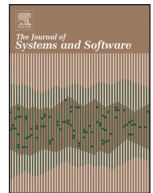




Contents lists available at ScienceDirect

## The Journal of Systems and Software

journal homepage: [www.elsevier.com/locate/jss](http://www.elsevier.com/locate/jss)METRIC: METamorphic Relation Identification based on the Category-choice framework<sup>☆</sup>Tsong Yueh Chen<sup>a</sup>, Pak-Lok Poon<sup>b,\*</sup>, Xiaoyuan Xie<sup>c</sup><sup>a</sup> Department of Computer Science and Software Engineering, Swinburne University of Technology, Hawthorn 3122, Australia<sup>b</sup> School of Business IT and Logistics, RMIT University, Melbourne 3000, Australia<sup>c</sup> State Key Laboratory of Software Engineering, Wuhan University, Wuhan, PR China

## ARTICLE INFO

## Article history:

Received 11 October 2014

Revised 20 July 2015

Accepted 23 July 2015

Available online xxx

## Keywords:

Metamorphic testing

Test oracle

Oracle problem

## ABSTRACT

Metamorphic testing is a promising technique for testing software systems when the oracle problem exists, and has been successfully applied to various application domains and paradigms. An important and essential task in metamorphic testing is the identification of metamorphic relations, which, due to the absence of a systematic and specification-based methodology, has often been done in an ad hoc manner—something which has hindered the applicability and effectiveness of metamorphic testing. To address this, a systematic methodology for identifying metamorphic relations based on the category-choice framework, called METRIC, is introduced in this paper. A tool implementing this methodology has been developed and examined in an experiment to determine the viability and effectiveness of METRIC, with the results of the experiment confirming that METRIC is both effective and efficient at identifying metamorphic relations.

© 2015 Elsevier Inc. All rights reserved.

## 1. Introduction

Many testing techniques exist for detecting software faults, with examples including random testing (Hamlet, 2002), coverage testing (Lyu et al., 1994), combinatorial testing (Nie and Leung, 2011), domain testing (White and Cohen, 1980), equivalence partitioning (Myers et al., 2004), the CHOICE reLATION fraMEwork (CHOC LATE) (Chen et al., 2003a), and the Classification-Tree Methodology (CTM) (Grochtmann and Grimm, 1993). However, except in some special situations such as detecting crashes and uncaught exceptions, application of most of these techniques requires some mechanism which can verify the correctness of the system output (Chen et al., 2003b), something known as a *test oracle* (or simply an *oracle*) in the absence of which, these testing techniques cannot be used.

When testing a software system, the *oracle problem* is said to occur when either an oracle does not exist for the tester to verify the correctness of the computed output; or an oracle does exist but cannot be used, perhaps due to feasibility or practicality issues. The oracle problem often occurs in software testing, rendering many testing techniques inapplicable (Barr et al., 2014; Hierons, 2012).

As an attempt to alleviate the oracle problem, Chen et al. (1998; 2003b) developed a technique known as metamorphic testing (MT), which involves multiple executions based on some necessary properties, formally known as *metamorphic relations* (MRs), of the software under test. The idea is that, when it is impossible or infeasible to check the correctness of output from an individual execution in the presence of the oracle problem, MT checks whether a relevant MR holds among multiple executions. If an MR does not hold, then it can be concluded that the software system is faulty.

MT has been receiving increasing attention in the software testing community (Barr et al., 2014; Gotlieb and Botella, 2003). The technique has been successfully applied in a number of application domains and paradigms, including healthcare (Murphy et al., 2011), bioinformatics (Chen et al., 2009), air traffic control (Hui et al., 2013), feature models (Segura et al., 2011), machine learning (Xie et al., 2011), spreadsheets (Poon et al., 2014), middleware (Chan et al., 2006), and the Web (Chan et al., 2007; Zhou et al., 2012). Furthermore, real-life faults have been detected by MT in some open source and comprehensively tested programs, including: the popular Siemens suites (Xie et al., 2013); a famous Machine Learning Framework (Weka) (Xie et al., 2011); feature model analysis systems (Segura et al., 2011); and compilers (Le et al., 2014; Tao et al., 2010). In fact, it was somewhat surprising that MT could detect these real-life faults, because these programs had been extensively used and tested. Detection of these faults emphasizes that MT is a new and promising test case selection strategy that complements the existing approaches (Chen et al., 2015).

<sup>☆</sup> The work described in this paper was partially supported by a Discovery Grant of the Australian Research Council (project no. ARC DP120104773).

\* Corresponding author. Tel.: +613-92144369.

E-mail addresses: [tychen@swin.edu.au](mailto:tychen@swin.edu.au) (T.Y. Chen), [drpoonpl@yahoo.com.hk](mailto:drpoonpl@yahoo.com.hk) (P.-L. Poon), [xxie@whu.edu.cn](mailto:xxie@whu.edu.cn) (X. Xie).

MT has also been used to extend the applicability of some testing and analysis methods which require an oracle. Morell's fault-based technique (Morell, 1990) assumes the presence of an oracle, but when integrated with MT, this need for an oracle is removed (Chen et al., 2003b). Spectrum-based fault localization (Jones and Harrold, 2005) is a very popular statistical approach that, until its recent integration with MT (Xie et al., 2013), also required the presence of an oracle. Xie et al. (2013) also developed a new type of program slices, metamorphic slices, which extend the applicability of program slices in some existing applications. Furthermore, MT has also been integrated with some current fault tolerance techniques to enhance their applicability to situations without oracles (Liu et al., 2014b).

