

Theoretical and Empirical Analyses of the Effectiveness of Metamorphic Relation Composition

Kun Qiu[✉], Zheng Zheng[✉], *Senior Member, IEEE*,
Tsong Yueh Chen[✉], *Member, IEEE*, and Pak-Lok Poon[✉], *Member, IEEE*

Abstract—Metamorphic Relations (MRs) play a key role in determining the fault detection capability of Metamorphic Testing (MT). As human judgement is required for MR identification, systematic MR generation has long been an important research area in MT. Additionally, due to the extra program executions required for follow-up test cases, some concerns have been raised about MT cost-effectiveness. Consequently, the reduction in testing costs associated with MT has become another important issue to be addressed. MR composition can address both of these problems. This technique can automatically generate new MRs by composing existing ones, thereby reducing the number of follow-up test cases. Despite this advantage, previous studies on MR composition have empirically shown that some composite MRs have lower fault detection capability than their corresponding component MRs. To investigate this issue, we performed theoretical and empirical analyses to identify what characteristics component MRs should possess so that their corresponding composite MR has at least the same fault detection capability as the component MRs do. We have also derived a convenient, but effective guideline so that the fault detection capability of MT will most likely not be reduced after composition.

Index Terms—Metamorphic testing, metamorphic relation, metamorphic relation composition, test oracle, fault detection capability

1 INTRODUCTION

TESTING is a prominent technique for software verification [1], [2]. This technique often requires the presence of a *test oracle* (or simply an *oracle*, which refers to some mechanism for the tester to verify the correctness of the software output). However, in many situations such as testing a complex numerical algorithm, the “expected” correct software output (i.e., the oracle) is often unavailable or infeasible to determine. This problem is known as the *oracle problem*, which refers to the situation where either an oracle does not exist, or an oracle does exist but cannot be practically used, possibly due to resource constraints.

Some approaches or techniques have been proposed to alleviate the oracle problem in testing [3]. Among them, metamorphic testing (MT) has been demonstrated by various studies to be a lightweight, yet effective technique. When applying MT, the necessary properties of the software under test are first identified from various sources, such as the software

specification. These properties are expressed in the form of relations among software inputs and outputs, formally known as *metamorphic relations* (MRs). Since its first publications in 1998 [4], [5], MT has been repeatedly found to be effective at alleviating the oracle problem in software testing across many different application domains and platforms, including biomedical applications [6], [7], [8], web services [9], [10], embedded systems [11], [12], component-based software [13], compilers [14], [15], [16], machine learning classifiers [17], [18], [19], [20], [21], online search engines [22], [23], [24], image processing [25], artificial intelligence (AI) systems [26], and autonomous car systems [27], [28], [29].

Although MT has demonstrated success in alleviating the oracle problem, its application consumes additional testing resources because it involves multiple program executions. In reality, due to testing resource constraints, software testers may be particularly concerned with the *cost-effectiveness* of testing, which can be roughly defined as the ratio of the number of faults detected to the number of program executions. Regarding this issue, some researchers [30], [31] have proposed increasing the cost-effectiveness of MT through *MR composition*, so that the number of MRs used for testing can be reduced. This was originally proposed to generate new MRs from existing ones [30], [31]. Since this can reduce the number of MRs, it will also reduce the number of program executions, and thereby reduce testing costs. However, whether or not MR composition can increase the cost-effectiveness of MT also depends on the answer to the question: “How does the fault detection capability of the composite MR compare with that of its corresponding component MRs?” In this paper, we refer to a newly-constructed MR after composition as a *composite MR*, and the MRs from which it is constructed as *component MRs*.

- Kun Qiu and Zheng Zheng are with the School of Automation Science and Electrical Engineering, Beihang University, Beijing 100191, China. E-mail: {qiuqun, zhengz}@buaa.edu.cn.
- Tsong Yueh Chen is with the Department of Computer Science and Software Engineering, Swinburne University of Technology, Hawthorn, VIC 3122, Australia. E-mail: tychen@swin.edu.au.
- Pak-Lok Poon is with the School of Engineering and Technology, Central Queensland University, Melbourne, VIC 3000, Australia. E-mail: p.poon@cqu.edu.au.

Manuscript received 29 Dec. 2019; revised 30 June 2020; accepted 9 July 2020.
Date of publication 20 July 2020; date of current version 15 Mar. 2022.
(Corresponding author: Zheng Zheng.)
Recommended for acceptance by G. Fraser.
Digital Object Identifier no. 10.1109/TSE.2020.3009698

Up to now, we have found mixed answers to our question. For example, Dong *et al.* [30] reported that the fault detection capability remains largely unchanged after MR composition. On the other hand, Liu *et al.* [31] reported that the fault detection capability may be reduced after MR composition, but without mentioning: (a) under what situations the fault detection capability will be reduced; or (b) how to avoid such a reduction happening.

In view of the above mixed findings, this paper aims to answer the question: *In what situations is testing with a composite MR more cost-effective than testing with its component MRs?* Our analysis has discovered some desirable characteristics in the component MRs, which can be easily verified by the tester, that indicate whether or not the composite MR should be used instead for testing. These characteristics are defined in terms of the bijectivity/injectivity of the input/output mappings of the component MRs. In addition to formally proving the validity of these characteristics, we have also checked their practicality through an empirical study. Furthermore, we propose a useful guideline for a tester to better decide whether or not to use MR composition. In brief, given a pair of metamorphic relations MR_x and MR_y that fall into a special class of MRs (as defined in Definition 1 of Section 3.1), where MR_x is composable with MR_y , they should be used to form a composite MR if both of the following conditions are met: (a) the output mapping of MR_x is injective; and (b) the input mapping of MR_y is a bijective mapping from the source inputs of MR_y to the source inputs of MR_x .

The rest of this paper is structured as follows. Section 2 outlines the concept of MT and gives the motivation for this study. Section 3 provides the basic concepts and terminology, which facilitate the subsequent discussion of our theoretical analysis of MR composition, which is discussed in detail in Section 4. Section 5 complements Section 4 by discussing our empirical analysis of MR composition. Based on the analyses in Sections 4 and 5, Section 6 presents our general guideline for MR composition and an analysis of its applicability. This section also discusses some related works. Finally, Section 7 summarizes and concludes the paper.

2 BACKGROUND

2.1 Metamorphic Testing (MT)

MT is a lightweight, elegant, and effective technique for alleviating the oracle problem. An intuition underlying MT is as follows: Even if we cannot verify the correctness of an individual output, it may still be possible to use the relations among multiple inputs and outputs for program verification.

Example 1 (Shortest Path in an Undirected Graph).

Consider an algorithm f for computing the length of the shortest path between any two nodes (a and b) in an undirected graph G . Let: (i) P denote an implementation of f ; (ii) $P[G, a, b]$ denote the length of the shortest path between a and b in G , which is computed by P ; and (iii) $f(G, a, b)$ denote the expected (and correct) length of the shortest path. When G contains many nodes and edges, for any a and b in G , it is resource-intensive and time-

consuming to compute $f(G, a, b)$ in a brute force manner for comparing with $P[G, a, b]$. This is because the computational complexity is of factorial order of the number of nodes in G .

With MT, this tedious verification task can be alleviated by checking, for example, two properties, which are expressed as metamorphic relations MR_1 and MR_2 , as follows:

- MR_1 : If a and b in G are swapped, then $f(G, a, b) = f(G, b, a)$;
- MR_2 : If G' is a permutation of G with a' and b' being the permuted counterparts of a and b , respectively, then $f(G', a', b') = f(G, a, b)$.

With respect to MR_1 and MR_2 , we should have $P[G, a, b] = P[G, b, a]$ and $P[G, a, b] = P[G', a', b']$. Otherwise, we can conclude that P is faulty. \square

Since its first publication in 1998 [4], [5], MT has been successfully applied across a wide range of application domains and platforms. Recently, funded by the UK Engineering and Physical Sciences Research Council and the Technology TRANSfer in COMputing systems (TETRACOM) EU project, a group of academics and researchers from the Department of Computing at Imperial College London (ICL) established GraphicsFuzz—a spinout company from ICL. GraphicsFuzz [16] combined fuzzing and MT to test graphics drivers. The company was acquired by Google in 2018.

A core concept of MT is the MR, which is a necessary property of a targeted function f . An MR of f is a relation over a sequence of two or more inputs $\langle t_1, t_2, \dots, t_n \rangle$ and their corresponding outputs $\langle f(t_1), f(t_2), \dots, f(t_n) \rangle$, where $n \geq 2$. An MR can be written as $\mathbb{R} \subseteq X^n \times Y^n$, where $X^n \times Y^n$ are Cartesian products of the n inputs and their corresponding n outputs. Generally, an MR can be represented as

$$\mathbb{R}(t_1, t_2, \dots, t_n, f(t_1), f(t_2), \dots, f(t_n)).$$

Consider, for instance, MR_1 in Example 1, which can be rewritten as

$$\mathbb{R}((G, a, b), (G, b, a), f(G, a, b), f(G, b, a)).$$

Give any MR, there exists a k , where $1 \leq k < n$, such that:

- t_1, t_2, \dots, t_k denote the *source inputs*;
- $f(t_1), f(t_2), \dots, f(t_k)$ denote the *source outputs*;
- $t_{k+1}, t_{k+2}, \dots, t_n$ denote the *follow-up inputs*;
- $f(t_{k+1}), f(t_{k+2}), \dots, f(t_n)$ denote the *follow-up outputs*.

Let P be an implementation of f . With respect to an MR, applying MT typically proceeds as follows:

- 1) Replace f with P in \mathbb{R} .
- 2) Execute P on a sequence of source inputs $\langle t_1, t_2, \dots, t_k \rangle$ to obtain the corresponding sequence of source outputs $\langle P[t_1], P[t_2], \dots, P[t_k] \rangle$.
- 3) Generate a sequence of follow-up inputs $\langle t_{k+1}, t_{k+2}, \dots, t_n \rangle$ in accordance with \mathbb{R} .
- 4) Execute P on the sequence of follow-up inputs to obtain the corresponding sequence of outputs $\langle P[t_{k+1}], P[t_{k+2}], \dots, P[t_n] \rangle$.

- 5) Compare the two sets of execution results with reference to \mathbb{R} . If \mathbb{R} is violated, then P is faulty.

The above steps are repeated for every identified MR.

2.2 MR Composition

The composition technique was originally proposed as a method of generating new MRs from existing ones [30], [31]. Example 2 explains the basic concept of MR composition.

Example 2 (MR Composition). Consider the following two MRs, corresponding to two well-known properties of the *sine* function:

- MR_1 : If $x' = -x$, then $\sin(x') = -\sin(x)$;
- MR_2 : If $x' = x + 2\pi$, then $\sin(x') = \sin(x)$.

We can compose MR_1 and MR_2 together to form a composite metamorphic relation MR_{12} , that is, $MR_1(MR_2)$. MR_{12} is formally expressed as: If $x' = -(x + 2\pi)$, then $\sin(x') = -\sin(x)$. \square

Example 2 above shows that MR composition can generate new MRs from existing ones. If we only generate test cases from MR_{12} (and ignore MR_1 and MR_2) for testing, the testing cost is obviously reduced. In Example 2, testing with both MR_1 and MR_2 involves three program executions: one for $\sin(x)$, one for $\sin(x + 2\pi)$, and one for $\sin(-x)$. On the other hand, testing with MR_{12} only involves two executions: one for $\sin(x)$ and another for $\sin(-(x + 2\pi))$. Thus, if we only consider the testing cost, testing with MR_{12} alone is definitely preferable to testing with both MR_1 and MR_2 . However, beyond savings in testing costs, we should also compare the fault detection capability of MR_{12} with that of MR_1 and MR_2 . This leads to the research question of this paper: **(RQ) Will testing a composite MR (e.g., MR_{12}) have the same chance of detecting faults when compared with testing its component MRs (e.g., MR_1 and MR_2)?**

Previous studies on MR composition do not provide a definite answer to RQ. For instance, the case study reported by Dong *et al.* [30] has provided a “Yes” answer to RQ. On the other hand, the study by Liu *et al.* [31] reported that this is not necessarily the case. To date, to the best of our knowledge, no systematic studies have been conducted to address our RQ. In view of this, we perform a theoretical analysis with the intention of providing a definite answer.

3 IMPORTANT CONCEPTS AND TERMINOLOGY

Before we present our theoretical analysis, we first formalize some important definitions and concepts.

3.1 Metamorphic Relations (MRs)

In this paper, without loss of generality, we assume that the targeted function (or algorithm) involves one single input and one single output. However, generalizing our results to functions with multiple inputs and outputs is straightforward.

Because composing any two MRs into their corresponding composite MR is not always feasible, our work only considers the special class of MRs defined in this subsection.

Before formally presenting this special class of MRs, let us revisit some basic concepts of mapping.

Basic Concepts of a Function

Let $f : A \rightarrow B$ be a function (or mapping) from A to B . Here:

- A is referred to as the *domain* of f ;
- B is referred to as the *codomain* of f ;
- If $f(a) = b$, then b is referred to as the *image* of a , and a is referred to as the *preimage* of b ;
- f is said to be *injective* if $\forall a, a' \in A, f(a) = f(a')$ implies $a = a'$;
- f is said to be *surjective* if $\forall b \in B, \exists a \in A$ such that $f(a) = b$;
- f is said to be *bijective* if f is both injective and surjective;
- For any $S \subseteq A$, $f(S)$ denotes the set of the images of elements in S under the function f , such that

$$f(S) = \bigcup_{s \in S} \{f(s)\};$$

- $f(A)$ is referred to as the *range* of f , and $f(A) \subseteq B$.

The specific class of MRs considered in our study is defined as follows:

Definition 1. A Special Class of Metamorphic Relations (MRs)

Let

- $f : T \rightarrow \mathcal{R}$ be a targeted function;
- $I : T \rightarrow T'$ (where $T \subseteq \mathcal{T}$, and $T' = I(T) \subseteq T$) be a mapping that takes in a source input and generates a follow-up input for f ;
- $O : R \rightarrow R'$ (where $R = f(T)$ and $R' = O(R) \subseteq \mathcal{R}$) be a mapping that takes in a source output (i.e., $f(t)$) and generates a follow-up output.

