set of categories by  $A = \{A_1, A_2, \dots, A_g\}$ , where  $g$  denotes the total number of categories. For each  $A_i$ , its choices are denotes by  $P_i^{A_i} = \{p_1^{A_i}, p_2^{A_i}, \dots, p_h^{A_i}\}$ , where  $h$  denotes the number of choices for  $A_i$ . For input  $x$ , let us denote the corresponding non-empty subset by  $A(x) = \{A_1^x, A_2^x, \dots, A_q^x\}$ , where  $q$  refers to the number of categories associated with  $x$ . Since categories are distinct and their choices are disjoint, input  $x$  in fact consists of values chosen from a non-empty subset of choices, denoted as  $P(x) = \{p_1^x, p_2^x, \dots, p_q^x\}$ , where  $p_i^x (i = 1, 2, \dots, q)$  is the choice of the category  $A_i^x$  for  $x$ .

For any two inputs  $x$  and  $y$ , we define  $DP(x, y)$  as the set that contains elements in either  $P(x)$  or  $P(y)$  but not both. That is,

$$DP(x, y) = (P(x) \cup P(y)) \setminus (P(x) \cap P(y)),$$

where “ $\setminus$ ” is the set difference operator. Now, we define

$$DA(x, y) = \{A_m | A_m \text{ if } p_j^i \in DP(x, y)\}.$$

In other words,  $DA(x, y)$  is the set of categories in which inputs  $x$  and  $y$  have different choices. Then, the distance measure between  $x$  and  $y$  is defined as  $|DA(x, y)|$  (the size of  $DA(x, y)$ ); that is, the number of categories that appear in either  $x$  or  $y$  but not both, or in which the choices in  $x$  and  $y$  differ.

obviously, a greater value of CP-distance representing more dissimilar execution. After selecting a source test case  $stc_i$  and obtaining a set of candidate MRs  $\mathcal{R} = \{r_1^{stc_i}, r_2^{stc_i}, \dots, r_g^{stc_i}\}$  ( $g$  is the number of MRs whose source test case could be  $stc_i$ ), PBMR generate a set of candidate follow-up test cases  $FC = \{ftc_{r_1^{stc_i}}, ftc_{r_2^{stc_i}}, \dots, ftc_{r_g^{stc_i}}\}$  according to every MR belonged to  $FC$ . Then, the distance  $CP_{i,h}$  ( $h \in \{1, 2, \dots, g\}$ )

between  $stc_i$  and each follow-up test case  $ftc_{r_h^{stc_i}}$  is calculated. Finally, the MR  $r_h^{stc_i}$  is selected as long as the following condition is hold:

$$CP_{i,h} = \max\{CP_{i,1}, CP_{i,2}, \dots, CP_{i,g}\}.$$

The details of PBMR is given in Algorithm 1.

---

#### Algorithm 1 PBMR

---

**Input:**  $stc_i, \mathcal{R} = \{r_1^{stc_i}, r_2^{stc_i}, \dots, r_g^{stc_i}\}, CArray = null$

**Output:**  $CP_{max}$  ( $max \in \{1, 2, \dots, g\}$ )

---

- 1: **for**  $h = 1 \rightarrow h = g$  **do**
  - 2:   Generate the follow-up test case  $ftc_h$  according to the  $stc_i$  and  $r_h^{stc_i}$
  - 3:   Calculate the distance  $CP_{i,h}$  between the  $stc_i$  and  $ftc_h$
  - 4:   Add  $CP_{i,h}$  to the list  $CArray$
  - 5: **end for**
  - 6: Calculate the maximal value  $CP_{max}$  ( $max \in \{1, 2, \dots, g\}$ ) in the list  $CArray$
- 

## 4 EMPIRICAL STUDY

We conducted a series of empirical studies to evaluate the performance of AMT.

### 4.1 Research Questions

In our experiments, we focused on addressing the following two research questions:

**RQ1 Could the proposed AMT strategy detect failures faster than MT?**

Fault-detection efficiency is a key criterion for evaluating the performance of a testing technique. In our study, we chose four Laboratory programs, a GNU program, and an Alibaba program to evaluate the fault-detecting efficiency.

**RQ2 What is the actual test case selection overhead when using the AMT technique?**

We evaluate the test case selection overhead of M-AMT and compare with traditional MT in detecting software faults.

### 4.2 Object Programs

In order to evaluate the fault-detection effectiveness of proposed methods in different scales, different implementation languages, and different fields, we chose to study three sets of object programs: four laboratorial programs that were developed according to corresponding specifications<sup>1</sup>, the regular expression processor component of the larger utility program GNU `grep`, and a Java library developed by Alibaba<sup>2</sup>, which can be used to convert Java Objects into their JSON representation, and convert a JSON string to an equivalent Java object. Table 1 summarized the basic information of the object programs, giving the developer of object programs (*Developer*), the name of object program

1. The implementations have been made available at: <https://github.com/phantomDai/subjects4tsc.git>

2. This project is available at <https://github.com/alibaba/fastjson>

(*Program*), the size of the object program (*LOC*), the number of total faults for each object program (*Number of All Faults*), the number of used faults in each object program (*Number of Used Faults*), test cases in the associated test pool (*Size of test suite*), the number of identified MRs for each object program (*Number of MRs*), and the number of partition for each object program (*Number of Partitions*).

Based on four real-life specifications, We developed four systems, which were implemented using Java programming language. A detailed description of each laboratory program is given in the following.

- 1 Expense Reimbursement System (ERS) provides an interface through which customers can know how much they need to pay according to plans, month charge, calls, and data usage.
- 2 Aviation Consignment Management System (ACMS) aims to help airline companies check the allowance (weight) of free baggage, and the cost of additional baggage. Based on the destination, flights are categorized as either domestic or international. For international flights, the baggage allowance is greater if the passenger is a student (30kg), otherwise it is 20kg. Each aircraft offers three cabins classes from which to choose (economy, business, and first), with passengers in different classes having different allowances.
- 3 Expense Reimbursement System (ERS) assists the sales Supervisor of a company with handling the following tasks: i) Calculating the cost of the employee who use the cars based on their titles and the number of miles actually traveled; ii) accepting the requests of reimbursement that include airfare, hotel accommodation, food and cell-phone expenses of the employee.
- 4 Meal Ordering System (MOS) helps catering provider to determine the quantity for every type of meal and other special requests (if any) that need to be prepared and loaded onto the aircraft served by this provider. For each flight, MOS generates a Meal Schedule Report (MSR), containing the number of various types of meals (first class, business class, economy class, vegetarian, child, crew member, and pilot) and the number of bundles of flowers.

The used version of the `grep` programs is 2.5.1a [39]. This program searches one or more input files for lines containing a match to a specified pattern. By default, `grep` prints the matching lines. We chose `grep` for our study for several reasons:

- 1 `grep` program is wide used in Unix system, providing a opportunity to demonstrate the real world relevance of our techniques.
- 2 `grep` program, and its input format, are of greater complexity than the programs in the other test sets, but still manageable as a target for automated test case generation.
- 3 `grep` program has been studied by several papers [20], [40], [41], that means, we can obtain its faults, source test cases and MRs freely and conveniently.

The inputs of the `grep` were categorized into three components: options, which consist of a list of commands to modify the searching process, pattern, which is the regular expression to be searched for, and files, which refers to

the input files to be searched. The scope of functionality of this program is larger, which leads to construct test infrastructure to test all of functionality would have been impractical. Therefore, we restricted our focus to the regular expression analyzer of the `grep`.

`FastJson` is a Java library that can be used to convert Java Objects into their JSON representation (known as serialization). It can also be used to convert a JSON string to an equivalent Java object (known as deserialization). We chose `FastJson` for our study for several reason:

- 1 `FastJson` is wide used in real-world projects, providing a opportunity to demonstrate the real world relevance of our techniques.
- 2 `FastJson` is more complex than other programs, and it is an open source project, that is, we can obtain its faults freely and conveniently.

The scope of functionality of this program is larger, which leads to construct test infrastructure to test all of functionality would have been impractical. Therefore, we restricted our focus to the deserialization of the `FastJson`. Two bigger version of `FastJson` are available, named `FastJson 1.1.*` (1.1.0 – 1.1.51) and `FastJson 1.2.*` (1.2.1 – 1.2.61), respectively. Since there are no failure reports for `FastJson 1.1.*`, we collected deserialization-related faults from the `FastJson 1.2.1–FastJson 1.2.48` (When we started this project, the latest version of `FastJson` was 1.2.48). We removed the faults that caused `FastJson` to throw an error information during execution.

There three sets of programs present complementary strengths and weaknesses as experiment objects. The Laboratory programs implement simple functions and their interfaces are easy to understand. The test engineers can easily generate source test cases, and identify MRs for testing laboratory programs. However, these programs are small and there are a limited number of faulty versions available for each programs. The `grep` program is a much larger system for which mutation faults could be generated. The `FastJson` is also a much larger system, and the faults of that are obtained on GitHub. We provide further details on each of those sets of object programs next.

### 4.3 Faults

We employed the tool `MuJava` [41] to conduct mutation analysis [42] for four laboratory programs, generating a total of 792 mutants. Each mutant was obtained by changing the original programs syntax (using one of all applicable mutation operators provided by `MuJava`). Equivalent mutants, and those that were too easily detected (requiring less than 20 randomly generated test cases), were removed. To ensure the statistical reliability, we obtained 30 different test suites using different random seeds, then tested all mutants with all test suites, calculating the average number of test cases needed to kill (detect) a mutant.

Barus et al. [40] have constructed faults for `grep`, including one real fault and 19 artificial faults (generated by `MuJava`). These mutants have failure rate between 0.001 to 0.1, corresponding to F-measures of between 10 and 1000 for RT. We removed easily detected faults (requiring less than

TABLE 1  
Twelve Programs as Experimental Objects

Developer	Program	Language	LOC	All Faults	Used Faults	Number of Test Cases	Number of MRs	Number of Partitions
Laboratory	CUBS	Java	107	187	3	—	184	8
	ACMS	Java	128	210	3	—	142	4
	ERS	Java	117	180	1	—	1130	12
	MOS	Java	135	215	1	—	3512	9
GNU	grep	C	10,068	20	8	101,193	12	550
Alibaba	FastJson_v31	Java	125,192	1	1	16,383	17	128
	FastJson_v36	Java	134,440	1	1	16,383	17	128
	FastJson_v40	Java	144,044	1	1	16,383	17	128
	FastJson_v48	Java	149,544	1	1	16,383	17	128

TABLE 2  
The Details of Independent Variables

Techniques	Test case selection strategy	MRs selection strategy
MR-AMT	MAPT*	RT
MP-AMT	MAPT*	PBMR
RR-AMT	RAPT*	RT
RP-AMT	RAPT*	PBMR
MT	RT	RT

20 randomly generated test cases), and finally obtained 8 faults.

All Fastjson fault information are available at <https://github.com/alibaba/fastjson>. We first found all the faults related to deserialization of `FastJson`, then removed the faults that caused the program to crash and error, and finally got four real, hard-to-detected faults.

## 4.4 Variables

### 4.4.1 Independent Variables

The independent variable in our experiment is the different strategies that improve the fault-detection efficiency. As choices for this variable, we include, of course, proposed adaptive metamorphic testing technique, which includes two test cases selection strategies (MAPT\* and RAPT\*) and two MRs selection strategies (RBMR and PBMR). As baseline techniques for use in comparison, we selected traditional MT, which randomly selects source test case and MRs.

Table 2 summarizes the details information of selecting source test cases and MRs of different techniques. Traditional MT is a natural baseline choice, because all proposed techniques is designed as an enhancement to MT, and assessing whether proposed techniques including MR-AMT (which uses MAPT\* and RT strategies to select source test case and MR, respectively), MP-AMT (which uses MAPT\* and PBMR strategies to select source test case and MR, respectively), RR-AMT (which uses RAPT\* and RT strategies to select source test case and MR, respectively), and RP-AMT (which uses RAPT\* and PBMR strategies to select source test case and MR, respectively) is more cost-effective than MT is important.

### 4.4.2 Dependent Variables

The choice of a metric to use in comparing the effectiveness of testing techniques is non-trivial.

The dependent variable for RQ1 is the metric for evaluating the fault-detection effectiveness. Several effectiveness metrics exist, including: the P-measure [42] (the probability of at least one fault being detected by a test suite); the E-measure [43] (the expected number of faults detected by a test suite); the F-measure [32] (the expected number of test case executions required to detect the first fault); and the T-measure [44] (the expected number of test cases required to detect all faults). Since the F-measures has been widely used for evaluating the fault-detection efficiency and effectiveness of MT, and APT testing techniques [20], [32], [41], they are also adopted in this study. In addition, as shown in Fig 1, the testing process may not terminate after the detection of the first fault. Furthermore, because the fault detection information can lead to different probability profile adjustment mechanisms, it is also important to see what would happen after revealing the first fault. Therefore, we introduce the F2-measure [32], [45] as the number of additional test cases required to reveal the second fault after detection of the first fault. We use  $F$  and  $F2$  to represent the F-measure and the F2-measure of a testing technique.

We define  $F$ -count as the number of test cases needed to detect a failure in a specific test run,  $F2$ -count as the number of additional test cases needed to detect the second failure after detecting the first failure. The F-measure is the expected  $F$ -count for a testing method:

$$F = \overline{F - count}. \quad (24)$$

The F2-measure is the expected  $F2 - count$  for a testing method:

$$F = \overline{F2 - count}. \quad (25)$$

The F-measure is particularly appropriate for measuring the failure-detection efficiency of MT. The value of F-measure indicates the speed of detecting failures for used testing techniques. The F2-measure can be further generalized to measure the number of test cases required for detecting the  $i + 1$ th fault in the context of having the  $i$ th fault being detected in an iterative way. The value of F2-measure reflects the AMT's ability to detect failures with test history information.

Holm-Bonferroni method [46] to determine which pair of testing techniques have significant difference in terms



of  $F$ . Across the whole study, for each pair of testing techniques, denoted by technique a and technique b, the null hypothesis ( $H_0$ ) was that a and b had similar performance in terms of one metric; whereas the alternative hypothesis ( $H_1$ ) was that a and b had significantly different performance in terms of  $F$ . All the null hypotheses were ordered by their corresponding p-values, from lowest to largest; in other words, for null hypotheses  $H_0^1, H_0^2, \dots, H_0^{h-1}$ , we had  $p_1 \leq p_2 \leq \dots$ . For the given confidence level  $\alpha = 0.05$ , we found the minimal index  $h$  such that  $p_h > \frac{\alpha}{N+1-h}$  (where  $N$  is the total number of null hypotheses). Then, we rejected  $H_0^1, H_0^2, \dots, H_0^{h-1}$ , that is, we regarded the pair of techniques involved in each of these hypotheses to have statistically significant difference in terms of a certain metric. On the other hand, we considered the pair of techniques involved in each of  $H_0^h, H_0^{h+1}, \dots$  not to have significant difference with respect to one metric, as these hypotheses were accepted.

An obvious metric for RQ2 is the time required to detect faults. Corresponding to the F-measure and F2-measure, in this study we used  $F$ -time and  $F2$ -time denote the time required (i.e. test cases generation, test cases selection, and test cases execution) to detect the first fault. Obviously, a smaller values of  $F$ -time indicate a better performance.

## 4.5 Experimental Settings

### 4.5.1 Generation of Categories and Choices for Object Programs

The categories and choices used for the laboratorial programs and Real-World Popular Programs (FastJson) considered in this study were designed by the authors. On the other hand, the categories and choices used for GNU Program was designed by Barus [41]. In large part, the selection of appropriate categories and choices is at a tester's discretion; we chose what we regarded as simple approaches for emulating that process. Precise details on the categories and choices used in our study are provided as following.

For each laboratorial program (ACMS, CUMS, ERS, MOS), categories and choices for its input domain was identified in advance based on its corresponding specification. Based on the identified categories and choices of each laboratorial program, we divided input domain into partitions, generated test cases, and identified MRs. The identified categories and corresponding choices for ACMS, CUMS, ERS, and MOS are shown in Table 3, 4, 5, and 6, respectively. Detailed description of categories and choices for laboratorial programs are available at <https://github.com/phantomDai/AMT-material.git>.

For `grep`, Barus et al. [40] have made use the existing testing information about `grep` recorded in the SIR [47], conducting categories and choices. Furthermore, they divided the categories into two groups — the independent and the dependent categories (The details of the dependent and independent categories are listed in Table 7). The independent group includes all categories that contain elements which can form a valid test case on their own. Dependent categories need the presence of elements of other categories to form valid test cases.