A metamorphic relation MR is a necessary property of f . MR is formally expressed as follows:

$$\forall t \in T (f(I(t)) = O(f(t))).$$

In the above:

- T is the *set of source inputs* for MR;
- I and O are the *input* and *output mappings*, respectively, of MR;
- t is a *source input* of MR, where $t \in T$;
- $I(t)$ is the *follow-up input* corresponding to t ;
- $f(t)$ is the *source output* corresponding to t ;
- $f(I(t))$ is the *follow-up output* corresponding to t . \square

Definition 1 has two assumptions: (a) an MR involves two separate mappings (the input mapping I and the output mapping O); and (b) the input and output mappings involve a single input and a single output, respectively. We argue that MRs of this special class are commonly observed across different application domains. Our argument will be verified in response to Question Q1 in Section 6.2.

Consider Example 1 again. Let \mathcal{T} and \mathcal{R} denote the domain and codomain of f , respectively. Furthermore, let t denote an input of f . Then, for any $t \in \mathcal{T}$, t is of the form of $\langle G, a, b \rangle$. Additionally,

- $MR_1: \forall t \in T_1(f(I_1(t)) = O_1(f(t)))$, where $T_1 = \mathcal{T}$, $I_1(\langle G, a, b \rangle) = \langle G, b, a \rangle$, and $O_1(x) = x$.
- $MR_2: \forall t \in T_2(f(I_2(t)) = O_2(f(t)))$, where $T_2 = \mathcal{T}$, $I_2(\cdot)$ is a permutation function that takes in $\langle G, a, b \rangle$ and permutes G , a , and b according to a certain pattern, such that $I_2(\langle G, a, b \rangle) = \langle G', a', b' \rangle$, where G' is the permuted G , a' and b' are the counterparts of a and b , respectively, and $O_2(x) = x$.

Next, consider Example 2. Let \mathcal{T} and \mathcal{R} be the domain and codomain of the *sine* function, respectively; we write *sine* as f and t as the input of f . Then:

- $MR_3: \forall t \in T_3(f(I_3(t)) = O_3(f(t)))$, where $T_3 = \mathcal{T}$, $I_3(t) = -t$, and $O_3(x) = -x$;
- $MR_4: \forall t \in T_4(f(I_4(t)) = O_4(f(t)))$, where $T_4 = \mathcal{T}$, $I_4(t) = t + 2\pi$, and $O_4(x) = x$.

3.2 Composable MR and Composite MR

Below we first define a composable MR, which will facilitate the definition of a composite MR.

Definition 2. Composable MR

Let

- $f: \mathcal{T} \rightarrow \mathcal{R}$ be a targeted function;
- MR_x and MR_y be two MRs of f .

MR_x is said to be **composable** with MR_y if:

- $I_y(T_y) \subseteq T_x$, that is, the range of I_y is a subset of the domain of I_x ;
- $O_y(R_y) \subseteq R_x$ (where $R_x = f(T_x)$ and $R_y = f(T_y)$), that is, the range of O_y is a subset of the domain of O_x . \square

For the rest of the paper, we use subscripts to link an MR and its related components. For instance, in Definition 2 above, T_x , I_x , and O_x denote the set of source inputs, input mapping, and output mapping corresponding to MR_x , respectively.

Refer to Example 1. It can be deduced that:

- $I_1(T_1) = T_1 = T_2$ because (i) $I_1(T_1) = T_1$ and (ii) $T_1 = T_2 = \mathcal{T}$;
- $O_1(R_1) = R_1 = R_2$ because (i) $R_1 = R_2$ as $f(T_1) = f(T_2)$ and (ii) $O_1(R_1) = R_1$ because $O_1(x) = x$;
- $I_2(T_2) = T_2 = T_1$ because (i) $I_2(T_2) = T_2$ as I_2 is a permutation function and (ii) $T_1 = T_2 = \mathcal{T}$;
- $O_2(R_2) = R_2 = R_1$ because (i) $R_2 = R_1$ as $f(T_2) = f(T_1)$ and (ii) $O_2(R_2) = R_2$ as we have $O_2(x) = x$.

With (a) and (b), it follows after Definition 2 that MR_2 is composable with MR_1 . Similarly, with (c) and (d), we conclude that MR_1 is composable with MR_2 . It should, however, be noted that for any two metamorphic relations MR_x and MR_y , if MR_x is composable with MR_y , it is not necessary that MR_y is also composable with MR_x .

We next formally define the construction of a composite MR:

Definition 3. Composite MR and Component MR

Let

- $f: \mathcal{T} \rightarrow \mathcal{R}$ be a targeted function;
- MR_x and MR_y be two MRs of f ;
- MR_x be composable with MR_y .

The **composite MR** (denoted by MR_{xy}), formed by composing MR_x with MR_y , is a necessary property of f . MR_{xy} is formally expressed as follows:

$$\forall t \in T_{xy}(f(I_{xy}(t)) = O_{xy}(f(t))),$$

where

- $T_{xy} = T_y$;
- $I_{xy}(t) = (I_x \circ I_y)(t) = I_x(I_y(t))$;
- $O_{xy}(f(t)) = (O_x \circ O_y)(f(t)) = O_x(O_y(f(t)))$.

We write $MR_{xy} = MR_x \circ MR_y$ or $MR_x(MR_y)$. Also, we refer to MR_x and MR_y as the **component metamorphic relations** of MR_{xy} . \square

Refer to Example 2. It is straightforward to conclude that MR_1 and MR_2 are composable with each other according to Definition 2. We write $f(t)$ as $\sin(t)$. Then, with respect to MR_1 and MR_2 , there are two possible composite MRs.

- $MR_{12} = MR_1(MR_2) = MR_1 \circ MR_2$:

$$\forall t \in T_{12}(f(I_{12}(t)) = O_{12}(f(t))),$$

where $T_{12} = T_2 = \mathcal{T}$, $I_{12}(t) = I_1(I_2(t)) = -(t + 2\pi)$ and $O_{12}(x) = O_1(O_2(x)) = -x$.

- $MR_{21} = MR_2(MR_1) = MR_2 \circ MR_1$:

$$\forall t \in T_{21}(f(I_{21}(t)) = O_{21}(f(t))),$$

where $T_{21} = T_1 = \mathcal{T}$, $I_{21}(t) = I_2(I_1(t)) = -t + 2\pi$, and $O_{21}(x) = O_2(O_1(x)) = -x$.

It should be noted that, when applying MT, the tester is not required to explicitly specify the composite MR in the format above: The task of composing composable MRs can be automated through programming in accordance with Definition 3. This automation can be implemented as follows. Two test scripts can be written — one calling function I_y and the other calling function I_x . If t is an input to I_y , then the returned value from I_y is used as an input to I_x . In this way, t and the returned value from I_x form a pair of source and follow-up inputs for the composite MR_{xy} . Similarly, the pair of source and follow-up outputs for MR_{xy} could be obtained by first executing a test script to call the function O_y , followed by executing another test script to call the function O_x . Here, the functions I_x , I_y , O_x , and O_y are implemented according to MR_x and MR_y .

When composing more than two MRs, the resultant composite MR can also be automatically obtained by recursively applying Definition 3. Since I and O are mappings, the composition of I s and the composition of O s are associative, that is: $I_{xyz} = (I_x \circ I_y) \circ I_z = I_x \circ (I_y \circ I_z)$ and $O_{xyz} = (O_x \circ O_y) \circ O_z = O_x \circ (O_y \circ O_z)$. Since an MR is

defined in terms of its own I and O , therefore, the composition of MRs is also associative. For example, $MR_{xyz} = (MR_x \circ MR_y) \circ MR_z = MR_x \circ (MR_y \circ MR_z)$. The order of composition, however, is important. For instance, MR_{xyz} may not be the same as MR_{yxz} . In other words, the composition of MRs is not commutative.

3.3 Evaluation of Fault Detection Capability

The evaluation and comparison of the fault detection capabilities of two different MRs (regardless of whether or not they are composite) requires some metric. We defined two — one qualitative, and one quantitative — to use in our study. To facilitate the definition of these two metrics, we first present the following definition:

Definition 4. Satisfiability of a Set of Source Inputs for an MR

Let

- f be a targeted function;
- P be an implementation of f ;
- MR be a metamorphic relation of f with T , I , and O being its set of source inputs, input mapping, and output mapping, respectively;
- S be a nonempty subset of T .

After executing all elements of S with P ,

- S is said to **satisfy** MR , if all elements of S satisfy MR , that is,

$$\forall t \in S(O(P[t]) = P[I(t)]);$$

- S is said to **violate** MR , if all elements of S violate MR , that is,

$$\forall t \in S(O(P[t]) \neq P[I(t)]).$$

□

With Definition 4, a qualitative metric is defined for measuring the fault detection capability of an MR as follows:

Definition 5. Satisfiability of an MR

Let

- f be a targeted function;
- P be an implementation of f ;
- MR be a metamorphic relation of f with T being its set of source inputs.

With respect to P , if T satisfies MR , MR is said to be **satisfiable**; otherwise, MR is said to be **violative**.

Let T^v denote the set of all elements in T that violate MR . In this case, T^v is referred to as the **set of violative source inputs** of MR . Obviously, $T^v = \emptyset$, iff MR is satisfiable. □

In Definition 5 above, given an implementation P and an MR, it is equivalent to say that P violates (or satisfies) the MR, when the MR is violative (or satisfiable).

Together Definitions 4 and 5 allow us to define the following quantitative metric for measuring the fault detection capability of an MR:

Definition 6. Fault Detection Rate of an MR

Let

- f be a targeted function;
- P be an implementation of f ;
- MR be a metamorphic relation of f with T being its set of source inputs and T^v being its set of violative source inputs.

Let θ denote the **fault detection rate** of MR with respect to P . Then, θ is defined as follows:

$$\theta = \frac{|T^v|}{|T|},$$

where $|T^v|$ and $|T|$ denote the sizes of T^v and T , respectively. □

The satisfiability of an MR (Definition 5) indicates whether or not a program under test can be revealed as faulty by this MR. Furthermore, the fault detection rate (Definition 6) indicates *how likely* it is that an MR will reveal a fault in P with only one *random* source input. Larger fault detection rates indicate higher fault detection capabilities.

4 THEORETICAL ANALYSIS OF FAULT DETECTION CAPABILITY

Given an implementation P , MR_x , and MR_y , there are four possible scenarios:

- 1) Both MR_x and MR_y are satisfiable;
- 2) MR_x is satisfiable and MR_y is violative;
- 3) MR_x is violative and MR_y is satisfiable;
- 4) Both MR_x and MR_y are violative.

For each of the above scenarios, we analyze the fault detection capability of MR_{xy} .

4.1 Scenario 1

In this scenario, both MR_x and MR_y are satisfiable, that is, $\theta_x = \theta_y = 0$. Although we intuitively expect θ_{xy} to be 0, let us formally prove it (Theorem 1). This proof needs the following lemma.

Lemma 1.

Let

- f be a targeted function;
- P be an implementation of f ;
- MR_x and MR_y be two MRs of f ;
- MR_x be composable with MR_y ;
- S_y be a nonempty subset of T_y .

If S_y satisfies MR_y and $I_y(S_y)$ satisfies MR_x , then S_y satisfies MR_{xy} . □

Proof (Lemma 1). Assume that S_y satisfies MR_y and $I_y(S_y)$ satisfies MR_x . It follows after Definition 4 that

$$\forall t \in S_y(O_y(P[t]) = P[I_y(t)]), \quad (1)$$

and

$$\forall t' \in I_y(S_y) (O_x(P[t']) = P[I_x(t')]). \quad (2)$$

By the definition of $I_y(S_y)$, for any $t' \in I_y(S_y)$, there exists a $t \in S_y$ such that $t' = I_y(t)$; and for any $t \in S_y$, there exists a $t' \in I_y(S_y)$ such that $t' = I_y(t)$. Therefore, Eq. (2) can be rewritten as follows:

$$\forall t \in S_y (O_x(P[I_y(t)]) = P[I_x(I_y(t))]). \quad (3)$$

Immediately after Eqs. (1) and (3), we have

$$\forall t \in S_y (O_x(O_y(P[t])) = P[I_x(I_y(t))]). \quad (4)$$

It follows after Definition 4 that S_y satisfies MR_{xy} . \square

Now we are ready to present Theorem 1 and its proof.

Theorem 1.

Let

- f be a targeted function;
- P be an implementation of f ;
- MR_x and MR_y be two MRs of f ;
- MR_x be composable with MR_y .

If MR_x and MR_y are satisfiable, then their composite metamorphic relation, MR_{xy} , is also satisfiable. \square

Proof (Theorem 1). Assume that MR_x and MR_y are both satisfiable. From Definition 5, we immediately have T_x satisfies MR_x and T_y satisfies MR_y . Since MR_x is composable with MR_y , it follows after Definition 2 that $I_y(T_y) \subseteq T_x$. Since T_x satisfies MR_x , $I_y(T_y)$ also satisfies MR_x .

It follows from Lemma 1 that, since T_y satisfies MR_y and $I_y(T_y)$ satisfies MR_x , T_y satisfies MR_{xy} . Because $T_{xy} = T_y$ (Definition 3), therefore, it follows after Definition 5 that MR_{xy} is satisfiable, that is, $\theta_{xy} = 0$. \square

Implication. Theorem 1 states that, if an implementation P under test does not violate any two component MRs (MR_x and MR_y), P will also not violate any composite MR constructed from MR_x and MR_y .

4.2 Scenario 2

In this scenario, MR_x is satisfiable ($\theta_x = 0$) and MR_y is violative ($\theta_y > 0$). Before analyzing the fault detection capability of MR_{xy} , we first introduce the following lemma to facilitate the proof of Theorem 2.

Lemma 2.

Let

- f be a targeted function;
- P be an implementation of f ;
- MR_x and MR_y be two MRs of f ;
- MR_x be composable with MR_y ;
- S_y be a nonempty subset of T_y .