To define the categories and choices of `FastJson`, the specification and user documentation were first examined,

TABLE 3  
Definition of Categories and Choices for ACMS

#	Categories	Associated Choices
1	class	First calss (1a) Business class (1b) Economy class (1c) Infant (1d)
2	region	Domestic flights (2a) International flights (2b)
3	isStudent	True (3a) False (3b)
4	luggage	Luggage < Free checked-in (4a) Luggage ≤ Free checked-in (4b)
5	fee	Fee = 0 (5a) Fee ≤ 0 (5b)

TABLE 4  
Definition of Categories and Choices for CUMS

#	Categories	Associated Choices
1	plan	A (1a) B (1b)
2	options	46CNY (2a) 96CNY (2b) 286CNY (2c) 886CNY (2d) 126CNY (2e) 186CNY (2f)
3	calls	calls < free calls (3a) calls ≥ free calls (3b)
4	data	data < free data (4a) data ≥ free data (4b)

however they were either not discovered or did not provide significant detail to construct appropriate categories and choices. Therefore, we also made use the existing best practices, the set of test cases and source codes about `FastJson`

TABLE 5  
Definition of Categories and Choices for ERS

#	Categories	Associated Choices
1	title	senior sales manager (1a) sales manager (1b) sales executive (1c)
2	mileage	$0 \leq \text{mileage} \leq 3000$ (2a) $3000 \leq \text{mileage} \leq 4000$ (2b) $\text{mileage} \geq 4000$ (2c)
3	sales	$0 \leq \text{sales} \leq 50,000$ (3a) $50,000 \leq \text{sales} \leq 80,000$ (3b) $0 \leq \text{sales} \leq 3000$ (3c) $80,000 \leq \text{sales} \leq 100,000$ (3d) $\text{sales} \geq 100,000$ (3e)
4	airline ticket(AT)	$AT = 0$ (4a) $AT \leq 0$ (4b)
5	other expenses(OE)	$OE = 0$ (5a) $OE \leq 0$ (5b)

TABLE 6  
Definition of Categories and Choices for *MOS*

#	Categories	Associated Choices
1	Aircraft Model	747200 (1a)
		747300 (1b)
		747400 (1c)
		002000 (1d)
		003000 (1e)
2	Change in the Number of Crews	Yes (2a)
		No (2b)
3	Number of Crews( <i>NC</i> )	<i>NC</i> > default value (3a)
		<i>NC</i> = default value (3b)
		<i>NC</i> < default value (3c)
4	Change in the Number of pilots	Yes (4a)
		No (4b)
5	Number of pilots( <i>NP</i> )	<i>NP</i> > default value (5a)
		<i>NP</i> = default value (5b)
		<i>NP</i> < default value (5c)
6	Number of Child Passengers	<i>childNum</i> > 0 (6a)
		<i>childNum</i> = 0 (6b)
7	Number of Bundles Flowers	<i>NF</i> > 0 (7a)
		<i>NF</i> = 0 (7b)

TABLE 7  
Table of Independent and Dependent categories for *grep*

Independent Categories	Dependent Categories
NormalChar	Bracket
WordSymbol	Iteration
DigitSymbol	Parentheses
SpaceSymbol	Line
NamedSymbol	Word
AnyChar	Edge
Range	Combine

recorded in the Alibaba repository<sup>3</sup>. As mentioned, we restrict our experimentation to testing the deserialization of *FastJson*. As a consequence, our categories-choices scheme considers only on this portion of *FastJson*. Deserialization is the process of translating a stored format (For example, a Json file or memory buffer) into an object state, which contains the values of all the variables in a program. Therefore, We consider the base member variables and their combinations. All identified categories and choices of *FastJson* are listed in Table 8. The Detailed description of categories and choices for *FastJson* are available at our repository: <https://github.com/phantomDai/AMT-material.git>.

#### 4.5.2 Generation of Test Cases for Object Programs

For *CUBS*, *ACMS*, and *ERS*, we used a generator that is based on the complete test frame devised for MT-related techniques selection. We systematically generated a test case for each complete test frame, which were collectively guaranteed to cover each category and choice. The final test suites contained 284, 1470, and 2260, respectively. All complete test frame are available at: <https://github.com/phantomDai/AMT-material.git>.

3. The best practices and test suite are available at: <https://github.com/alibaba/fastjson/wiki>

TABLE 8  
Definition of Categories and Choices for *FastJson*

#	Categories	Associated Choices
1	Float	Exist
		Inexistence
2	Double	Exist
		Inexistence
3	Short	Exist
		Inexistence
4	Byte	Exist
		Inexistence
5	Int	Exist
		Inexistence
6	Long	Exist
		Inexistence
7	Boolean	Exist
		Inexistence
8	Char	Exist
		Inexistence
9	Date	Exist
		Inexistence
10	String	Exist
		Inexistence
11	Enum	Exist
		Inexistence
12	Map	Exist
		Inexistence
13	List	Exist
		Inexistence
14	Set	Exist
		Inexistence

For *grep*, we used the test suite generated by Barus, who has described the test case generation process in detail [41]. The original test suite generated by Barus contains 171,634 inputs. During the test process, we found that some of the test frames contain the same choices, but in a different order. In practice, each test frame corresponds to a partition in which test cases cover some paths, thus, those test frames contain the same choices cover same or similar paths. we removed those test frames, and obtained 101,193 test frames, which are available at our repository: <https://github.com/phantomDai/AMT-material.git>.

As for *FastJson*, we used a generator that is based on the complete test frame devised for MT-related techniques selection. We systematically generated a test case for each complete test frame, which were collectively guaranteed to cover each category and choice. The final test suites contained 16383 Java object and 16383 corresponding Json files. All complete test frames are available at: <https://github.com/phantomDai/AMT-material.git>.

#### 4.5.3 Partitioning

The proposed technique AMT based on the partition testing that is a mainstream family of software testing techniques, which can be realized in different ways, such as Intuitive Similarity, Equivalent Paths, Risk-Based, and Specified As Equivalent (two test values are equivalent if the specification says that the program handles them in the same way). In this

TABLE 9  
Categories and choices of ACMS

Category	Associated choices
aircraft cabin	<i>Economy, Business, First Class</i>
flight region	<i>International, Domestic</i>
baggage weight	<i>Above Limit, Below Limit</i>

TABLE 10  
Partitions of ACMS

Partition	Complete Test Frame
$c_1$	$\{International, Economy, weight_*\}$
$c_2$	$\{International, Business, weight_*\}$
$c_3$	$\{International, FirstClass, weight_*\}$
$c_4$	$\{Domestic, Economy, weight_*\}$
$c_5$	$\{Domestic, Business, weight_*\}$
$c_6$	$\{Domestic, FirstClass, weight_*\}$

study, we made use of the CPM to conduct the partitioning. Our previous work [32] found that there is not a strong correlation between the granularity level in partitioning and the performance of APT.

We randomly select two categories for ACMS, CUBS, ERS, and MOS, and partition their input domain according to the combinations of choices. We use an explanatory example to illustrate the method of partitioning. According to the description of Section 4.2, we identify categories and associated choices of ACMS, as shown in Table 9. With these categories/choices and constraints among choices, we further derive a set of complete test frames and each of them corresponds to a partition. Note that the number of partitions may vary with the granularity level. To ease the illustration, we derive partitions without considering the baggage weight category. Table 10 shows the resulting partitions for ACMS, where  $weight_*$  indicates no classification on the baggage weight.

Barus has divided the categories into two groups — the independent and the dependent categories [40]. The independent group includes all categories that contain elements which can form a valid test case on their own. Dependent categories need the presence of elements of other categories to form valid test cases. We obtained two partitioning schemes by choosing independent categories to ignore dependent categories and choosing dependent categories to ignore independent categories. As a consequence we constructed 550 and 3380 partitions, respectively. We finally selected the first partition scheme (550 partitions), since the second partition scheme tends to result in the repeated selection of the same test case during testing.

#### 4.5.4 Initial Test Profile

The proposed techniques M-AMT and R-AMT make use of feedback information to adjust the test profile. Then the updated test profile guide those techniques to select next source test case and MR. Before executing the test, a concrete test profile should be initialized. Because test cases may be generated randomly during the test process, a feasible method is to use a uniform probability distribution as the initial testing profile. On the other hand, in our previous work [32], we have compared the equal and proportional

TABLE 11  
Initial Test Profile for Subject Program

Program	Number of partitions	Initial test profile
ACMS	8	$\{ \langle s_1, \frac{1}{8} \rangle, \langle s_2, \frac{1}{8} \rangle, \dots, \langle s_8, \frac{1}{8} \rangle \}$
CUBS	4	$\{ \langle s_1, \frac{1}{4} \rangle, \langle s_2, \frac{1}{4} \rangle, \dots, \langle s_4, \frac{1}{4} \rangle \}$
ERS	12	$\{ \langle s_1, \frac{1}{12} \rangle, \langle s_2, \frac{1}{12} \rangle, \dots, \langle s_{12}, \frac{1}{12} \rangle \}$
MOS	10	$\{ \langle s_1, \frac{1}{10} \rangle, \langle s_2, \frac{1}{10} \rangle, \dots, \langle s_{10}, \frac{1}{10} \rangle \}$
grep	550	$\{ \langle s_1, \frac{1}{550} \rangle, \langle s_2, \frac{1}{550} \rangle, \dots, \langle s_{550}, \frac{1}{550} \rangle \}$
FastJson	128	$\{ \langle s_1, \frac{1}{128} \rangle, \langle s_2, \frac{1}{128} \rangle, \dots, \langle s_{550}, \frac{1}{128} \rangle \}$

TABLE 12  
Theoretical Values of  $\delta$

Program	$\theta_M$	$\theta_\Delta$	$\delta$
ACMS	2.86E-1	1.04E-1	1.72E-2
CUBS	2.50E-1	2.78E-2	2.60E-2
ERS	6.00E-1	2.93E-1	1.98E-2
MOS	1.59E-1	1.57E-1	9.42E-3
grep	9.41E-1	8.82E-1	4.20E-1
FastJson	5.80E-2	5.33E-2	2.77E-3

initial test profile in terms of associated performance metrics (F-measure, F2-measure [32], and T-measure). The comparison results that there was no significant difference between these two types of initial test profiles. In our experiment, we used a uniform probability distribution for the initial test profile. The initial test profiles of each subject program is summarized in Table 11, where  $\langle s_i, p_i \rangle$  means that the probability of selecting partition  $s_i$  is  $p_i$ .

#### 4.5.5 Constants

Previous studies [32], [48]–[50] have given some guidelines on how to set  $\varepsilon$  and  $\delta$  of RAPT. We followed these studies to set  $\varepsilon = 0.05$  and calculate the proper  $\delta$  according to the following formula (which is extracted from [50]):

$$\frac{1}{\theta_M} - 1 < \frac{\varepsilon}{\delta} < \frac{1}{\theta_\Delta} - 1, \quad (26)$$

where  $\theta_M$  and  $\theta_\Delta$  are the largest and the second largest failure rates amongst all partitions, respectively. Note that the value of  $\theta_\Delta$  was different for different versions of FastJson. According to Formula 26, the theoretically values of  $\delta$  in each program are shown in Table 12.

To set the values of  $\gamma$  and  $\tau$  for MAPT, we followed our previous study [32], in which we conducted a series of preliminary experiments. We observed that  $\gamma$  and  $\tau$  could neither be too large nor too small, which imply that the change in probability profile would be either very dramatic or very marginal. Both situations might result in the “un-fair” (either too big or too small) award/punishment to certain partitions. After several rounds of trials, we concluded that  $\gamma = \tau = 0.1$  were fair settings. All faults in our experiments are non-trivial and thus not easy to be killed. Thus, the value of Boui could not be too small. We set  $Boui = 70\% \times k_i$ , where  $k_i$  is the number of test cases selected from  $c_i$ .

## 4.6 Identification of MRs for Object Programs

### 4.6.1 The MRs of Laboratorial Programs

In our study, to obtain metamorphic relations of ACMS, CUBS, ERS, and MOS, we made use of METRIC [24], which is based on the category choice framework [33]. In METRIC, only two distinct complete test frames that are abstract test cases defining possible combinations of inputs, are considered by the tester to generate an MR. In general, METRIC has the following steps to identify MRs:

- Step1. Users select two relevant and distinct complete test frames as a *candidate pair*.
- Step2. Users determine whether or not the selected candidate pair is useful for identifying an MR, and if it is, then provide the corresponding MR description.
- Step3. Restart from step 1, and repeat until all candidate pairs are exhausted, or the predefined number of MRs to be generated is reached.

Following the above guidelines, we identified MRs for CUBS, ACMS, and ERS. Table 13 presents a part of MRs of those programs. In the table,  $O_1$  and  $O_2$  denote the output of source test case and follow-up test case, respectively. *NMC.0* indicates the number of crew meals have not changed; *NMP.0* indicates the number of passenger meals have not changed; *NMP.1* indicates the number of passenger meals have increased; *NMCP.0* indicates the number of child meals have not changed; *NMCP.1* indicates the number of child meals have increased; *NMCP. -1* indicates the number of passenger meals have decreased; *NF.0* indicates the number of flowers have not changed; and *NF.1* indicates the number of flowers have increased.

The development of the MRs for `grep` are restricted to the test cases of `grep`. Specifically, these MRs are only associated with the regular expression parameter of `grep`'s input. The details of MRs are provided in [51].

Based on the generated test cases, best practice, and all faults reported in Alibaba repository, we conducted MRs for `FastJson`. Accordingly, we identified twenty MRs particularly for testing the deserialization of `FastJson`. The details of MRs are available at: <https://github.com/phantomDai/AMT-material.git>.

## 4.7 Experimental Environment

Our experiments were conducted on a virtual machine running the Ubuntu 18.04 64-bit operating system. In this system, there were two CPUs and a memory of 4GB. The test scripts were generated using Java, bash shell, and Python. In the experiments, we repeatedly ran the testing using each technique 30 times [52] with different random seeds to guarantee the statistically reliable mean values of the metrics.

## 5 EXPERIMENTAL RESULTS

### 5.1 RQ1: Fault-detection efficiency

F-, and F2-results for ACMS, CUBS, PBS, MOS, `grep`, and `FastJson` are shown using boxplots in Figures 2 and 3, where  $AMT_1$ ,  $AMT_2$ ,  $AMT_3$ , and  $AMT_4$  denotes  $MR - AMT$ ,  $MP - AMT$ , and  $RR - AMT$ , and  $RP - AMT$ ,

respectively. In each boxplot, the upper and lower bounds of the box represent the third and first quartiles of the metric, respectively; the middle line represents the median value; the upper and lower whiskers mark, respectively, the largest and smallest data within the range of  $\pm 1.5 \times IQR$  (where  $IQR$  is the interquartile range); outliers beyond the  $IQR$  are denoted with hollow circles; and each solid circle represents the mean value of the metric.

From Figures 2 and 3, we can observe that in general, RP-AMT was the best performer in terms of F-measure and F2-measure, followed by MP-AMT, MR-AMT, RR-AMT, and MT in descending order. We further conducted hypothesis testing to verify the statistical significance of this observation. We used the Holm-Bonferroni method [53] to determine which pair of testing techniques have significant difference in terms of each metric. Across the whole study, for each pair of testing techniques, denoted by technique  $a$  and technique  $b$ , the null hypothesis ( $H_0$ ) was that  $a$  and  $b$  had the similar performance in terms of one metric. All the null hypotheses were ordered by their corresponding p-values, from lowest to largest; in other words, for null hypotheses  $H_0^1, H_0^2, \dots$ , we had  $p_1 \leq p_2 \leq \dots$ . For the given confidence level  $\alpha = 0.05$ , we found the minimal index  $h$  such that  $p_h > \frac{\alpha}{N+1-h}$  (where  $N$  is the total number of null hypotheses). Then, we rejected  $H_0^1, H_0^2, \dots, H_0^{h-1}$ , that is, we regarded the pair of techniques involved in each of these hypotheses to have statistically significant difference in terms of a certain metric. On the other hand, we considered the pair of techniques involves in each of  $H_0^h, H_0^{h+1}, \dots$  not to have significant difference with respect to one metric, as these hypotheses were accepted. The statistical testing results are shown in Tables 14 and 15. Each cell in the tables gives the number of scenarios where the technique on the top row performed better than that on the left column. If the difference is significant, the number will be displayed in bold. For example, the "131" in the top right corner in Table 14 indicates that out of 180 scenarios (6 programs repeatedly ran), RP-AMT had smaller F-measure than MT for 131 scenarios. Correspondingly, the "49" in the bottom left corner in Table 14 means that the F-measure of MT was smaller than that of RP-AMT for only 49 scenarios. The corresponding null hypothesis was that MT and RP-AMT had similar performance in terms of F-measure. Since this hypothesis was rejected, we could say that the fault-detection capabilities of RAPT and RPT were significantly different in terms of F-measure. This statistically significance difference is represented by bold font of "131" and "49", which further indicates that RAPT was significantly better than RPT.

### 5.2 RQ2: Selection Overhead

The detailed experimental results of F-time and F2-time are shown in Figures 4 and 5. From Figures 4 and 5, we can observe that in general, RP-AMT had the best performance, followed by RR-AMT, MR-AMT, and MP-AMT. Note that MR-AMT, MP-AMT, RR-AMT, and RP-AMT do not consistently delivered the best performance for all object programs in terms of both F-time and F2-time, while the larger the program, the better the AMT performs. As we can observe from figure 2, the execution of the test cases accounts for