Suppose that S_y violates MR_y and $I_y(S_y)$ satisfies MR_x . If O_x is an injective mapping, then S_y violates MR_{xy} . \square

Proof (Lemma 2). Since S_y violates MR_y , it follows after Definition 4 that

$$\forall t \in S_y (O_y(P[t]) \neq P[I_y(t)]). \quad (5)$$

Because $I_y(S_y)$ satisfies MR_x , it follows after Definition 4 that

$$\forall t' \in I_y(S_y) (O_x(P[t']) = P[I_x(t')]). \quad (6)$$

By the definition of $I_y(S_y)$, for any $t' \in I_y(S_y)$, there exists a $t \in S_y$ such that $t' = I_y(t)$; and for any $t \in S_y$, there exists a $t' \in I_y(S_y)$ such that $t' = I_y(t)$. Therefore, Eq. (6) can be rewritten as follows:

$$\forall t \in S_y (O_x(P[I_y(t)]) = P[I_x(I_y(t))]). \quad (7)$$

Assume that O_x is an injective mapping. Immediately after Eqs. (5) and (7), we have

$$\forall t \in S_y (O_x(O_y(P[t])) \neq P[I_x(I_y(t))]). \quad (8)$$

Therefore, S_y violates MR_{xy} . \square

With Lemma 2, we can now introduce Theorem 2 and its proof.

Theorem 2.

Let

- f be a targeted function;
- P be an implementation of f ;
- MR_x and MR_y be two MRs of f ;
- MR_x be composable with MR_y .

Suppose that MR_x is satisfiable ($\theta_x = 0$) and MR_y is violative ($\theta_y > 0$). If O_x is an injective mapping, then MR_{xy} is violative and $\theta_{xy} = \theta_y$. \square

Proof (Theorem 2). To determine θ_{xy} , we need to know the set of source inputs (T_{xy}) and the set of violative source inputs (T_{xy}^v) of MR_{xy} . It follows from Definition 3 that $T_{xy} = T_y$. In what follows, we will prove that, if O_x is an injective mapping, then $T_{xy}^v = T_y^v$.

Since MR_y is violative, we have $T_y = T_y^v \cup \overline{T_y^v}$,¹ where $T_y^v \neq \emptyset$. Since MR_x is composable with MR_y , we have $I_y(T_y) \subseteq T_x$. In turn, we have $I_y(T_y^v) \subseteq T_x$ and $I_y(\overline{T_y^v}) \subseteq T_x$. Since MR_x is satisfiable, it follows after Definition 5 that T_x satisfies MR_x . Therefore, we have

$$I_y(T_y^v) \text{ satisfies } MR_x; \quad (i)$$

and

$$I_y(\overline{T_y^v}) \text{ satisfies } MR_x, \text{ if } I_y(\overline{T_y^v}) \neq \emptyset. \quad (ii)$$

Next, let us assume that O_x is injective. Since we have (i) above and T_y^v violates MR_y , it follows from Lemma 2 that

$$T_y^v \text{ violates } MR_{xy}. \quad (iii)$$

Furthermore, we have

$$T_y^v \subseteq T_{xy}^v, \quad (iv)$$

because T_{xy}^v contains all the elements that violate MR_{xy} .

Next, we consider the following two exhaustive cases:

1. In this paper, we use $\overline{T_y^v}$ to denote the complementary set of T_y^v over T_y . For instance, $\overline{T_y^v}$ is the complementary set of T_y^v over T_y .

Case (a): $\overline{T}_y^v \neq \emptyset$. Since we have (ii) and \overline{T}_y^v satisfies MR_y , it follows from Lemma 1 that \overline{T}_y^v satisfies MR_{xy} . Furthermore, since we have (iii) and $T_{xy} = T_y = T_y^v \cup \overline{T}_y^v$, therefore, we have $T_{xy}^v = T_y^v$.

Case (b): $\overline{T}_y^v = \emptyset$. Since we have (iv), $T_{xy}^v \subseteq T_{xy}$ (where $T_{xy} = T_y$), and $T_y = T_y^v \cup \overline{T}_y^v = T_y^v$, therefore, we have $T_{xy}^v = T_y^v$.

In view of the above two cases, regardless of whether or not \overline{T}_y^v is empty, we have $T_{xy}^v = T_y^v$. Then, it follows after Definition 6 that

$$\theta_{xy} = \frac{|T_{xy}^v|}{|T_{xy}|} = \frac{|T_y^v|}{|T_y|} = \theta_y > 0.$$

In other words, MR_{xy} is violative. \square

Implication. Theorem 2 gives a sufficient condition for MR_{xy} and MR_y having the same fault detection capability if MR_x is satisfiable.

A previous study [31] reported that composing some “loose” MRs may result in a composite MR with a lower fault detection capability. However, that study has not formally defined the meaning of “loose” MRs. By means of our theoretical analysis, we found that a “loose” MR in fact refers to one whose output mapping is not injective.

4.3 Scenario 3

In this scenario, MR_x is violative ($\theta_x > 0$) and MR_y is satisfiable ($\theta_y = 0$). Before introducing Theorems 3 and 4, we need the following lemma to facilitate their proofs.

Lemma 3.

Let

- f be a targeted function;
- P be an implementation of f ;
- MR_x and MR_y be two MRs of f ;
- MR_x be composable with MR_y ;
- S_y be a nonempty subset of T_y .

If S_y satisfies MR_y and $I_y(S_y)$ violates MR_x , then S_y violates MR_{xy} . \square

Proof (Lemma 3). Assume that there exists a nonempty $S_y \subseteq T_y$ such that S_y satisfies MR_y and $I_y(S_y)$ violates MR_x . Since S_y satisfies MR_y , it follows after Definition 4 that

$$\forall t \in S_y (O_y(P[t]) = P[I_y(t)]). \quad (9)$$

Because $I_y(S_y)$ violates MR_x , it follows after Definition 4 that

$$\forall t' \in I_y(S_y) (O_x(P[t']) \neq P[I_x(t')]). \quad (10)$$

By the definition of $I_y(S_y)$, for any $t' \in I_y(S_y)$, there exists a $t \in S_y$ such that $t' = I_y(t)$; and for any $t \in S_y$, there exists a $t' \in I_y(S_y)$ such that $t' = I_y(t)$. Therefore, Eq. (10) can be rewritten as follows:

$$\forall t \in S_y (O_x(P[I_y(t)]) \neq P[I_x(I_y(t))]). \quad (11)$$

Immediately after Eqs. (9) and (11), we have

$$\forall t \in S_y (O_x(O_y(P[t])) \neq P[I_x(I_y(t))]). \quad (12)$$

Therefore, S_y violates MR_{xy} . \square

Theorem 3.

Let

- f be a targeted function;
- P be an implementation of f ;
- MR_x and MR_y be two MRs of f ;
- MR_x be composable with MR_y .

Suppose that MR_x is violative ($\theta_x > 0$) and MR_y is satisfiable ($\theta_y = 0$). If $I_y(T_y) = T_x$, then MR_{xy} is violative with $\theta_{xy} = \frac{|T_x|}{|T_y|} > 0$ (where $T_a \neq \emptyset$, $T_a \subseteq T_y$, and $I_y(T_a) = T_x^v$). \square

Proof (Theorem 3). To determine θ_{xy} , we need to know the set of source inputs (T_{xy}) and the set of violative source inputs (T_{xy}^v) of MR_{xy} . It follows from Definition 3 that $T_{xy} = T_y$. In what follows, we will prove that, if $I_y(T_y) = T_x$, then $T_{xy}^v = T_x$ (where $T_x \neq \emptyset$, $T_x \subseteq T_y$, and $I_y(T_x) = T_x^v$).

Since MR_x is violative, we have $T_x = T_x^v \cup \overline{T}_x^v$, where $T_x^v \neq \emptyset$. Let us assume that: (a) $I_y(T_y) = T_x$; and (b) $T_a, T_b \subseteq T_y$, such that $I_y(T_a) = T_x^v$ and $I_y(T_b) = \overline{T}_x^v$. Since $T_x^v \neq \emptyset$, we have $T_a \neq \emptyset$ and

$$T_y = T_a \cup T_b. \quad (i)$$

Furthermore, since MR_y is satisfiable, it follows after Definition 5 that T_y satisfies MR_y . Therefore, we have

$$T_a \text{ satisfies } MR_y, \quad (ii)$$

and

$$T_b \text{ satisfies } MR_y, \text{ if } T_b \neq \emptyset. \quad (iii)$$

Since we have (ii) and $I_y(T_a)$ violates MR_x (as $I_y(T_a) = T_x^v$ by definition), it follows from Lemma 3 that

$$T_a \text{ violates } MR_{xy}. \quad (iv)$$

Furthermore, we have

$$T_a \subseteq T_{xy}^v, \quad (v)$$

because T_{xy}^v contains all the elements that violate MR_{xy} .

Next, we consider the following two exhaustive cases:

Case (a): $T_b \neq \emptyset$. Since we have (iii) and $I_y(T_b)$ satisfies MR_x (because $I_y(T_b) = \overline{T}_x^v$ by definition), it follows from Lemma 1 that T_b satisfies MR_{xy} . Furthermore, since we have (iv) and $T_{xy} = T_y = T_a \cup T_b$, therefore, $T_{xy}^v = T_a$.

Case (b): $T_b = \emptyset$. Since we have (v), $T_{xy}^v \subseteq T_{xy} = T_y$, and $T_y = T_a \cup T_b = T_a$, therefore, we have $T_{xy}^v = T_a$.

In view of the above two cases, regardless of whether or not T_b is empty, we have $T_{xy}^v = T_a$ (where $T_a \neq \emptyset$, $T_a \subseteq T_y$, and $I_y(T_a) = T_x^v$). It follows after Definition 6 that

$$\theta_{xy} = \frac{|T_{xy}^v|}{|T_{xy}|} = \frac{|T_a|}{|T_y|} > 0.$$

In other words, MR_{xy} is violative. \square

Theorem 4.

Let

- f be a targeted function;
- P be the implementation of f ;
- MR_x and MR_y be two MRs of f ;
- MR_x be composable with MR_y .

Suppose that MR_x is violative ($\theta_x > 0$) and MR_y is satisfiable ($\theta_y = 0$). If $I_y(T_y) = T_x$ and I_y is bijective, then MR_{xy} is violative and $\theta_{xy} = \theta_x$. \square

Proof (Theorem 4). Assume that $I_y(T_y) = T_x$. Since MR_x is violative and MR_y is satisfiable, it follows from Theorem 3 that $\theta_{xy} = \frac{|T_a|}{|T_y|}$, where $T_a \neq \emptyset$, $T_a \subseteq T_y$, and $I_y(T_a) = T_x^v$. Furthermore, assume that I_y is bijective. Immediately, we have $|T_a| = |T_x^v|$ and $|T_x| = |T_y|$. Therefore, $\theta_{xy} = \frac{|T_x^v|}{|T_x|}$. After Definition 6, we have $\theta_x = \frac{|T_x^v|}{|T_x|}$. In other words, MR_{xy} is violative and $\theta_{xy} = \theta_x$. \square

Implication. Theorem 3 provides a sufficient condition for MR_{xy} to be violative. By contrast, Theorem 4 gives a sufficient condition for MR_{xy} to be violative and $\theta_{xy} = \theta_x$.

4.4 Scenario 4

In this scenario, when MR_x is violative ($\theta_x > 0$) and MR_y is violative ($\theta_y > 0$), in theory, there may exist an implementation P for which the composite MR (MR_{xy}) is satisfiable, because all the errors generated by the violating source inputs of MR_x and MR_y may coincidentally offset each other. However, intuitively speaking, such situation will rarely occur because it is undoubtedly an extraordinary coincidence. This has led us to propose the following hypothesis:

Let f be a targeted function, MR_x and MR_y be its two MRs, and MR_y be composable with MR_x . If O_x is injective, $I_y(T_y) = T_x$, $\theta_x > 0$, and $\theta_y > 0$, then it is highly likely that MR_{xy} is violative ($\theta_{xy} > 0$).

Below we present a theoretical analysis to show why the above hypothesis is strongly held. Suppose that we have the following four assumptions:

- 1) O_x is injective;
- 2) $I_y(T_y) = T_x$;
- 3) $\theta_x > 0$;
- 4) $\theta_y > 0$.

Immediately after assumptions 3 and 4, we have $T_x^v \neq \emptyset$ and $T_y^v \neq \emptyset$. Also, in view of assumption 2 and the definition of $T_x = T_x^v \cup \overline{T_x^v}$, we can define T_a and T_b which are subsets of T_y , such that $I_y(T_a) = T_x^v$ and $I_y(T_b) = \overline{T_x^v}$. Therefore, $T_y = T_a \cup T_b$. Since $T_x^v \neq \emptyset$, we have $T_a \neq \emptyset$.

There are two schemes to partition T_y : $T_y = T_y^v \cup \overline{T_y^v}$ and $T_y = T_a \cup T_b$. If we combine both schemes together, then T_y can be partitioned into the following four sets:

- $A_1 = T_b \cap \overline{T_y^v}$;
- $A_2 = T_b \cap T_y^v$;
- $A_3 = T_a \cap \overline{T_y^v}$;
- $A_4 = T_a \cap T_y^v$.

Obviously, $T_y = A_1 \cup A_2 \cup A_3 \cup A_4$.

Next, we prove two important propositions on the above four assumptions:

Proposition 1. If $A_2 \neq \emptyset$, then A_2 violates MR_{xy} .

Proposition 2. If $A_3 \neq \emptyset$, then A_3 violates MR_{xy} .