TABLE 13  
Part of MRs of Laboratorial Programs

Program	MRs		
	Source Test Cases	Follow-up Test Cases	Relations
CUBS	{1a, 2a, 3a, 4a}	{1a, 2a, 3a, 4b}	$O_1 \geq O_2$
	{1a, 2a, 3a, 4a}	{1a, 2a, 3b, 4a}	$O_1 \geq O_2$
	{1a, 2a, 3a, 4a}	{1a, 2a, 3b, 4b}	$O_1 \geq O_2$
	{1a, 2a, 3a, 4a}	{1a, 2b, 3a, 4a}	$O_1 \geq O_2$
	{1a, 2a, 3a, 4a}	{1a, 2b, 3a, 4b}	$O_1 \geq O_2$
ACMS	{1a, 2a, 3a, 4a, 5a}	{1a, 2a, 3a, 4a, 5b}	$O_1 = O_2 = 0$
	{1a, 2a, 3a, 4a, 5a}	{1a, 2a, 3a, 4b, 5a}	$O_1 = O_2 = 0$
	{1a, 2a, 3a, 4a, 5a}	{1a, 2b, 3a, 4a, 5a}	$O_1 = O_2 = 0$
	{1a, 2a, 3a, 4a, 5a}	{1b, 2a, 3a, 4a, 5a}	$O_1 = O_2 = 0$
	{1a, 2a, 3a, 4a, 5a}	{1c, 2a, 3a, 4a, 5a}	$O_1 = O_2 = 0$
CUBS	{1a, 2a, 3a, 4a}	{1a, 2a, 3a, 4b}	$O_1 \leq O_2$
	{1a, 2a, 3a, 4a}	{1a, 2a, 3b, 4a}	$O_1 = O_2 = 0$
	{1a, 2a, 3a, 4a}	{1a, 2a, 3c, 4a}	$O_1 = O_2 = 0$
	{1a, 2a, 3a, 4a}	{1a, 2b, 3a, 4a}	$O_1 = O_2 = 0$
	{1a, 2a, 3a, 4a}	{1a, 2c, 3a, 4a}	$O_1 \geq O_2$
MOS	{1a, 2b, 3b, 4b, 5b, 6b, 7b}	{1a, 2b, 3b, 4b, 5b, 6b, 7a}	$NMC.0, NMP.0, NMCP.0, NF.1$
	{1a, 2b, 3b, 4b, 5b, 6b, 7b}	{1a, 2b, 3b, 4b, 5b, 6a, 7b}	$NMC.0, NMP.0, NMCP.1, NF.0$
	{1a, 2b, 3b, 4b, 5b, 6b, 7b}	{1a, 2b, 3b, 4b, 5b, 6a, 7a}	$NMC.0, NMP.0, NMCP.1, NF.1$
	{1a, 2b, 3b, 4b, 5b, 6b, 7b}	{1a, 2b, 3b, 4a, 5c, 6b, 7b}	$NMC.0, NMP.0, NMCP. - 1, NF.0$
	{1a, 2b, 3b, 4b, 5b, 6b, 7b}	{1a, 2b, 3b, 4a, 5a, 6b, 7b}	$NMC.0, NMP.1, NMCP.0, NF.0$

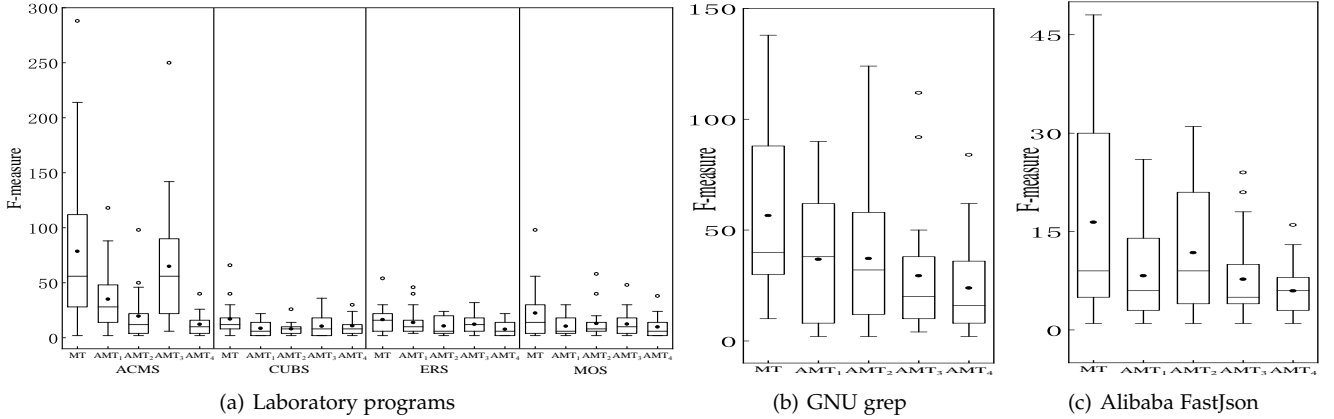


Fig. 2. F-measure boxplots for each program

TABLE 14

Number of Scenarios where the Technique on the top row has a lower metric (F-measure) score than the technique on the left column

	MT	MR-AMT	MP-AMT	RR-AMT	RP-AMT
MT	—	<b>124</b>	<b>128</b>	<b>119</b>	<b>131</b>
MR-AMT	56	—	105	86	106
MP-AMT	52	75	—	92	102
RR-AMT	61	94	88	—	110
RP-AMT	49	74	78	70	—

TABLE 15

Number of Scenarios where the Technique on the top row has a lower metric (F2-measure) score than the technique on the left column

	MT	MR-AMT	MP-AMT	RR-AMT	RP-AMT
MT	—	<b>142</b>	<b>148</b>	<b>136</b>	<b>156</b>
MR-AMT	38	—	135	118	145
MP-AMT	32	45	—	121	139
RR-AMT	44	62	59	—	143
RP-AMT	24	35	41	37	—

most of the overhead. AMT compensates for the overhead of selecting test cases by reducing the number of test cases

executed.

We also conducted hypothesis testing to determine

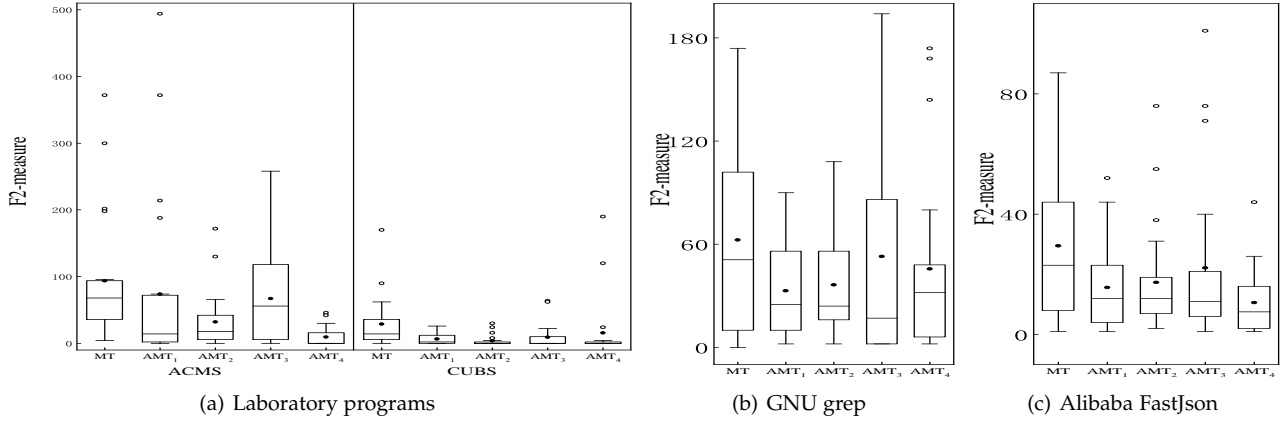


Fig. 3. F2-measure boxplots for each program

which pair of testing techniques have significant difference in terms of F-time and F2-time, as shown in Tables 16 and 17.

From Table 16 and 17, we can observe that fourteen entries (“102” & “78” for MR-AMT versus MT, “112” & “68” for MP-AMT versus MT, “91” & “89” for RR-AMT versus MT, and “120” & “60” for RP-AMT versus MT, “96” & “84” for MR-AMT versus MT, “114” & “66” for RR-AMT versus MT, and “128” & “52” for RP-AMT versus MT) are in bold font. These observations imply that in terms of F-time, AMT performed better than MT (“87” & “93” for MP-AMT versus MT). On the other hand, the performance of RP-AMT was significantly better than that of the other four techniques. In other words, the additional computation incurred in AMT for updating the probability profiles and selecting MRs is compensated with the saving of test executions.

TABLE 16

Number of Scenarios where the Technique on the top row has a lower metric (F-time) score than the technique on the left column

	MT	MR-AMT	MP-AMT	RR-AMT	RP-AMT
MT	—	<b>102</b>	<b>112</b>	<b>91</b>	<b>120</b>
MR-AMT	<b>78</b>	—	98	71	103
MP-AMT	<b>68</b>	82	—	69	97
RR-AMT	<b>89</b>	109	111	—	113
RP-AMT	<b>60</b>	77	83	67	—

TABLE 17

Number of Scenarios where the Technique on the top row has a lower metric (F2-time) score than the technique on the left column

	MT	MR-AMT	MP-AMT	RR-AMT	RP-AMT
MT	—	<b>96</b>	87	<b>114</b>	<b>128</b>
MR-AMT	<b>84</b>	—	105	116	104
MP-AMT	93	75	—	103	107
RR-AMT	<b>66</b>	64	77	—	118
RP-AMT	<b>52</b>	76	73	62	—

### 5.3 Threats To Validity

#### 5.3.1 Internal Validity

A threat to internal validity is related to the implementations of the testing techniques, which involved a moderate

amount of programming work. However, our code was cross-checked by different individuals, and we are confident that all techniques were correctly implemented.

#### 5.3.2 External Validity

The possible threat to external validity is related to the subject programs and seeded faults under evaluation. Although the four laboratory programs are not very complex, they do implement real-life business scenarios of diverse application domains. Furthermore, we chose two open source programs (grep and FastJson). Although we have tried to improve the generalisability of the findings by applying subject programs in different areas, we anticipate that the evaluation results may vary slightly with different subject web services.

#### 5.3.3 Construct Validity

The metrics used in our study are simple in concept and straightforward to apply, and hence there should be little threat to the construct validity.

#### 5.3.4 Conclusion Validity

As reported for empirical studies in the field of software engineering [52], at least 30 observations are necessary to ensure the statistical significance of results. Accordingly, we have run a sufficient number of trials to ensure the reliability of our experimental results. Furthermore, as will be discussed in Section 5, we also conducted statistical tests to confirm the significance of the results.

## 6 RELATED WORK

When testing a software system, the oracle problem appears in some situations where either an oracle does not exist for the tester to verify the correctness of the computed results; or an oracle does exist but cannot be used. The oracle problem often occurs in software testing, which renders many testing techniques inapplicable [2]. To alleviate the oracle problem, Chen et al. [5] proposed a technique named metamorphic testing (MT) that has been receiving increasing attention in the software testing community [2], [6], [9]. To improve the fault-detecting efficiency of MT, researchers mainly focus on the following aspects: i) MR; ii) source test case.



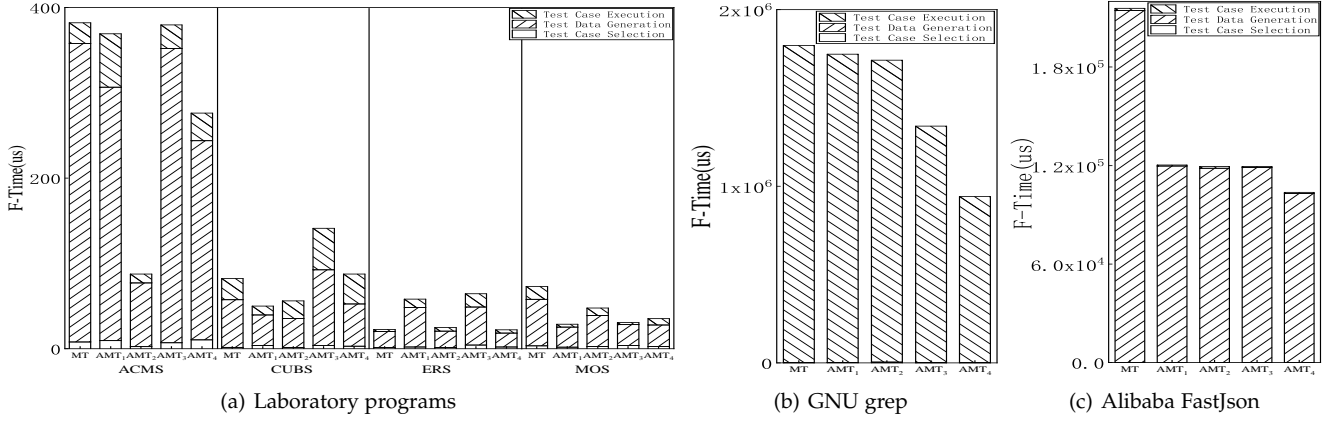


Fig. 4. The F-time results of test case selection, generation, and execution

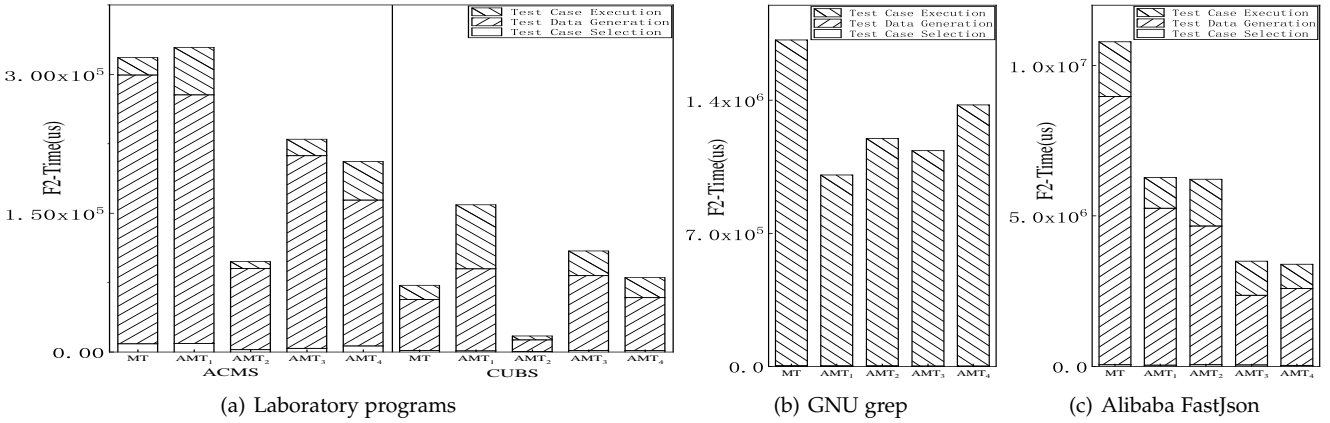


Fig. 5. The F2-time results of test case selection, generation, and execution

The MRs and the source test cases are the most important components of MT. However, defining MRs can be difficult. Chen et al. [24] proposed a specification-based method and developed a tool called MR-Generator for identifying MRs based on category-choice framework [33]. Zhang et al. [58] proposed a search-based approach to automatic inference of polynomial MRs for a software under test, where a set of parameters is used to represent polynomial MRs, and the problem of inferring MRs is turn into a problem of searching for suitable values of the parameters. Then, particle swarm optimization is used to solve the search problem. Sun et al. [59] proposes a data-mutation directed metamorphic relation acquisition methodology, in which data mutation is employed to construct input relations and the generic mapping rule associated with each mutation operator to construct output relations. Liu et al. [60] proposed to systematically construct MRs based on some already identified MRs.

With the development of various identification MR methods, how to select efficient MR from a large number of MR has become an urgent problem to be solved. Chen et al. [19] initially suggested all available metamorphic relations, of which each one might have sensitivity to different faults, should be used as part of the test strategy. However, the resources for software development are always limited. Chen et al. [21] proposed good MRs should be those that can make the multiple executions of the program as different as

possible. Furthermore, they proposed that good MRs should be selected with regard to the algorithm that the program follows. However, Mayer et al. confirmed the usefulness of whitebox considerations proposed in [21] and refined the informal criterion (“different execution path”) to criteria underlying or used in the implementation/algorithm. Asrafi et al. [citeasrafi2011testing] conducted a case study to analyze the relationship between the execution behavior and the fault-detection effectiveness of metamorphic relations, and some code coverage criteria are used to reflect the execution behavior. It is shown that the higher the combined code coverage of the source and follow-up test cases, the more different are the executions, and the more effective is the metamorphic relation. Just et al. [28] assessed the applicability of metamorphic testing for system and integration testing in the context of an image encoder, and concluded that the metamorphic relations derived from the components of a system are usually better at detecting faults than those metamorphic relations derived from the whole system. This finding was later confirmed by Xie et al. [61].

The existing metamorphic relationship selection strategies require the execution of test cases for MR selection, which increases the resource overhead of software testing. Since AMT selects efficient MRs through analyzing historical information in the testing process, it saves testing resources.

Source test cases also have a important impact on the

fault detection efficiency of MT. However, 57% of existing studies employed RT to select test cases, and 34% of existing studies used existing test suites according to a survey report by Segura et al. [9]. In other words, investigation of the impact of source test cases on MR (and MT) effectiveness and efficiency are an area yet to be explored. Chen et al. [19] compared the effects of source test cases generated by special value testing and random testing on the effectiveness of MT, and found that MT can be used as a complementary test method to special value testing. Batra and Sengupta [62] integrated genetic algorithms into MT to select source test cases maximising the the paths traversed in the software under test. Dong et al. [63] proposed a Path-Combination-Based MT method that first generates symbolic input for each executable paths and minis relationships among these symbolic inputs and their outputs, then constructs MRs on the basis of these relationships, and generates actual test cases corresponding to the symbolic inputs. Barus et al. [20] employed adaptive random testing [29] that is a class of testing method aimed to improve the performance of RT by increasing the diversity across a program's input domain of the test cases executed to generate source test cases for MT.