Proof (Proposition 1). Assume that

$$A_2 \neq \emptyset. \quad (i)$$

By the definition of A_2 , we have $A_2 \subseteq T_y^v$. In other words, we have

$$A_2 \text{ violates } MR_y. \quad (ii)$$

Since $I_y(A_2) = I_y(T_b \cap T_y^v) \subseteq I_y(T_b)$ and $I_y(T_b)$ satisfies MR_x (as $I_y(T_b) = \overline{T_x^v}$ by definition), we have

$$I_y(A_2) \text{ satisfies } MR_x. \quad (iii)$$

With (i), (ii), (iii) above and assumption 1, it follows from Lemma 2 that A_2 violates MR_{xy} . \square

Proof (Proposition 2). Assume that

$$A_3 \neq \emptyset. \quad (iv)$$

By definition of A_3 , we have $A_3 \subseteq \overline{T_y^v}$. In other words, we have

$$A_3 \text{ satisfies } MR_y. \quad (v)$$

Since $I_y(A_3) = I_y(T_a \cap \overline{T_y^v}) \subseteq I_y(T_a)$ and $I_y(T_a)$ violates MR_x (as $I_y(T_a) = T_x^v$ by definition), we have

$$I_y(A_3) \text{ violates } MR_x. \quad (vi)$$

With (iv), (v), and (vi), it follows from Lemma 3 that A_3 violates MR_{xy} . \square

Next, we will show that there are two tight and restrictive relations that are necessary for MR_{xy} to be satisfiable. By Definition 5, T_{xy}^v is the set of all elements in T_{xy} that violate MR_{xy} . Therefore, after Proposition 1, we have

$$A_2 \subseteq T_{xy}^v, \quad (vii)$$

and after Proposition 2, we have

$$A_3 \subseteq T_{xy}^v. \quad (viii)$$

Let us assume that MR_{xy} is satisfiable, that is, $T_{xy}^v = \emptyset$ or $\theta_{xy} = 0$. Immediately, with (vii) and (viii) above, we have $A_2 = A_3 = \emptyset$. Because $T_b = A_1 \cup A_2 = A_1 = T_b \cap \overline{T_y^v}$, we have $T_b \subseteq \overline{T_y^v}$. Since $T_a = A_3 \cup A_4 = A_4 = T_a \cap T_y^v$, we have $T_a \subseteq T_y^v$. Furthermore, since $T_y = T_a \cup T_b = T_y^v \cup \overline{T_y^v}$, we have $T_a = T_y^v$ and $T_b = \overline{T_y^v}$. Since $I_y(T_a) = T_x^v$ and $I_y(T_b) = \overline{T_x^v}$ by definition, we have

$$I_y(T_y^v) = T_x^v, \quad (ix)$$

and

$$I_y(\overline{T_y^v}) = \overline{T_x^v}. \quad (x)$$

TABLE 1
Fault Detection Rates (θ_{xy}) of MR_{xy} in Different Testing Scenarios

| Scenario | θ_x | θ_y | $\max\{\theta_x, \theta_y\}$ | θ_{xy} |
|----------|------------|------------|--------------------------------|---|
| 1 | = 0 | = 0 | = 0 | $\theta_{xy} = 0$ (Theorem 1) |
| 2 | = 0 | > 0 | = θ_y | If O_x is injective, then $\theta_{xy} = \theta_y$ (Theorem 2) |
| 3 | > 0 | = 0 | = θ_x | If $I_y(T_y) = T_x$, then $\theta_{xy} > 0$ (Theorem 3) If $I_y(T_y) = T_x$ and I_y is bijective, then $\theta_{xy} = \theta_x$ (Theorem 4) |
| 4 | > 0 | > 0 | Varies in different situations | If $I_y(T_y) = T_x$ and O_x is injective, then it is very likely to have $\theta_{xy} > 0$ |

Since (ix) and (x) follow after the assumption that $\theta_{xy} = 0$, they are necessary relations for θ_{xy} to be 0. Relation (ix) implies that, for every violating source test input for MR_y , its corresponding follow-up input must be a violating source input for MR_x . Similarly, relation (x) implies that, for every non-violating source input for MR_y , its corresponding follow-up input must be a non-violating source input for MR_x . Obviously, these two relations are very tight and restrictive and, hence, they are unlikely to be satisfied simultaneously. Since the situation of $\theta_{xy} = 0$ requires the simultaneous satisfaction of relations (ix) and (x) (which is very rare, as explained above), it can be comfortably concluded that the situation of $\theta_{xy} > 0$ is very likely to occur. In summary, the above theoretical analysis has showed that our hypothesis will be strongly held because of two very tight and restrictive relations (ix) and (x). We performed an empirical study to support the above theoretical analysis for Scenario 4. Details will be given in Section 5.

Table 1 summarizes our theoretical analysis on the fault detection rates of MR_{xy} in the four different testing scenarios.

5 EMPIRICAL ANALYSIS OF FAULT DETECTION CAPABILITY

In the first three testing scenarios (Scenarios 1, 2, and 3) discussed in Sections 4.1, 4.2, and 4.3, we are able to obtain a definite answer after a theoretical analysis. In other words, for each of those three scenarios, we have found the characteristics that component MRs should possess to guarantee that a composite MR has the same chance of detecting the faults as its component MRs do. On the other hand, Scenario 4 is too complicated to have a definite answer solely based on a theoretical analysis. Nevertheless, our theoretical analysis of Scenario 4 has led to a hypothesis, as stated at the beginning of Section 4.4, which was further verified by an empirical study to be discussed in this section. We next discuss the settings and observations of our empirical study.

5.1 Subject Programs

Our empirical study involved the following four subject programs:

- *TriangleSquare (TSQ)*. It accepts three numbers, corresponding to the three edges of a triangle, and calculates its area if a legitimate triangle can be formed [30].

- *SparseMatrixMultiplication (SMM)*. It accepts two sparse matrices² as inputs and computes their product matrix [30].
- *Dnapars (DNA)*. It is commonly used in bioinformatics [8]. It takes in a matrix containing a set of species' DNAs and generates an evolution tree.
- *K-Nearest Neighbors (KNN)*. It is a machine learning classifier algorithm, which takes in a training data set and a testing data, and then predicts the label for the latter based on the former [19].

Table 2 gives more details about these subject programs, in terms of their inputs, outputs, and the approaches we used to assert two equivalent outputs.

5.2 Experimental Procedures

We applied the following steps to each subject program:

- 1) For each identified MR, manually check its compliance with Definition 1. This check is required because Definition 1 specifies a special (but common) class of MRs that our study has assumed.
- 2) For any tuple of two MRs (MR_x, MR_y), use Definition 2 to manually check whether MR_x is composable with MR_y . If yes, then check whether: (a) O_x is injective; and (b) $I_y(T_y) = T_x$. If yes to both (a) and (b), then generate the composite MR (MR_{xy}) from MR_x and MR_y .
- 3) Apply mutation analysis and random testing to estimate individual fault detection rates for MR_x , MR_y , and MR_{xy} .

5.3 Component and Composite MRs

Table 3 lists all the MRs used in our study; they were sourced from previous MT-related studies [8], [19], [30]. All these MRs were then checked and confirmed to comply with Definition 1 (see step 1). There were 23 ($7 + 7 + 6 + 3$) such MRs for the four subject programs. The 3rd column of this table gives the details on these MRs. For each such MR, we also explicitly list its set of source inputs (4th column), input mapping (5th column), and output mapping (6th column) in the table.

For any tuple of two MRs (MR_x, MR_y), checks were performed to ensure that conditions (a) and (b) in step 2 were fulfilled. After checking for all the tuples of two MRs, a total of 106 ($TSQ = 42$, $SMM = 42$, $DNA = 16$, and $KNN =$

2. A sparse matrix is a matrix in which most of the elements are zero.

TABLE 2
Inputs, Outputs, and Equivalent Output Assertions of Subject Programs

| Subject Programs | Inputs (t) | Outputs ($f(t)$) | Assertion of Equivalent Outputs |
|------------------|--|--|--|
| TSQ | $t = \langle a, b, c \rangle$, where a , b , and c denote the three edges of a triangle | If the three edges form a legitimate triangle, then $f(t) = s$, where s denotes the area of the triangle | Let s and s' denote two outputs. Then, $s = s'$ iff $ s - s' < \epsilon^\dagger$. |
| SMM | $t = \langle A, B \rangle$, where A and B denote two sparse matrices in the Compressed Sparse Row (CSR) format | $f(t) = C$, where C denotes the product of A and B | Let $C = [e_{ij}]$ and $C' = [e'_{ij}]$ denote two outputs. Then, $C = C'$ iff for all $ e_{ij} - e'_{ij} < \epsilon^\dagger$. |
| DNA | $t = X$, where X denotes the species' DNAs and is in the format of a $(n \times m)$ matrix (n = number of species; m = length of a DNA sequence) | $f(t) = \langle tl, tree \rangle$, where tl and $tree$ denote the length and the structural description of the generated evolution tree, respectively | Let $\langle tl, tree \rangle$ and $\langle tl', tree' \rangle$ denote two outputs. Then, $\langle tl, tree \rangle = \langle tl', tree' \rangle$ iff $ tl - tl' < \epsilon^\dagger$ and $tree$ is identical to $tree'$. |
| KNN | $t = \langle X, C, S \rangle$, where X denotes the attributes of the training data set in the format of a $(k \times m)$ matrix (k = number of entries of the training data set; m = number of attributes of each entry), C denotes the class labels of the training data with a size of k , and S denotes the testing item's attributes with a size of m elements | $f(t) = cl$, where cl denotes the calculated class label for S | Let cl and cl' denote two class labels. Then, $cl = cl'$ iff cl and cl' are exactly the same. |

(\dagger) ϵ is an extremely small value and is set to 10^{-6} .

6) eligible pairs of component MRs were found, thereby resulting in the construction of 106 composite MRs. As an example, let us consider MR_1 and MR_2 of TSQ in Table 3. We write the function corresponding to TSQ as f . It can be deduced that: (a) $T_2 = I_1(T_1) = T_1 = I_2(T_2) = T$ because both I_1 and I_2 are bijective mappings from T to T ; and (b) $R_2 = O_1(R_1) = R_1 = O_2(R_2) = f(T)$ because both O_1 and O_2 are bijective mappings from $f(T)$ to $f(T)$. It then follows from Definition 2 that MR_1 and MR_2 are composable with each other. Therefore, MR_{12} and MR_{21} were formed and tested in our study.

5.4 Measurement of Fault Detection Rates and Generation of Mutants

According to Definition 6, the fault detection rate of an MR with respect to a program P is the ratio of $|T^v|$ (the size of the MR's set of violative source inputs) to $|T|$ (the size of the MR's set of source inputs). However, because the size of T is often very large, therefore it is practically infeasible to conduct exhaustive testing to determine the value of T^v . In turn, we cannot compute an MR's fault detection rate based on Definition 6. Therefore, in this study, we used the following equation as the "estimator" of an MR's fault detection rate (θ):

$$\theta \simeq \frac{N^v}{N},$$

where N^v denotes the number of tested source inputs that caused violations to a given MR, and N denotes the total number of source inputs used in testing.

Mutation analysis [32] has long been used in MT to evaluate the fault detection rate of MRs (e.g., in [19]). Thus, we also used the mutation technique (together with random testing) to estimate the fault detection rates of component and composite MRs. Table 4 shows the mutants of the four subject programs

with injected faults. We randomly generated 10 000 ($N = 10\,000$) source inputs for each of the four subject programs.

For each mutant of every subject program, we then performed two operations: (a) used each identified component MR and each constructed composite MR to generate a separate set of follow-up inputs from the set of source inputs (with a size of 10 000); and (b) executed these source inputs and follow-up inputs with the mutants and checked for violations to MRs.

5.5 Experimental Observations

Table 5 shows the estimated fault detection rates of component and composite MRs when $\theta_x > 0$ and $\theta_y > 0$ (a total of 108 such cases). The complete set of experimental data of the estimated fault detection rates for the four subject programs are given in Tables 7, 8, and 9 in the online Appendix, which can be found on the Computer Society Digital Library at <http://doi.ieeecomputersociety.org/10.1109/TSE.2020.3009698>.

For ease of comparison and analysis, based on the data in Table 5, Fig. 1 shows the values of $\min\{\theta_x, \theta_y\}$, $\max\{\theta_x, \theta_y\}$, and θ_{xy} (where $\theta_x > 0$ and $\theta_y > 0$) for each Test Reference Number. Each subfigure of Fig. 1 corresponds to a subject program, and its x -axis corresponds to the relevant column "Test Reference No." of Table 5.

Fig. 1 and Table 5 together result in the following three observations:

- 1) For all cases where $\theta_x > 0$ and $\theta_y > 0$, we have $\theta_{xy} > 0$.
- 2) The situation where $\theta_{xy} \geq \min\{\theta_x, \theta_y\}$ was found in all cases.
- 3) In a large majority of the cases ($79/108 = 73.15\%$), the situation where $\theta_{xy} \geq \max\{\theta_x, \theta_y\}$ was found.

In summary, our empirical study showed that, θ_{xy} is never smaller than θ_x or θ_y , and θ_{xy} is very likely to be at

TABLE 3
List of MRs Used in the Empirical Study

| Sub. Pro. | MRs | Description | T | Input Mapping | Output Mapping |
|-----------|--------|---|---------------------------|--|---|
| TSQ | MR_1 | If $\langle a', b', c' \rangle = \langle b, a, c \rangle$, then $s' = s$ | $T_1 = \mathcal{T}$ | $I_1(\langle a, b, c \rangle) = \langle b, a, c \rangle$ | $O_1(s) = s$ |
| | MR_2 | If $\langle a', b', c' \rangle = \langle a, c, b \rangle$, then $s' = s$ | $T_2 = \mathcal{T}$ | $I_2(\langle a, b, c \rangle) = \langle a, c, b \rangle$ | $O_2(s) = s$ |
| | MR_3 | If $\langle a', b', c' \rangle = \langle c, b, a \rangle$, then $s' = s$ | $T_3 = \mathcal{T}$ | $I_3(\langle a, b, c \rangle) = \langle c, b, a \rangle$ | $O_3(s) = s$ |
| | MR_4 | If $\langle a', b', c' \rangle = \langle 2a, 2b, 2c \rangle$, then $s' = 4s$ | $T_4 = \mathcal{T}$ | $I_4(\langle a, b, c \rangle) = \langle 2a, 2b, 2c \rangle$ | $O_4(s) = 4s$ |
| | MR_5 | If $\langle a', b', c' \rangle = \langle \sqrt{2b^2 + 2c^2 - a^2}, b, c \rangle$, then $s' = s$ | $T_5 = \mathcal{T}$ | $I_5(\langle a, b, c \rangle) = \langle \sqrt{2b^2 + 2c^2 - a^2}, b, c \rangle$ | $O_5(s) = s$ |
| | MR_6 | If $\langle a', b', c' \rangle = \langle a, \sqrt{2a^2 + 2c^2 - b^2}, c \rangle$, then $s' = s$ | $T_6 = \mathcal{T}$ | $I_6(\langle a, b, c \rangle) = \langle a, \sqrt{2a^2 + 2c^2 - b^2}, c \rangle$ | $O_6(s) = s$ |
| | MR_7 | If $\langle a', b', c' \rangle = \langle a, b, \sqrt{2a^2 + 2b^2 - c^2} \rangle$, then $s' = s$ | $T_7 = \mathcal{T}$ | $I_7(\langle a, b, c \rangle) = \langle a, b, \sqrt{2a^2 + 2b^2 - c^2} \rangle$ | $O_7(s) = s$ |
| SMM | MR_1 | If $\langle A', B' \rangle = \langle B^T, A^T \rangle$, then $A'B' = AB$ | $T_1 = \mathcal{T}$ | $I_1(\langle A, B \rangle) = \langle B^T, A^T \rangle$ | $O_1(C) = C^T$ |
| | MR_2 | If $\langle A', B' \rangle = \langle PA, B \rangle$, then $A'B' = P(AB)$ | $T_2 = \mathcal{T}$ | $I_2(\langle A, B \rangle) = \langle PA, B \rangle$ | $O_2(C) = PC$ |
| | MR_3 | If $\langle A', B' \rangle = \langle A, BP \rangle$, then $A'B' = (AB)P$ | $T_3 = \mathcal{T}$ | $I_3(\langle A, B \rangle) = \langle A, PB \rangle$ | $O_3(C) = CP$ |
| | MR_4 | If $\langle A', B' \rangle = \langle QA, B \rangle$, then $A'B' = Q(AB)$ | $T_4 = \mathcal{T}$ | $I_4(\langle A, B \rangle) = \langle QA, B \rangle$ | $O_4(C) = QC$ |
| | MR_5 | If $\langle A', B' \rangle = \langle A, BQ \rangle$, then $A'B' = (AB)Q$ | $T_5 = \mathcal{T}$ | $I_5(\langle A, B \rangle) = \langle A, BQ \rangle$ | $O_5(C) = CQ$ |
| | MR_6 | If $\langle A', B' \rangle = \langle cA, B \rangle$, then $A'B' = c(AB)$ | $T_6 = \mathcal{T}$ | $I_6(\langle A, B \rangle) = \langle cA, B \rangle$ | $O_6(C) = cC$ |
| | MR_7 | If $\langle A', B' \rangle = \langle A, cB \rangle$, then $A'B' = c(AB)$ | $T_7 = \mathcal{T}$ | $I_7(\langle A, B \rangle) = \langle A, cB \rangle$ | $O_7(C) = cC$ |
| DNA | MR_1 | If X' is generated by swapping two sites (two columns) of X , then $tl' = tl$ and $tree' = tree$ | $T_1 = \mathcal{T}$ | $I_1(X) = XP$ | $O_1(\langle tl, tree \rangle) = \langle tl, tree \rangle$ |
| | MR_2 | If X' is generated by inserting k , say 5, uninformative sites to X , then $tl' = tl$ and $tree' = tree$ | $T_2 = \mathcal{T}$ | $I_2(X) = [X, U]$, where U is a $(n \times k)$ matrix with all its elements identical | $O_2(\langle tl, tree \rangle) = \langle tl, tree \rangle$ |
| | MR_3 | If X' is generated by removing all uninformative sites from X , then $tl' = tl$ and $tree' = tree$ | $T_3 = \mathcal{T}$ | $I_3(X)$ removes those columns in X , which have identical elements | $O_3(\langle tl, tree \rangle) = \langle tl, tree \rangle$ |
| | MR_4 | If X' is generated by concatenating the duplication X to X itself, then $tl' = 2tl$ and $tree' = tree$ | $T_4 = \mathcal{T}$ | $I_4(X) = [X, X]$ | $O_4(\langle tl, tree \rangle) = \langle 2tl, tree \rangle$ |
| | MR_5 | If X has only four rows, and X' is generated by adding a hyper-variable site to X , then $tree' = tree$ | $T_5 \subset \mathcal{T}$ | $I_5(X) = [X, S]$, where S denotes the hyper-variable site | $O_5(\langle tl, tree \rangle) = tree$ |
| | MR_6 | If X' is constructed by permutating the alphabets in X , then $tl' = tl$ and $tree' = tree$ | $T_6 = \mathcal{T}$ | $I_6(X)$ is an alphabet permutation function | $O_6(\langle tl, tree \rangle) = \langle tl, tree \rangle$ |
| KNN | MR_1 | If X' and S' are generated by affine transformation, such that $g(X) = g([x_{ij}]) = [\alpha x_{ij} + \beta]$, then $cl' = cl$ | $T_1 = \mathcal{T}$ | $I_1(\langle X, C, S \rangle) = \langle g(X), C, g(S) \rangle$ | $O_1(cl) = cl$ |
| | MR_2 | If C' is generated by a permutation function on C (that is, $C' = Perm(C)$), then $cl' = Perm(cl)$ | $T_2 = \mathcal{T}$ | $I_2(\langle X, C, S \rangle) = \langle X, Perm(C), S \rangle$ | $O_2(cl) = Perm(cl)$ |
| | MR_3 | If X' and S' are generated by permutating m columns of X and S , respectively, then $cl' = cl$ | $T_3 = \mathcal{T}$ | $I_3(\langle X, C, S \rangle) = \langle XP, C, SP \rangle$ | $O_3(cl) = cl$ |

NOTE:

- T is the set of source inputs for an MR.
- \mathcal{T} denotes the domain of a targeted function f .
- For SMM, DNA, and KNN, P denotes a matrix which is obtained by swapping two rows of the identity matrix I , and Q denotes a matrix which is obtained by multiplying one of the main principal elements of I with a scalar c .

least equal to the maximum of θ_x and θ_y (where θ_x and θ_y are both positive). Thus, the study results strongly supported our hypothesis as stated at the beginning of Section 4.4.

5.6 Threats to Validity

We discuss the possible threats to validity from the following four aspects:

MR Identification. As mentioned in Section 5.3, all the (component) MRs used in our empirical study were sourced from previous MT-related studies [8], [19], [30]. Thus, the validity of those MRs as the necessary properties of the relevant programs was already established in published works. Therefore, no threat to validity should exist in this aspect.

MT Implementation. This task involved the following steps: (a) generating composite MRs from the sourced component MRs (if applicable); (b) generating source inputs for each subject program; (c) generating follow-up inputs based on the component MRs and the composite MRs; (d) executing the mutants of each subject program with the relevant sets of source inputs and follow-up inputs; and (e) verifying the relationships between the outputs for the source and follow-up inputs. Steps (b) to (e) were fairly simple, and their implementations were straightforward. As a result, the probability of introducing mistakes in these four steps should be very low. Nevertheless, the correctness of these steps was thoroughly checked. On the other hand, for step (a), the composition of MRs should not be complex.

However, as a precaution to ensure that the generated

TABLE 4
Mutants of Subject Programs

| Sub. Pro. | Mu. ID | Code Change for Mutant Generation |
|-----------|------------|---|
| TSQ | μ_1 | Swap lines 6 and 8 |
| | μ_2 | Replace " $p=(a+b+c)/2$ " by " $p=(a+b+c)*2$ " in line 22 |
| | μ_3 | Replace " $/2$ " by " $*2$ " in lines 32, 40, and 48 |
| | μ_4 | Replace " $(\text{math.sqrt}(3.0)*a*a)/4.0$ " by " $(\text{math.sqrt}(3.0)*a*a)/2.0$ " in line 102 |
| SMM | μ_1 | Replace " n " by " 1 " in line 43 |
| | μ_2 | Replace " $c[nz] = aij * b[k]$ " by " $c[nz] = aij$ " in line 52 |
| | μ_3 | Replace " $c[nz] = aij * b[k]$ " by " $c[nz] = b[k]$ " in line 52 |
| | μ_4 | Replace " $c[col] += aij*b[k]$ " by " $c[col] += aij$ " in line 56 |
| | μ_5 | Replace " $c[col] += aij*b[k]$ " by " $c[col] += b[k]$ " in line 56 |
| DNA | μ_1 | Replace " $ns=1 < < G$ " by " $ns=1 < < C$ " in line 720 in file "seq.c" |
| | μ_2 | Replace " $ally[alias[i-1]-1] != alias[i-1]$ " by " $ally[alias[i-1]-1] > = alias[i-1]$ " in line 553 in file "seq.c" |
| | μ_3 | Replace " $i < b$ " by " $i < =$ " in line 1097 in file "seq.c" |
| | μ_4 | Replace " $j=i+1$ " by " $j=i-1$ " in line 555 in file "seq.c" |
| | μ_5 | Replace " $i < =(\text{long})O$ " by " $i > (\text{long})O$ " in line 994 in file "seq.c" |
| | μ_6 | Replace " $itemp=alias[i-1]$ " by " $itemp=alias[i+1]$ " in line 563 in file "seq.c" |
| | μ_7 | Replace " $j < =(\text{long})O$ " by " $j > =(\text{long})O$ " in line 1119 in file "seq.c" |
| | μ_8 | Replace " $p > \text{numsteps}[i] += \text{weight}[i]$ " by " $p > \text{numsteps}[i] = \text{weight}[i]$ " in line 946 in file "seq.c" |
| | μ_9 | Replace " $j < =(\text{long})O$ " by " $j > (\text{long})O$ " in line 1126 in file "seq.c" |
| | μ_{10} | Replace " $(i=)$ " by " $(i!=)$ " in line 2700 in file "seq.c" |
| KNN | μ_1 | Replace " $-$ " by " $/$ " in line 29 in function "euclideanDistance" |
| | μ_2 | Replace " $-$ " by " $+$ " in line 29 in function "euclideanDistance" |
| | μ_3 | Replace " $+=$ " by " $=$ " in line 29 in function "euclideanDistance" |
| | μ_4 | Replace " $\text{math.sqrt}(\text{distance})$ " by " distance " in line 32 in function "euclideanDistance" |
| | μ_5 | Replace " $\text{reverse}=\text{True}$ " by " $\text{reverse}=\text{False}$ " in line 41 in function "getNeighbors" |
| | μ_6 | Swap lines 57 and 59 in function "getResponse" |
| | μ_7 | Replace " $+=$ " by " $*=$ " in line 52 in function "getResponse" |
| | μ_8 | Replace " $+=$ " by " $=$ " in line 52 in function "getResponse" |
| | μ_9 | Replace " $\text{reverse}=\text{True}$ " by " $\text{reverse}=\text{False}$ " in line 55 in function "getResponse" |

composite MRs were valid with respect to their corresponding subject programs, we performed two tasks. First, we conducted a desk check on these composite MRs, during which we detected no abnormality. Second, we tested all the four subject programs (their "original" versions; not their mutants) against these composite MRs. The testing results did not reveal any MR violation. To a large extent, these two tasks provide assurance that the composite MRs were correctly generated.

Subject Program Selection. Undoubtedly, it would be desirable to have a large set of programs for our empirical study. However, using a large set of programs was prohibited due to resource constraints. Nevertheless, using these four subject programs still provided a good insight into the validity of the hypothesis as stated at the beginning of Section 4.4 because these programs cover: (a) different application domains, including numerical calculations (TSQ and SMM), bioinformatics (DNA), and machine learning classifiers (KNN); and (b) different levels of complexity (TSQ and KNN are relatively less complex in logic, whereas SMM and DNA are relatively more complex).

Mutant Generation and Selection. The mutants used for TSQ, SMM, and DNA were sourced from other previous studies [8], [30], whereas the mutants used for KNN were generated by us (because we could not find mutants for this program from the published work). To a large extent, selecting mutants from previous studies for TSQ, SMM, and DNA

helped reduce experimental bias. Whereas for KNN, mutants were randomly generated in accordance with the different types of mutation operators published in [32]. Thus, these generated mutants should form a representative set for our empirical study.

6 DISCUSSION

6.1 A General Guideline for MR Composition

Based on our theoretical and empirical analyses, we have the following convenient yet effective general guideline for performing MR composition:

A General Guideline for MR Composition

Let

- f be a targeted function;
- MR_x and MR_y be two MRs of f .

To improve the cost-effectiveness of MT, MR_{xy} should be used instead of MR_x and MR_y if:

- both MR_x and MR_y belong to the special class of MRs in accordance with Definition 1;
- MR_x is composable with MR_y according to Definition 2;
- $I_y(T_y) = T_x$, I_y is bijective, and O_x is injective.