Different from the above investigates, we focused on performing test cases and MRs with fault revealing capabilities as quickly as possible by making use of feedback information. We first divided the input domain into disjoint partitions, and randomly selected an MR to generate follow-up test cases depended on source test case of related input partitions, then updated the test profile of input partitions according to the results of test execution. Next, a partition was selected according to updated test profile, and an MR was randomly selected from the set of MRs whose source test cases belong to selected partition.

## ACKNOWLEDGMENT

This research is supported by the National Natural Science Foundation of China (Grant Nos. 61872039 and 61872167), the Beijing Natural Science Foundation (Grant No. 4162040), the Aeronautical Science Foundation of China (Grant No. 2016ZD74004), and the Fundamental Research Funds for the Central Universities (Grant No. [FRF-GF-17-B29](#)).

## REFERENCES

- [1] E. J. Weyuker, "On testing non-testable programs," *The Computer Journal*, vol. 25, no. 4, pp. 465–470, 1982.
- [2] E. T. Barr, M. Harman, P. McMinn, M. Shahbaz, and S. Yoo, "The oracle problem in software testing: A survey," *IEEE Transactions on Software Engineering*, vol. 41, no. 5, pp. 507–525, 2015.
- [3] K. Patel and R. M. Hierons, "A mapping study on testing non-testable systems," *Software Quality Journal*, vol. 26, no. 4, pp. 1373–1413, 2018.
- [4] S. S. Brilliant, J. C. Knight, and P. Ammann, "On the performance of software testing using multiple versions," in *Proceedings of the 20th International Symposium on Fault-Tolerant Computing (FTCS'90)*, 1990, pp. 408–415.
- [5] T. Y. Chen, S. C. Cheung, and S. M. Yiu, "Metamorphic testing: A new approach for generating next test cases," Department of Computer Science, Hong Kong University of Science and Technology, Hong Kong, Tech. Rep. HKUST-CS98-01, Tech. Rep., 1998.
- [6] T. Y. Chen, F.-C. Kuo, H. Liu, P.-L. Poon, D. Towey, T. Tse, and Z. Q. Zhou, "Metamorphic testing: A review of challenges and opportunities," *ACM Computing Surveys*, vol. 51, no. 1, p. 4, 2018.
- [7] K. Y. Sim, C. S. Low, and F.-C. Kuo, "Eliminating human visual judgment from testing of financial charting software," *Journal of Software*, vol. 9, no. 2, pp. 298–312, 2014.
- [8] W. K. Chan, S.-C. Cheung, J. C. Ho, and T. Tse, "Pat: A pattern classification approach to automatic reference oracles for the testing of mesh simplification programs," *Journal of Systems and Software*, vol. 82, no. 3, pp. 422–434, 2009.
- [9] S. Segura, G. Fraser, A. B. Sanchez, and A. Ruiz-Cortés, "A survey on metamorphic testing," *IEEE Transactions on Software Engineering*, vol. 42, no. 9, pp. 805–824, 2016.
- [10] C. Murphy, M. S. Raunak, A. King, S. Chen, C. Imbriano, G. Kaiser, I. Lee, O. Sokolsky, L. Clarke, and L. Osterweil, "On effective testing of health care simulation software," in *Proceedings of the 3rd workshop on software engineering in health care (SEHC'11), Co-located with the 33th International Conference on Software Engineering (ICSE'11)*, 2011, pp. 40–47.
- [11] T. Y. Chen, J. W. Ho, H. Liu, and X. Xie, "An innovative approach for testing bioinformatics programs using metamorphic testing," *BMC bioinformatics*, vol. 10, no. 1, pp. 24–32, 2009.
- [12] Z. Hui, S. Huang, Z. Ren, and Y. Yao, "Metamorphic testing integer overflow faults of mission critical program: A case study," *Mathematical Problems in Engineering*, vol. 2013, no. 4, pp. 6–12, 2013.
- [13] Z. Q. Zhou, S. Zhang, M. Hagenbuchner, T. Tse, F.-C. Kuo, and T. Y. Chen, "Automated functional testing of online search services," *Software Testing, Verification and Reliability*, vol. 22, no. 4, pp. 221–243, 2012.
- [14] S. Segura, J. A. Parejo, J. Troya, and A. Ruiz-Cortés, "Metamorphic testing of restful web apis," *IEEE Transactions on Software Engineering*, vol. 44, no. 11, pp. 1083–1099, 2018.
- [15] Y. Tian, K. Pei, S. Jana, and B. Ray, "Deeptest: Automated testing of deep-neural-network-driven autonomous cars," in *Proceedings of the 40th International Conference on Software Engineering (ICSE'18)*, 2018, pp. 303–314.
- [16] D. Marijan, A. Gotlieb, and M. K. Ahuja, "Challenges of testing machine learning based systems," in *2019 IEEE International Conference On Artificial Intelligence Testing (AITest'19)*, 2019, pp. 101–102.
- [17] J. M. Zhang, M. Harman, L. Ma, and Y. Liu, "Machine learning testing: Survey, landscapes and horizons," *IEEE Transactions on Software Engineering*, 2020, DOI:10.1109/TSE.2019.2962027.
- [18] C.-A. Sun, A. Fu, P.-L. Poon, X. Xie, and T. Y. Chen, "Metric+: A metamorphic relation identification technique based on input plus output domains," *IEEE Transactions on Software Engineering*, 2019, available online: <https://doi.org/10.1109/TSE.2019.2934848>.
- [19] T. Y. Chen, F.-C. Kuo, Y. Liu, and A. Tang, "Metamorphic testing and testing with special values," in *Proceedings of the 5th ACIS International Conference on Software Engineering, Artificial Intelligence, Networking and Parallel/Distributed Computing (SNPD'04)*, 2004, pp. 128–134.
- [20] A. C. Barus, T. Y. Chen, F.-C. Kuo, H. Liu, and H. W. Schmidt, "The impact of source test case selection on the effectiveness of metamorphic testing," in *Proceedings of the 1st International Workshop on Metamorphic Testing (MET'16), Co-located with the 38th International Conference on Software Engineering (ICSE'16)*, 2016, pp. 5–11.
- [21] T. Y. Chen, D. Huang, T. Tse, and Z. Q. Zhou, "Case studies on the selection of useful relations in metamorphic testing," in *Proceedings of the 4th IberoAmerican Symposium on Software Engineering and Knowledge Engineering (JIISIC'04)*, 2004, pp. 569–583.
- [22] Y. Cao, Z. Q. Zhou, and T. Y. Chen, "On the correlation between the effectiveness of metamorphic relations and dissimilarities of test case executions," in *Proceedings of the 13th International Conference on Quality Software (QSIC'13)*, 2013, pp. 153–162.
- [23] C.-A. Sun, G. Wang, B. Mu, H. Liu, Z. Wang, and T. Y. Chen, "Metamorphic testing for web services: Framework and a case study," in *Proceedings of the 9th IEEE International Conference on Web Services (ICWS'11)*, 2011, pp. 283–290.
- [24] T. Y. Chen, P.-L. Poon, and X. Xie, "Metric: Metamorphic relation identification based on the category-choice framework," *Journal of Systems and Software*, vol. 116, pp. 177–190, 2016.
- [25] X. Xie, J. Li, C. Wang, and T. Y. Chen, "Looking for an mr? try metwiki today," in *Proceedings of the 1st International Workshop on Metamorphic Testing (MET'16), Co-located with the 38th International Conference on Software Engineering (ICSE'16)*, 2016, pp. 1–4.
- [26] J. Mayer and R. Guderlei, "An empirical study on the selection of good metamorphic relations," in *Proceedings of the 30th Annual*