TABLE 5
Estimated Fault Detection Rates of Component and Composite MRs (Where $\theta_x > 0$ and $\theta_y > 0$)

| Sub. Pro. | Test Reference No. | Mu. ID | θ_x | θ_y | θ_{xy} | Sub. Pro. | Test Reference No. | Mu. ID | θ_x | θ_y | θ_{xy} |
|-----------|--------------------|---------|--------------------|--------------------|-----------------------|-----------|--------------------|------------|--------------------|--------------------|-----------------------|
| TSQ | 1 | μ_1 | $\theta_1: 0.4482$ | $\theta_3: 0.4356$ | $\theta_{13}: 0.4482$ | SMM | 17 | μ_4 | $\theta_1: 0.9962$ | $\theta_7: 0.9969$ | $\theta_{17}: 0.9969$ |
| | 2 | μ_1 | $\theta_3: 0.4356$ | $\theta_1: 0.4482$ | $\theta_{31}: 0.4356$ | | 18 | μ_4 | $\theta_7: 0.9969$ | $\theta_1: 0.9962$ | $\theta_{71}: 0.9962$ |
| | 3 | μ_1 | $\theta_1: 0.4482$ | $\theta_5: 0.4372$ | $\theta_{15}: 0.6605$ | | 19 | μ_4 | $\theta_5: 0.9969$ | $\theta_7: 0.9969$ | $\theta_{57}: 0.9969$ |
| | 4 | μ_1 | $\theta_5: 0.4372$ | $\theta_1: 0.4482$ | $\theta_{51}: 0.4372$ | | 20 | μ_4 | $\theta_7: 0.9969$ | $\theta_5: 0.9969$ | $\theta_{75}: 0.9969$ |
| | 5 | μ_1 | $\theta_1: 0.4482$ | $\theta_6: 0.2123$ | $\theta_{16}: 0.4372$ | | 21 | μ_5 | $\theta_1: 0.9962$ | $\theta_4: 0.9968$ | $\theta_{14}: 0.9962$ |
| | 6 | μ_1 | $\theta_6: 0.2123$ | $\theta_1: 0.4482$ | $\theta_{61}: 0.6605$ | | 22 | μ_5 | $\theta_4: 0.9968$ | $\theta_1: 0.9962$ | $\theta_{41}: 0.9968$ |
| | 7 | μ_1 | $\theta_1: 0.4482$ | $\theta_7: 0.2249$ | $\theta_{17}: 0.2249$ | | 23 | μ_5 | $\theta_1: 0.9962$ | $\theta_6: 0.9968$ | $\theta_{16}: 0.9962$ |
| | 8 | μ_1 | $\theta_7: 0.2249$ | $\theta_1: 0.4482$ | $\theta_{71}: 0.2249$ | | 24 | μ_5 | $\theta_6: 0.9968$ | $\theta_1: 0.9962$ | $\theta_{61}: 0.9969$ |
| | 9 | μ_1 | $\theta_3: 0.4356$ | $\theta_5: 0.4372$ | $\theta_{35}: 0.6605$ | | 25 | μ_5 | $\theta_4: 0.9968$ | $\theta_6: 0.9968$ | $\theta_{46}: 0.9968$ |
| | 10 | μ_1 | $\theta_5: 0.4372$ | $\theta_3: 0.4356$ | $\theta_{53}: 0.4372$ | | 26 | μ_5 | $\theta_6: 0.9968$ | $\theta_4: 0.9968$ | $\theta_{64}: 0.9968$ |
| | 11 | μ_1 | $\theta_3: 0.4356$ | $\theta_6: 0.2123$ | $\theta_{36}: 0.2123$ | | 1 | μ_2 | $\theta_1: 0.0774$ | $\theta_2: 0.2085$ | $\theta_{12}: 0.2071$ |
| | 12 | μ_1 | $\theta_6: 0.2123$ | $\theta_3: 0.4356$ | $\theta_{63}: 0.2123$ | | 2 | μ_2 | $\theta_2: 0.2085$ | $\theta_1: 0.0774$ | $\theta_{21}: 0.205$ |
| | 13 | μ_1 | $\theta_3: 0.4356$ | $\theta_7: 0.2249$ | $\theta_{37}: 0.4372$ | | 3 | μ_2 | $\theta_3: 0.1116$ | $\theta_1: 0.0774$ | $\theta_{31}: 0.1444$ |
| | 14 | μ_1 | $\theta_7: 0.2249$ | $\theta_3: 0.4356$ | $\theta_{73}: 0.6605$ | | 4 | μ_2 | $\theta_1: 0.0774$ | $\theta_4: 0.9936$ | $\theta_{14}: 0.9943$ |
| | 15 | μ_1 | $\theta_5: 0.4372$ | $\theta_6: 0.2123$ | $\theta_{56}: 0.4372$ | | 5 | μ_2 | $\theta_4: 0.9936$ | $\theta_1: 0.0774$ | $\theta_{41}: 0.9938$ |
| | 16 | μ_1 | $\theta_6: 0.2123$ | $\theta_5: 0.4372$ | $\theta_{65}: 0.4372$ | | 6 | μ_2 | $\theta_1: 0.0774$ | $\theta_6: 0.2522$ | $\theta_{16}: 0.2543$ |
| | 17 | μ_1 | $\theta_5: 0.4372$ | $\theta_7: 0.2249$ | $\theta_{57}: 0.4372$ | | 7 | μ_2 | $\theta_6: 0.2522$ | $\theta_1: 0.0774$ | $\theta_{61}: 0.2515$ |
| | 18 | μ_1 | $\theta_7: 0.2249$ | $\theta_5: 0.4372$ | $\theta_{75}: 0.4372$ | | 8 | μ_2 | $\theta_3: 0.1116$ | $\theta_2: 0.2085$ | $\theta_{32}: 0.1116$ |
| | 19 | μ_1 | $\theta_6: 0.2123$ | $\theta_7: 0.2249$ | $\theta_{67}: 0.4372$ | | 9 | μ_2 | $\theta_2: 0.2085$ | $\theta_4: 0.9936$ | $\theta_{24}: 0.9869$ |
| | 20 | μ_1 | $\theta_7: 0.2249$ | $\theta_6: 0.2123$ | $\theta_{76}: 0.4372$ | | 10 | μ_2 | $\theta_4: 0.9936$ | $\theta_2: 0.2085$ | $\theta_{42}: 0.9799$ |
| | 21 | μ_2 | $\theta_5: 0.6652$ | $\theta_6: 0.6636$ | $\theta_{56}: 1.0$ | | 11 | μ_2 | $\theta_2: 0.2085$ | $\theta_6: 0.2522$ | $\theta_{26}: 0.2434$ |
| | 22 | μ_2 | $\theta_6: 0.6636$ | $\theta_5: 0.6652$ | $\theta_{65}: 1.0$ | | 12 | μ_2 | $\theta_6: 0.2522$ | $\theta_2: 0.2085$ | $\theta_{62}: 0.8355$ |
| | 23 | μ_2 | $\theta_5: 0.6652$ | $\theta_7: 0.6762$ | $\theta_{57}: 1.0$ | | 13 | μ_2 | $\theta_3: 0.1116$ | $\theta_4: 0.9936$ | $\theta_{34}: 0.9948$ |
| | 24 | μ_2 | $\theta_7: 0.6762$ | $\theta_5: 0.6652$ | $\theta_{75}: 1.0$ | | 14 | μ_2 | $\theta_3: 0.1116$ | $\theta_6: 0.2522$ | $\theta_{36}: 0.2567$ |
| | 25 | μ_2 | $\theta_6: 0.6636$ | $\theta_7: 0.6762$ | $\theta_{67}: 1.0$ | | 15 | μ_2 | $\theta_4: 0.9936$ | $\theta_6: 0.2522$ | $\theta_{46}: 0.995$ |
| | 26 | μ_2 | $\theta_7: 0.6762$ | $\theta_6: 0.6636$ | $\theta_{76}: 1.0$ | | 16 | μ_2 | $\theta_6: 0.2522$ | $\theta_4: 0.9936$ | $\theta_{64}: 0.995$ |
| | 27 | μ_3 | $\theta_5: 0.5487$ | $\theta_6: 0.5471$ | $\theta_{56}: 0.6605$ | DNA | 17 | μ_3 | $\theta_3: 0.07$ | $\theta_2: 0.8059$ | $\theta_{32}: 0.07$ |
| | 28 | μ_3 | $\theta_6: 0.5471$ | $\theta_5: 0.5487$ | $\theta_{65}: 0.6605$ | | 18 | μ_4 | $\theta_3: 0.0698$ | $\theta_2: 0.237$ | $\theta_{32}: 0.0698$ |
| | 29 | μ_3 | $\theta_5: 0.5487$ | $\theta_7: 0.5597$ | $\theta_{57}: 0.6605$ | | 19 | μ_4 | $\theta_2: 0.237$ | $\theta_4: 1.0$ | $\theta_{24}: 1.0$ |
| | 30 | μ_3 | $\theta_7: 0.5597$ | $\theta_5: 0.5487$ | $\theta_{75}: 0.6605$ | | 20 | μ_4 | $\theta_4: 1.0$ | $\theta_2: 0.237$ | $\theta_{42}: 1.0$ |
| | 31 | μ_3 | $\theta_6: 0.5471$ | $\theta_7: 0.5597$ | $\theta_{67}: 0.6605$ | | 21 | μ_4 | $\theta_2: 0.237$ | $\theta_6: 0.261$ | $\theta_{26}: 0.2851$ |
| | 32 | μ_3 | $\theta_7: 0.5597$ | $\theta_6: 0.5471$ | $\theta_{76}: 0.6605$ | | 22 | μ_4 | $\theta_6: 0.261$ | $\theta_2: 0.237$ | $\theta_{62}: 0.9999$ |
| | 33 | μ_4 | $\theta_5: 0.1115$ | $\theta_6: 0.1115$ | $\theta_{56}: 0.1115$ | | 23 | μ_4 | $\theta_3: 0.0698$ | $\theta_4: 1.0$ | $\theta_{34}: 1.0$ |
| | 34 | μ_4 | $\theta_6: 0.1115$ | $\theta_5: 0.1115$ | $\theta_{65}: 0.1115$ | | 24 | μ_4 | $\theta_3: 0.0698$ | $\theta_6: 0.261$ | $\theta_{36}: 0.2559$ |
| | 35 | μ_4 | $\theta_5: 0.1115$ | $\theta_7: 0.1115$ | $\theta_{57}: 0.1115$ | | 25 | μ_4 | $\theta_4: 1.0$ | $\theta_6: 0.261$ | $\theta_{46}: 1.0$ |
| | 36 | μ_4 | $\theta_7: 0.1115$ | $\theta_5: 0.1115$ | $\theta_{75}: 0.1115$ | | 26 | μ_4 | $\theta_6: 0.261$ | $\theta_4: 1.0$ | $\theta_{64}: 1.0$ |
| | 37 | μ_4 | $\theta_6: 0.1115$ | $\theta_7: 0.1115$ | $\theta_{67}: 0.1115$ | | 27 | μ_5 | $\theta_3: 0.214$ | $\theta_2: 0.9984$ | $\theta_{32}: 0.214$ |
| | 38 | μ_4 | $\theta_7: 0.1115$ | $\theta_6: 0.1115$ | $\theta_{76}: 0.1115$ | | 28 | μ_6 | $\theta_3: 0.068$ | $\theta_2: 0.2307$ | $\theta_{32}: 0.068$ |
| SMM | 1 | μ_1 | $\theta_1: 0.9996$ | $\theta_2: 1.0$ | $\theta_{12}: 0.9996$ | KNN | 29 | μ_6 | $\theta_2: 0.2307$ | $\theta_4: 0.9735$ | $\theta_{24}: 0.9753$ |
| | 2 | μ_1 | $\theta_2: 1.0$ | $\theta_1: 0.9996$ | $\theta_{21}: 0.9997$ | | 30 | μ_6 | $\theta_4: 0.9735$ | $\theta_2: 0.2307$ | $\theta_{42}: 0.9753$ |
| | 3 | μ_2 | $\theta_1: 1.0$ | $\theta_4: 1.0$ | $\theta_{14}: 1.0$ | | 31 | μ_6 | $\theta_2: 0.2307$ | $\theta_6: 0.2704$ | $\theta_{26}: 0.2809$ |
| | 4 | μ_2 | $\theta_4: 1.0$ | $\theta_1: 1.0$ | $\theta_{41}: 1.0$ | | 32 | μ_6 | $\theta_6: 0.2704$ | $\theta_2: 0.2307$ | $\theta_{62}: 0.2891$ |
| | 5 | μ_2 | $\theta_1: 1.0$ | $\theta_6: 1.0$ | $\theta_{16}: 1.0$ | | 33 | μ_6 | $\theta_3: 0.068$ | $\theta_4: 0.9735$ | $\theta_{34}: 0.9748$ |
| | 6 | μ_2 | $\theta_6: 1.0$ | $\theta_1: 1.0$ | $\theta_{61}: 1.0$ | | 34 | μ_6 | $\theta_3: 0.068$ | $\theta_6: 0.2704$ | $\theta_{36}: 0.267$ |
| | 7 | μ_2 | $\theta_4: 1.0$ | $\theta_6: 1.0$ | $\theta_{46}: 1.0$ | | 35 | μ_6 | $\theta_4: 0.9735$ | $\theta_6: 0.2704$ | $\theta_{46}: 0.986$ |
| | 8 | μ_2 | $\theta_6: 1.0$ | $\theta_4: 1.0$ | $\theta_{64}: 1.0$ | | 36 | μ_6 | $\theta_6: 0.2704$ | $\theta_4: 0.9735$ | $\theta_{64}: 0.986$ |
| | 9 | μ_3 | $\theta_1: 1.0$ | $\theta_5: 1.0$ | $\theta_{15}: 1.0$ | | 37 | μ_9 | $\theta_3: 0.214$ | $\theta_2: 1.0$ | $\theta_{32}: 0.214$ |
| | 10 | μ_3 | $\theta_5: 1.0$ | $\theta_1: 1.0$ | $\theta_{51}: 1.0$ | | 38 | μ_{10} | $\theta_3: 0.0753$ | $\theta_2: 0.6398$ | $\theta_{32}: 0.0753$ |
| | 11 | μ_3 | $\theta_1: 1.0$ | $\theta_7: 1.0$ | $\theta_{17}: 1.0$ | | 1 | μ_3 | $\theta_1: 0.0032$ | $\theta_3: 0.4465$ | $\theta_{13}: 0.4413$ |
| | 12 | μ_3 | $\theta_7: 1.0$ | $\theta_1: 1.0$ | $\theta_{71}: 1.0$ | | 2 | μ_3 | $\theta_3: 0.4465$ | $\theta_1: 0.0032$ | $\theta_{31}: 0.4479$ |
| | 13 | μ_3 | $\theta_5: 1.0$ | $\theta_7: 1.0$ | $\theta_{57}: 1.0$ | | 3 | μ_6 | $\theta_1: 0.0022$ | $\theta_2: 0.9401$ | $\theta_{12}: 0.9454$ |
| | 14 | μ_3 | $\theta_7: 1.0$ | $\theta_5: 1.0$ | $\theta_{75}: 1.0$ | | 4 | μ_6 | $\theta_2: 0.9401$ | $\theta_1: 0.0022$ | $\theta_{21}: 0.9418$ |
| | 15 | μ_4 | $\theta_1: 0.9962$ | $\theta_5: 0.9969$ | $\theta_{15}: 0.9969$ | | 5 | μ_9 | $\theta_1: 0.0013$ | $\theta_2: 0.0995$ | $\theta_{12}: 0.1024$ |
| | 16 | μ_4 | $\theta_5: 0.9969$ | $\theta_1: 0.9962$ | $\theta_{51}: 0.9962$ | | 6 | μ_9 | $\theta_2: 0.0995$ | $\theta_1: 0.0013$ | $\theta_{21}: 0.1022$ |

In the above guideline, condition (c) involves applying Theorems 2, 3, and 4 as discussed in Section 4. This condition also involves the hypothesis stated at the beginning of Section 4.4, which was confirmed by our empirical analysis (Section 5) to be highly likely to be held true. When the injectivity/bijectivity requirement of our guideline does not hold, testers should consider other information about the program under test, and the amount of testing resources available, to inform their own decisions on MR composition. For example, if the program has a long execution time, then

the reduction in program executions achieved by using composite MRs (even at the expense of a slight deterioration in fault detection effectiveness) may still be a better choice.

Table 1 summaries the fault detection capability of MR_{xy} in four different testing scenarios. It can be seen from Table 1 that, in three of the four scenarios (Scenarios 1, 2, and 3), the fault detection capability of the composite MR (i.e., MR_{xy}) is identical to those of applying both MR_x and MR_y , if the three preconditions of the general guideline are satisfied. Since fewer test cases are required for testing MR_{xy} when compared with using

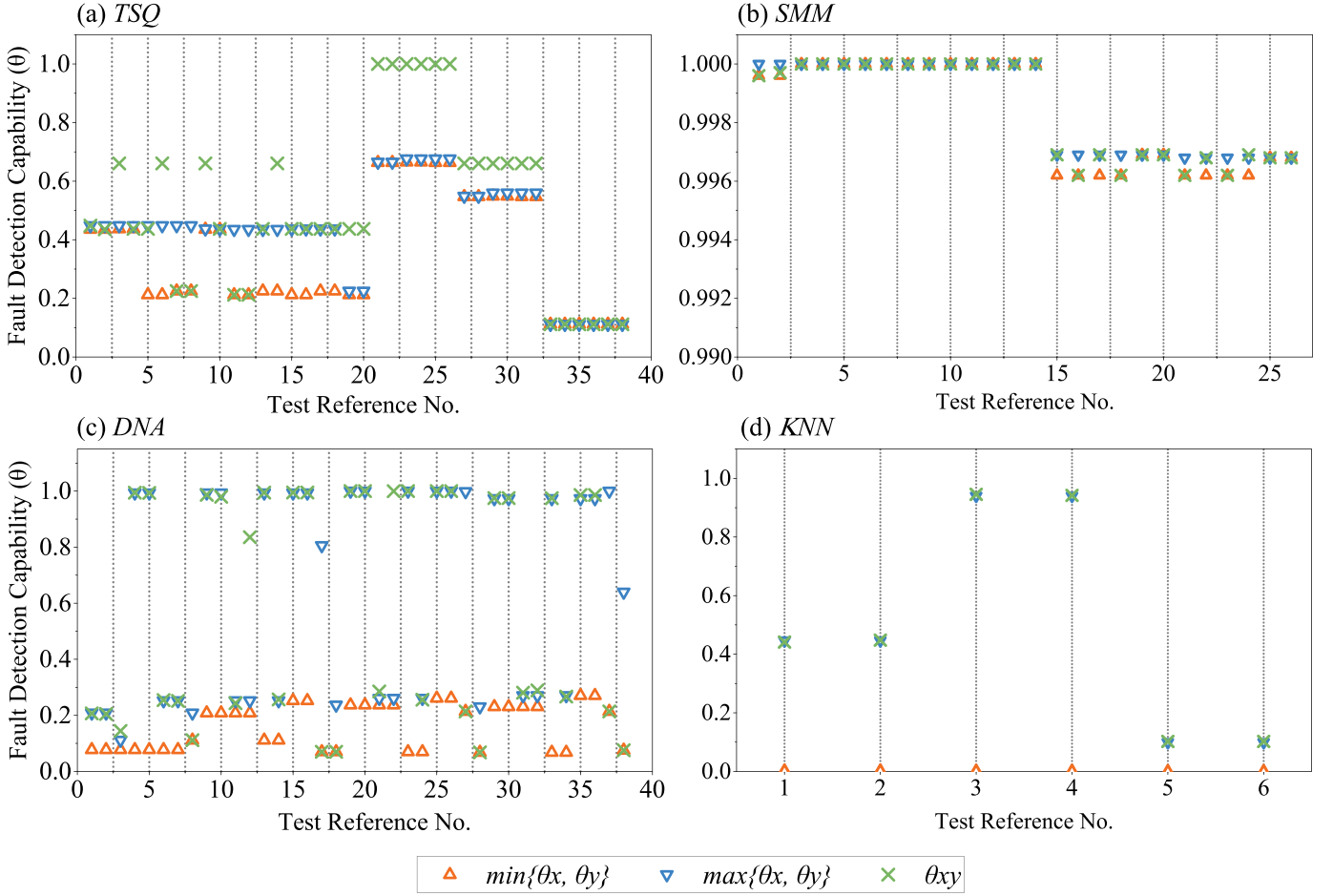


Fig. 1. Comparison among $\min\{\theta_x, \theta_y\}$, $\max\{\theta_x, \theta_y\}$, and θ_{xy} .

both MR_x and MR_y , the cost-effectiveness of MT in using MR_{xy} (instead of MR_x and MR_y) is obviously improved in these three scenarios. Even in Scenario 4 where a definite conclusion on the fault detection capability of MR_{xy} (when compared with MR_x and MR_y) cannot be drawn, we argue that it is very likely that the cost-effectiveness of MR_{xy} is higher than MR_x and MR_y because of two reasons: (a) the situation where $\theta_{xy} = 0$ should rarely occur according to our theoretical analysis (see Section 4.4), which was further confirmed to be true by our empirical analysis (see observation 1 in Section 5.5); and (b) our empirical analysis showed that $\theta_{xy} \geq \min\{\theta_x, \theta_y\}$ for all cases and, in about three-quarter of the cases, we have $\theta_{xy} \geq \max\{\theta_x, \theta_y\}$ (see observations 2 and 3 in Section 5.5).

6.2 Applicability of Our General Guideline

To evaluate the practicality and usefulness of our general guideline on MR composition, we reviewed a set of published papers on MT through which the following two questions could be answered:

- Q1:** How likely is it that a given MR belongs to the special class, in accordance with Definition 1?
- Q2:** Given an MR that belongs to the special class, according to Definition 1, how likely is this MR to have a bijective input mapping I and an injective output mapping O ?

We first studied in detail two recent survey papers on MT [37], [38], and then identified some other works on MT

that have been developed after publishing those two surveys. After these exercises, we identified 10 popular application domains for MT as shown in Table 6. For each of these domains, we found some relevant published works related to MT. After a close examination, we compiled a list of MRs which were mentioned in these published works. Further checking of the list allowed us to identify some “common” MRs with similar types or characteristics within the same domain and even across different domains. These “common” MRs were counted only once in our review. For example, among the three MT-related papers on compilers, after tallying the count for “common” MRs, we found eight distinct MRs. We caution readers that, although our review did not (which was also infeasible to) involve every MT-related paper, we argue that the compiled list of MRs from our collected papers was fairly comprehensive because the list of MRs covered 10 different and popular domains for MT.

For Q1, we found 54.88 percent of MRs belong to the special class, in accordance with Definition 1. For Q2, we found that, among those MRs complying with Definition 1, 91.11 percent of them have a bijective I and an injective O . Based on these findings for Q1 and Q2, we can conclude that the three preconditions in our general MR composition guideline can easily be met ($50.00\% \approx 54.88\% \times 91.11\%$). This shows that our general guideline should be widely applicable to many testing scenarios and application domains.

We also noted that the applicability of our guideline varies across different application domains. More specifically,

TABLE 6
Applicability of the Guideline

| Application Domains | References | No. of Identified MRs | No. of MRs complying with Definition 1 | No. of MRs whose I is bijective and O is injective | Answer to Q1 (%) | Answer to Q2 (%) |
|--|------------------------------|-----------------------|--|--|------------------|------------------|
| Biomedical applications | [6], [7], [8] | 42 | 12 | 12 | 28.57 | 100.00 |
| Web services | [9], [10] | 9 | 3 | 3 | 33.33 | 100.00 |
| Embedded systems | [11], [12] | 3 | 0 | 0 | 0.00 | 0.00 |
| Component-based software | [13] | 3 | 2 | 2 | 66.67 | 100.00 |
| Compilers | [14], [15], [16] | 8 | 8 | 5 | 100.00 | 62.50 |
| Machine learning classifiers | [17], [18], [19], [20], [21] | 16 | 7 | 7 | 43.75 | 100.00 |
| Online search engines | [22], [23], [24] | 9 | 0 | 0 | 0.00 | 0.00 |
| Assorted computer science algorithms | [33], [34], [35] | 27 | 16 | 11 | 59.26 | 68.75 |
| Numerical and scientific programs | [30], [36] | 17 | 14 | 14 | 82.35 | 100.00 |
| AI systems (e.g., image processing and autonomous car systems) | [25], [26], [27], [28] | 30 | 28 | 28 | 93.33 | 100.00 |
| Total | | 164 | 90 | 82 | 54.88 | 91.11 |

according to the results in Table 6, the guideline is largely applicable to compilers, numeric and scientific programs, and AI systems (e.g., image processing and autonomous car systems), and relatively less applicable to biomedical applications, web services, embedded systems, and online search engines. We caution readers that this observation is based on the MRs reported in the literature. To date, only a relatively small set of MRs for compilers have been studied, and they all involve equality relations between source and follow-up outputs. This makes the guideline applicable to most of these MRs. Due to the nature of numeric and scientific application domains, the input and output mappings of the reported MRs are often composable and bijective, and therefore the guideline is also mostly applicable to these domains. Furthermore, image processing and autonomous driving systems are currently the main AI systems using MT. Most of the MRs defined for these systems are based on spatial transformation mappings that are composable and bijective, resulting in the applicability of the guideline in these domains. On the other hand, the reported MRs for biomedical applications, web services, embedded systems, and online search engines often involve subset or substring relations, which are less likely to have bijective input and output mappings. The guideline, therefore, is less applicable to such systems. However, we also wish to highlight that this observation may vary as more MRs emerge.

6.3 Related Work

Two major challenges for MT are: (a) identification of MRs; and (b) additional computations of the program executions for follow-up test cases. Although MR composition can address both challenges, only few papers primarily focused on the effectiveness of MR composition — we found only two papers [30], [31] in this area. Obviously, MR composition reduces the number of program executions and, hence, lowers the computation costs — this is indisputable. However, these two studies on MR composition [30], [31] do not have a consensus on whether or not the fault detection capability after composition will be jeopardized. Furthermore, both studies [30], [31] adopted a purely empirical approach. Therefore, their observations

could not provide a full picture of this issue and are dependent on the subjects being investigated. Understandably, some of their observations may look contradictory, that is, they reach different conclusions on the fault detection effectiveness of the composite MRs. On the other hand, with our theoretical results, such illusive contradictions are clarified. In summary, their results [30], [31] motivated our study, which in turn provides a more comprehensive interpretation of their results.

MR composition is an obvious and straightforward method to generate new MRs that is easily implemented and automated. However, this method requires the existence of some MRs for generation of new ones. Recently, with the increasing recognition and acceptance of MT by the software testing community, a growing number of research studies on MR generation/identification has emerged. Examples of these studies include machine-learning-based techniques [39], [40], [41], search-based techniques [42], [43], data-mutation-based techniques [26], [44], [45], pattern-based techniques [46], [47], and the category-choice approach [48], [49]. Since the main focus of this paper is not on MR generation, comparing our work with the above studies is beyond the scope of this paper.

7 SUMMARY AND CONCLUSION

Two major advantages of MR composition are the facilitation of automatic MR generation, and the reduction in testing costs (by reducing the number of program executions for follow-up test cases). However, MR composition has a potential drawback: The fault detection capability of the composite MR may be lower than that of its component MRs, jeopardizing the overall effectiveness of MT. This issue motivated us to perform theoretical and empirical analyses, with a goal of identifying characteristics that the component MRs should possess so that the fault detection capability of the generated composite MR will likely not be less than that of its component MRs. In short, given a pair of metamorphic relations MR_x and MR_y that belong to the special class defined in Definition 1, where MR_x is composable with MR_y , they should be used to form a composite MR_{xy} if both the following

conditions are met: (a) the output mapping of MR_x is injective; and (b) the input mapping of MR_y is a bijective mapping from the source inputs of MR_y to the source inputs of MR_x . This result is produced based on Theorems 1 to 4, Propositions 1 and 2, and the empirical analysis discussed in the paper. This result provides a solid foundation for MT and has paved a path for future studies on MR composition.

Based on our theoretical and empirical analyses, a convenient yet effective general guideline on MR composition has been developed. We further performed studies (by using a sample of MRs extracted from previously published works on MT) to confirm the applicability of the guideline across a range of application domains.

Our study has focused on a special class of MRs, as stated in Definition 1. Although this class of MRs is common, it would be worthwhile to extend our study to examine MRs that are outside this class. A hierarchy of composable MRs and their composite MRs could then be built. It would be interesting to investigate what further information and insights could be exploited from this hierarchy, using the theoretical results reported in this paper. This future study would enhance the foundation of MT research. Another potentially fruitful direction is a large-scale empirical analysis of the relationships among θ_{xy} , $\min\{\theta_x, \theta_y\}$, and $\max\{\theta_x, \theta_y\}$. The results of this empirical analysis would help practitioners better estimate their testing costs.