- International Computer Software and Applications Conference (COMP-SAC'06)*, 2006, pp. 475–484.
- [27] M. Asrafi, H. Liu, and F.-C. Kuo, "On testing effectiveness of metamorphic relations: A case study," in *Proceedings of the 15th International Conference on Secure Software Integration and Reliability Improvement (SSIRI'11)*, 2011, pp. 147–156.
  - [28] R. Just and F. Schweiggert, "Automating software tests with partial oracles in integrated environments," in *Proceedings of the 5th Workshop on Automation of Software Test (AST'10)*, 2010, pp. 91–94.
  - [29] T. Y. Chen, H. Leung, and I. Mak, "Adaptive random testing," in *Proceedings of the 9th Asian Computing Science Conference (ASIAN'04)*, 2004, pp. 320–329.
  - [30] J. W. Cangussu, K.-Y. Cai, S. D. Miller, and A. P. Mathur, "Software cybernetics," *Wiley Encyclopedia of Computer Science and Engineering*, 2007.
  - [31] H. Yang, F. Chen, and S. Aliyu, "Modern software cybernetics: New trends," pp. 157–169, 2017.
  - [32] C.-A. Sun, H. Dai, H. Liu, T. Y. Chen, and K.-Y. Cai, "Adaptive partition testing," *IEEE Transactions on Computers*, vol. 68, no. 2, pp. 157–169, 2019.
  - [33] T. J. Ostrand and M. J. Balcer, "The category-partition method for specifying and generating functional tests," *Communications of the ACM*, vol. 31, no. 6, pp. 676–686, 1988.
  - [34] 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.
  - [35] P. E. Ammann and J. C. Knight, "Data diversity: An approach to software fault tolerance," *IEEE Transactions on Computers*, no. 4, pp. 418–425, 1988.
  - [36] G. B. Finelli, "Nasa software failure characterization experiments," *Reliability Engineering & System Safety*, vol. 32, no. 1–2, pp. 155–169, 1991.
  - [37] E. J. Weyuker and B. Jeng, "Analyzing partition testing strategies," *IEEE Transactions on Software Engineering*, vol. 17, no. 7, pp. 703–711, 1991.
  - [38] T. Y. Chen and Y.-T. Yu, "On the relationship between partition and random testing," *IEEE Transactions on Software Engineering*, vol. 20, no. 12, pp. 977–980, 1994.
  - [39] T. G. Project, "Grep home page." <http://www.gnu.org/software/grep/>, 2006.
  - [40] A. C. Barus, T. Y. Chen, F.-C. Kuo, H. Liu, R. Merkel, and G. Rothermel, "A novel linear-order algorithm for adaptive random testing of programs with non-numeric inputs," *Technical Report TR-UNL-CSE-2014-0004*, 2014.
  - [41] —, "A cost-effective random testing method for programs with non-numeric inputs," *IEEE Transactions on Computers*, vol. 65, no. 12, pp. 3509–3523, 2016.
  - [42] J. W. Duran and S. C. Ntafos, "An evaluation of random testing," *IEEE transactions on Software Engineering*, no. 4, pp. 438–444, 1984.
  - [43] T. Y. Chen and Y.-T. Yu, "Optimal improvement of the lower bound performance of partition testing strategies," *IEEE Proceedings-Software Engineering*, vol. 144, no. 5, pp. 271–278, 1997.
  - [44] L. Zhang, B.-B. Yin, J. Lv, K.-Y. Cai, S. S. Yau, and J. Yu, "A history-based dynamic random software testing," in *Proceedings of the 38th International Computer Software and Applications Conference Workshops (COMPSACW'14)*, 2014, pp. 31–36.
  - [45] C.-A. Sun, H. Dai, G. Wang, D. Towey, T. Y. Chen, and K.-Y. Cai, "Dynamic random testing of web services: A methodology and evaluation," *IEEE Transactions on Services Computing*, 2019, available online: <https://doi.org/10.1109/TSC.2019.2960496>.
  - [46] S. Holm, "A simple sequentially rejective multiple test procedure," *Scandinavian journal of statistics*, pp. 65–70, 1979.
  - [47] H. Do, S. Elbaum, and G. Rothermel, "Supporting controlled experimentation with testing techniques: An infrastructure and its potential impact," *Empirical Software Engineering*, vol. 10, no. 4, pp. 405–435, 2005.
  - [48] J. Lv, H. Hu, and K.-Y. Cai, "A sufficient condition for parameters estimation in dynamic random testing," in *Proceedings of the 35th IEEE Annual International Computer Software and Applications Conference Workshops (COMPSACW'11)*, 2011, pp. 19–24.
  - [49] Z. Yang, B. Yin, J. Lv, K.-Y. Cai, S. S. Yau, and J. Yu, "Dynamic random testing with parameter adjustment," in *Proceedings of the 38th IEEE Annual International Computer Software and Applications Conference Workshops (COMPSACW'14)*, 2014, pp. 37–42.
  - [50] Y. Li, B.-B. Yin, J. Lv, and K.-Y. Cai, "Approach for test profile optimization in dynamic random testing," in *Proceedings of the 39th International Computer Software and Applications Conference (COMPSAC15)*, vol. 3, 2015, pp. 466–471.
  - [51] A. C. Barus, "An in-depth study of adaptive random testing for testing program with complex input types," Ph.D. dissertation, Ph. D. Thesis. Faculty of Information and Communication Technologies, 2010.
  - [52] A. Arcuri and L. Briand, "A practical guide for using statistical tests to assess randomized algorithms in software engineering," in *Proceedings of the 33rd International Conference on Software Engineering (ICSE'11)*, 2011, pp. 1–10.
  - [53] H. Sture, "A simple sequentially rejective multiple test procedure," *Scandinavian Journal of Statistics*, vol. 6, pp. 65–70, 1979.
  - [54] X. Xie, W. E. Wong, T. Y. Chen, and B. Xu, "Metamorphic slice: An application in spectrum-based fault localization," *Information and Software Technology*, vol. 55, no. 5, pp. 866–879, 2013.
  - [55] G. Dong, S. Wu, G. Wang, T. Guo, and Y. Huang, "Security assurance with metamorphic testing and genetic algorithm," in *Proceedings of the IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology (WI-IAT'10)*, vol. 3, 2010, pp. 397–401.
  - [56] H. Liu, I. I. Yusuf, H. W. Schmidt, and T. Y. Chen, "Metamorphic fault tolerance: An automated and systematic methodology for fault tolerance in the absence of test oracle," in *Proceedings of the 36th International Conference on Software Engineering (ICSE'14)*, 2014, pp. 420–423.
  - [57] C.-A. Sun, G. Wang, B. Mu, H. Liu, Z. Wang, and T. Y. Chen, "A metamorphic relation-based approach to testing web services without oracles," *International Journal of Web Services Research*, vol. 9, no. 1, pp. 51–73, 2012.
  - [58] J. Zhang, J. Chen, D. Hao, Y. Xiong, B. Xie, L. Zhang, and H. Mei, "Search-based inference of polynomial metamorphic relations," in *Proceedings of the 29th ACM/IEEE International Conference on Automated Software Engineering (ASE'14)*, 2014, pp. 701–712.
  - [59] C.-A. Sun, Y. Liu, Z. Wang, and W. Chan, "μmt: A data mutation directed metamorphic relation acquisition methodology," in *2016 IEEE/ACM the 1st International Workshop on Metamorphic Testing (MET'16)*, Co-located with the 38th International Conference on Software Engineering (ICSE'16), 2016, pp. 12–18.
  - [60] H. Liu, X. Liu, and T. Y. Chen, "A new method for constructing metamorphic relations," in *Proceedings of the 12th International Conference on Quality Software (QSIC'12)*, 2012, pp. 59–68.
  - [61] X. Xie, J. Tu, T. Y. Chen, and B. Xu, "Bottom-up integration testing with the technique of metamorphic testing," in *Proceedings of 14th International Conference on Quality Software (QSIC'14)*, 2014, pp. 73–78.
  - [62] G. Batra and J. Sengupta, "An efficient metamorphic testing technique using genetic algorithm," in *Proceedings of the 3rd International Conference on Information Intelligence, Systems, Technology and Management (ICISTM'11)*, 2011, pp. 180–188.
  - [63] G. Dong, T. Guo, and P. Zhang, "Security assurance with program path analysis and metamorphic testing," in *Proceedings of the 4th IEEE International Conference on Software Engineering and Service Science (ICSESS'13)*, 2013, pp. 193–197.



Computer Science from Beihang University, China. His research interests include software testing, program analysis, and Service-Oriented Computing.



**Chang-ai Sun** is a Professor in the School of Computer and Communication Engineering, University of Science and Technology Beijing. Before that, he was an Assistant Professor at Beijing Jiaotong University, China, a postdoctoral fellow at the Swinburne University of Technology, Australia, and a postdoctoral fellow at the University of Groningen, The Netherlands. He received the bachelor degree in Computer Science from the University of Science and Technology Beijing, China, and the PhD degree in Computer Science from Beihang University, China. His research interests include software testing, program analysis, and Service-Oriented Computing.

**Hepeng Dai** is a PhD student in the School of Computer and Communication Engineering, University of Science and Technology Beijing, China. He received the master degree in Software Engineering from University of Science and Technology Beijing, China and the bachelor degree in Information and Computing Sciences from China University of Mining and Technology, China. His current research interests include software testing and debugging.





**Tsong Yueh Chen** is a Professor of Software Engineering at the Department of Computer Science and Software Engineering in Swinburne University of Technology. He received his PhD in Computer Science from The University of Melbourne, the MSc and DIC from Imperial College of Science and Technology, and BSc and MPhil from The University of Hong Kong. He is the inventor of metamorphic testing and adaptive random testing.



**Huai Liu** is a Lecturer of Information Technology at the College of Engineering & Science in Victoria University, Melbourne, Australia. Prior to joining VU, he worked as a research fellow at RMIT University and a research associate at Swinburne University of Technology. He received the BEng in physioelectronic technology and MEng in communications and information systems, both from Nankai University, China, and the PhD degree in software engineering from the Swinburne University of Technology,

Australia. His current research interests include software testing, cloud computing, and end-user software engineering.