ACKNOWLEDGMENTS

We are indebted to Dr. Dave Towey of the University of Nottingham Ningbo China for his valuable comments and suggestions on improving this paper. This work was supported in part by the National Natural Science Foundation of China (Grant Nos. 61772055 and 61872169), the Technical Foundation Project of Ministry of Industry and Information Technology of China (Grant No. JSZL2016601B003), and the Equipment Preliminary R&D Project of China (Grant No. 41402020102).

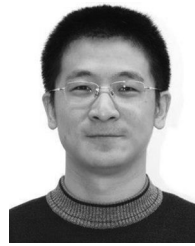
REFERENCES

- [1] P. Ammann and J. Offutt, *Introduction to Software Testing*. New York, NY: Cambridge Univ. Press, 2018.
- [2] M. Pezzè and M. Young, *Software Testing and Analysis: Process, Principles, and Techniques*. Noida, India: Wiley, 2007.
- [3] E. T. Barr, M. Harman, P. McMinn, M. Shahbaz, and S. Yoo, "The oracle problem in software testing: A survey," *IEEE Trans. Softw. Eng.*, vol. 41, no. 5, pp. 507–525, May 2015.
- [4] T. Y. Chen, S. C. Cheung, and S. M. Yiu, "Metamorphic testing: A new approach for generating next test cases," *Hong Kong Univ. Sci. Technol.*, Hong Kong, *Tech. Rep. HKUST-CS98-01*, 1998.
- [5] F. T. Chan, T. Y. Chen, S. C. Cheung, M. F. Lau, and S. M. Yiu, "Application of metamorphic testing in numerical analysis," in *Proc. IASTED Int. Conf. Softw. Eng.*, 1998, pp. 191–197.
- [6] T. Y. Chen, J. W. K. Ho, H. Liu, and X. Xie, "An innovative approach for testing bioinformatics programs using metamorphic testing," *BMC Bioinf.*, vol. 10, no. 1, 2009, Art. no. 24.
- [7] L. L. Pullum and O. Ozmen, "Early results from metamorphic testing of epidemiological models," in *Proc. ASE/IEEE Int. Conf. Biomed. Comput.*, 2012, pp. 62–67.
- [8] M. S. Sadi, F.-C. Kuo, J. W. K. Ho, M. A. Charleston, and T. Y. Chen, "Verification of phylogenetic inference programs using metamorphic testing," *J. Bioinf. Comput. Biol.*, vol. 9, no. 6, pp. 729–747, 2011.
- [9] W. K. Chan, S. C. Cheung, and K. R. P. Leung, "A metamorphic testing approach for online testing of service-oriented software applications," *Int. J. Web Serv. Res.*, vol. 4, no. 2, pp. 61–81, 2007.
- [10] C.-A. Sun, G. Wang, B. Mu, H. Liu, Z. S. Wang, and T. Y. Chen, "Metamorphic testing for web services: Framework and a case study," in *Proc. IEEE Int. Conf. Web Serv.*, 2011, pp. 283–290.
- [11] W. K. Chan, T. Y. Chen, H. Lu, T. H. Tse, and S. S. Yau, "Integration testing of context-sensitive middleware-based applications: A metamorphic approach," *Int. J. Softw. Eng. Knowl. Eng.*, vol. 16, no. 5, pp. 677–703, 2006.
- [12] F.-C. Kuo, T. Y. Chen, and W. K. Tam, "Testing embedded software by metamorphic testing: A wireless metering system case study," in *Proc. 36th Conf. Local Comput. Netw.*, 2011, pp. 291–294.
- [13] X.-L. Lu, Y.-W. Dong, and C. Luo, "Testing of component-based software: A metamorphic testing methodology," in *Proc. 7th Int. Conf. Ubiquitous Intell. Comput. 7th Int. Conf. Autonomic Trusted Comput.*, 2010, pp. 272–276.
- [14] V. Le, M. Afshari, and Z. Su, "Compiler validation via equivalence modulo inputs," *ACM SIGPLAN Notices*, vol. 49, no. 6, pp. 216–226, 2014.
- [15] Q. Tao, W. Wu, C. Zhao, and W. Shen, "An automatic testing approach for compiler based on metamorphic testing technique," in *Proc. Asia-Pacific Softw. Eng. Conf.*, 2010, pp. 270–279.
- [16] A. F. Donaldson, H. Evrard, A. Lascu, and P. Thomson, "Automated testing of graphics shader compilers," in *Proc. ACM Program. Lang.*, 2017, Art. no. 93.
- [17] C. Murphy, G. E. Kaiser, and L. Hu, "Properties of machine learning applications for use in metamorphic testing," Columbia Univ., New York, NY, USA, *Tech. Rep. CUICS-011-08*, 2011.
- [18] X. Xie, J. Ho, C. Murphy, G. Kaiser, B. Xu, and T. Y. Chen, "Application of metamorphic testing to supervised classifiers," in *Proc. 9th Int. Conf. Qual. Softw.*, 2009, pp. 135–144.
- [19] X. Xie, J. W. K. Ho, C. Murphy, G. Kaiser, B. Xu, and T. Y. Chen, "Testing and validating machine learning classifiers by metamorphic testing," *J. Syst. Softw.*, vol. 84, no. 4, pp. 544–558, 2011.
- [20] S. Nakajima and H. N. Bui, "Dataset coverage for testing machine learning computer programs," in *Proc. 23rd Asia-Pacific Softw. Eng. Conf.*, 2016, pp. 297–304.
- [21] P. Saha and U. Kanewala, "Fault detection effectiveness of metamorphic relations developed for testing supervised classifiers," in *Proc. IEEE Int. Conf. Artif. Intell. Testing*, 2019, pp. 157–164.
- [22] Z. Q. Zhou, S. Zhang, M. Hagenbuchner, T. H. Tse, F.-C. Kuo, and T. Y. Chen, "Automated functional testing of online search services," *Softw. Testing Verification Rel.*, vol. 22, no. 4, pp. 221–243, 2012.
- [23] Z. Q. Zhou, S. Xiang, and T. Y. Chen, "Metamorphic testing for software quality assessment: A study of search engines," *IEEE Trans. Softw. Eng.*, vol. 42, no. 3, pp. 264–284, Mar. 2016.
- [24] J. Brown, Z. Q. Zhou, and Y.-W. Chow, "Metamorphic testing of navigation software: A pilot study with Google maps," in *Proc. 51st Hawaii Int. Conf. Syst. Sci.*, 2018, pp. 5687–5696.
- [25] J. Mayer and R. Guderlei, "On random testing of image processing applications," in *Proc. 6th Int. Conf. Qual. Softw.*, 2006, pp. 85–92.
- [26] H. Zhu, D. Liu, I. Bayley, R. Harrison, and F. Cuzzolin, "Datamorphic testing: A method for testing intelligent applications," in *Proc. IEEE Int. Conf. Artif. Intell. Testing*, 2019, pp. 149–156.
- [27] M. Zhang, Y. Zhang, L. Zhang, C. Liu, and S. Khurshid, "DeepRoad: GAN-based metamorphic testing and input validation framework for autonomous driving systems," in *Proc. 33rd ACM/IEEE Int. Conf. Automated Softw. Eng.*, 2018, pp. 132–142.
- [28] Y. Tian, K. Pei, S. Jana, and B. Ray, "DeepTest: Automated testing of deep-neural-network-driven autonomous cars," in *Proc. 40th Int. Conf. Softw. Eng.*, 2018, pp. 303–314.
- [29] Z. Q. Zhou and L. Sun, "Metamorphic testing of driverless cars," *Commun. ACM*, vol. 62, no. 3, pp. 61–67, 2019.
- [30] G. Dong, B. Xu, L. Chen, C. Nie, and L. Wang, "Case studies on testing with compositional metamorphic relations," *J. Southeast Univ.*, vol. 24, no. 4, pp. 437–443, 2008.
- [31] H. Liu, X. Liu, and T. Y. Chen, "A new method for constructing metamorphic relations," in *Proc. 12th Int. Conf. Qual. Softw.*, 2012, pp. 59–68.
- [32] B. H. Smith and L. Williams, "On guiding the augmentation of an automated test suite via mutation analysis," *Empir. Softw. Eng.*, vol. 14, no. 3, pp. 341–369, 2009.
- [33] M. Jiang, T. Y. Chen, F.-C. Kuo, and Z. Ding, "Testing central processing unit scheduling algorithms using metamorphic testing," in *Proc. 4th Int. Conf. Softw. Eng. Service Sci.*, 2013, pp. 530–536.

- [34] P. Rao, Z. Zheng, T. Y. Chen, N. Wang, and K. Cai, "Impacts of test suite's class imbalance on spectrum-based fault localization techniques," in *Proc. 13th Int. Conf. Qual. Softw.*, 2013, pp. 260–267.
- [35] X. Xie, W. E. Wong, T. Y. Chen, and B. Xu, "Metamorphic slice: An application in spectrum-based fault localization," *Inf. Softw. Technol.*, vol. 55, no. 5, pp. 866–879, 2013.
- [36] T. Y. Chen, J. Feng, and T. H. Tse, "Metamorphic testing of programs on partial differential equations: A case study," in *Proc. 26th Annu. Int. Comput. Softw. Appl.*, 2002, pp. 327–333.
- [37] S. Segura, G. Fraser, A. B. Sanchez, and A. Ruiz-Cortés, "A survey on metamorphic testing," *IEEE Trans. Softw. Eng.*, vol. 42, no. 9, pp. 805–824, Sep. 2016.
- [38] T. Y. Chen *et al.*, "Metamorphic testing: A review of challenges and opportunities," *ACM Comput. Surv.*, vol. 51, no. 1, 2018, Art. no. 4.
- [39] U. Kanewala and J. M. Bieman, "Using machine learning techniques to detect metamorphic relations for programs without test oracles," in *Proc. IEEE 24th Int. Symp. Softw. Rel. Eng.*, 2013, pp. 1–10.
- [40] U. Kanewala, J. M. Bieman, and A. Ben-Hur, "Predicting metamorphic relations for testing scientific software: A machine learning approach using graph kernels," *Softw. Testing Verification Rel.*, vol. 26, no. 3, pp. 245–269, 2016.
- [41] A. Nair, K. Meinke, and S. Eldh, "Leveraging mutants for automatic prediction of metamorphic relations using machine learning," in *Proc. 3rd ACM SIGSOFT Int. Workshop Mach. Learn. Techn. Softw. Qual. Eval.*, 2019, pp. 1–6.
- [42] J. Zhang *et al.*, "Search-based inference of polynomial metamorphic relations," in *Proc. 29th ACM/IEEE Int. Conf. Automated Softw. Eng.*, 2014, pp. 701–712.
- [43] B. Zhang, H. Zhang, J. Chen, D. Hao, and P. Moscato, "Automatic discovery and cleansing of numerical metamorphic relations," in *Proc. IEEE Int. Conf. Softw. Maintenance Evol.*, 2019, pp. 235–245.
- [44] H. Zhu, "JFuzz: A tool for automated Java unit testing based on data mutation and metamorphic testing methods," in *Proc. 2nd Int. Conf. Trustworthy Syst. Appl.*, 2015, pp. 8–15.
- [45] C.-A. Sun, Y. Liu, Z. Wang, and W. K. Chan, "μMT: A data mutation directed metamorphic relation acquisition methodology," in *Proc. IEEE/ACM 1st Int. Workshop Metamorphic Testing*, 2016, pp. 12–18.
- [46] S. Segura, J. A. Parejo, J. Troya, and A. Ruiz-Cortés, "Metamorphic testing of RESTful web APIs," *IEEE Trans. Softw. Eng.*, vol. 44, no. 11, pp. 1083–1099, Nov. 2018.
- [47] Z. Q. Zhou, L. Sun, T. Y. Chen, and D. Towey, "Metamorphic relations for enhancing system understanding and use," *IEEE Trans. Softw. Eng.*, to be published, doi: [10.1109/TSE.2018.2876433](https://doi.org/10.1109/TSE.2018.2876433).
- [48] T. Y. Chen, P.-L. Poon, and X. Xie, "METRIC: METAmorphic Relation Identification based on the Category-choice framework," *J. Syst. Softw.*, vol. 116, pp. 177–190, 2016.
- [49] C.-A. Sun, A. Fu, P.-L. Poon, X. Xie, H. Liu, and T. Y. Chen, "METRIC+: A metamorphic relation identification technique based on input plus output domains," *IEEE Trans. Softw. Eng.*, to be published, doi: [10.1109/TSE.2019.2934848](https://doi.org/10.1109/TSE.2019.2934848).



Kun Qiu received the BS degree in automation from the Hefei University of Technology, Hefei, China, in 2013. He is currently working toward the PhD degree in the School of Automation Science and Electrical Engineering, Beihang University, Beijing, China. His research interests include software testing and software reliability analysis.



Zheng Zheng (Senior Member, IEEE) received the PhD degree in computer software and theory from the Chinese Academy of Science, Beijing, China. In 2014, he worked as a research scholar with the Department of Electrical and Computer Engineering, Duke University, Durham, North Carolina. He is currently a professor with the Beihang University, Beijing, China. His research interests include software dependability modeling, software testing, and software fault localization.



Tsong Yueh Chen (Member, IEEE) received the PhD degree from the University of Melbourne, Melbourne, Australia. He is currently a professor of software engineering with Swinburne University of Technology, Australia. Prior to joining Swinburne, he taught with The University of Hong Kong and the University of Melbourne. He is the inventor of metamorphic testing and adaptive random testing.



Pak-Lok Poon (Member, IEEE) received the PhD degree in software engineering from the University of Melbourne, Melbourne, Australia. He is an associate professor with the School of Engineering and Technology, Central Queensland University, Australia. His research interests include software testing, requirements engineering and inspection, electronic commerce, and computers in education. He was a guest editor of the special issue of the *Journal of Systems and Software* on test oracles in 2018, and will be a guest editor of the special issue of the same journal on metamorphic testing in 2021. He was also an organizer for the 3rd, 4th, and 5th International Workshops on Metamorphic Testing in 2018, 2019, and 2020, respectively.

▷ For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/csdl.