In addition to integration with other testing and analysis methods, MT fundamentals and foundations have also been developing. With symbolic execution, MT has been extended from testing whether MRs are satisfied for selected inputs as test cases, to being able to prove whether MRs are satisfied for all possible inputs and to aiding debugging through identification of the constraints on inputs that violate MRs (Chen et al., 2002; 2011). Using the framework of constraint logic programming, Gotlieb and Botella (2003) have developed a technique for a specific class of programs which can find test cases leading to the violation of an MR, or prove that an MR is satisfied. Mayer and Guderlei (2006) conducted an empirical study to classify MRs into several different types using mutation analysis. Furthermore, a recently conducted study has examined how effectively MT can be used as a substitute for the test oracle (Liu et al., 2014a). These studies have enriched the theoretical and practical foundations of MT.

In spite of the successes in applying MT to multiple application domains and paradigms, the innovative integration of MT with other testing and analysis methods, and the development of the theoretical foundations for MT, concerns remain regarding its limitations in terms of MR generation. This is understandable, given that the set of identified MRs is the essence of MT (Mayer and Guderlei, 2006), and that a common opinion has been that identification of these MRs appears to be a manual and difficult task (Groce et al., 2014; Mishra and Kaiser, 2012): "[E]ven experts may find it difficult to discover metamorphic properties [relations]" (Groce et al., 2014, p.320). Although MR identification has represented one of the most challenging and core problems in MT, very limited work has been done on this problem, with current MR identification methods either requiring the availability of some MRs for identification of further ones (Kanewala and Bieman, 2013; Liu et al., 2012), or focusing on a very specific type of MRs (Zhang et al., 2014). In other words, there does not yet exist a systematic and effective identification methodology that does not require the pre-existence of some MRs and can be applied for many different types of MRs. The absence of such a methodology undoubtedly affects the applicability and effectiveness of MT. A fundamental question therefore is: **How can MRs be systematically and effectively identified?**

Because MRs represent specific properties of the developed software, they actually form part of the software requirements (either explicitly or implicitly), which will then be included in the software specification. Thus, intuitively, a specification document should contain important and useful information that can be used for identifying the MRs. Accordingly, a specification-based methodology has been developed: **METAMorphic Relation Identification** based on the **Category-choice framework** (METRIC). METRIC is built within the specification-based category-choice framework, which includes the concepts of categories, choices, and complete test frames. **Examples of this framework include CHOC'LATE and CTM, and their variants** (Chen et al., 2012; 2003a; Grochtmann and Grimm, 1993; Hierons et al., 2003; Singh et al., 1997). METRIC enables the tester to systematically identify MRs from complete test frames generated by the category-choice framework to test a software system, of which the oracle problem exists (that is, an oracle cannot be derived from the specification or the derived oracle cannot be feasibly applied). In

this paper, our METRIC methodology is not a replacement or enhancement of the category-choice framework. Instead, METRIC makes use of the complete test frames generated by the category-choice framework to identify MRs for MT.

The rest of this paper is organized as follows: **Sections 2 and 3** introduce the concepts of MT/MR and the category-choice framework, respectively, both of which are essential for understanding METRIC. **Section 4** explains the motivations of the study and discusses details of METRIC. **Section 5** describes the main features of an associated generator tool developed to support METRIC. **Sections 6 and 7** describe the setting and results of an experiment to demonstrate the viability and effectiveness of METRIC. **Section 8** further elaborates this research, discussing limitations, and highlighting potential future work. **Section 9** discusses some related work done by other researchers, and finally, **Section 10** summarizes and concludes the paper.

## 2. Metamorphic testing (MT) and metamorphic relation (MR)

MT is an innovative approach to alleviating the oracle problem in testing (Chen et al., 1998; 2003b): Instead of relying on the existence and applicability of an oracle, MT uses a set of MRs (corresponding to some specific properties of the software) to generate test cases and verify test results. Chen et al. (1998) define an *oracle* as a mechanism or procedure that can decide whether or not an output is correct. Their definition follows those of Weyuker (1982), and Peters and Paras (1998). This paper adopts their definition of an oracle. Readers are reminded, however, that some researchers define an oracle in a slightly different manner (for example, see Barr et al. (2014)).

To outline the main concepts of MT and MR, assume a program  $P$ , and its input domain  $D$ . The size of  $D$  is often very large, or even infinite, making exhaustive testing infeasible. In such a situation, a test suite  $U \subseteq D$  may be constructed. Suppose  $U = \{t_1, t_2, \dots, t_n\}$ , where  $t_i$  is a test case for any  $i = 1, 2, \dots, n$ , and  $P(t_i)$  represents the result of executing  $P$  with  $t_i$ . If an oracle exists, and can be feasibly applied, the correctness of  $P(t_i)$  can be verified. The following example illustrates the concepts and the process of MT.

**Example 1 (SINE function).** Consider an implemented SINE function. Except for some special values such as  $90^\circ$  ( $\sin(90^\circ) = 1$ ), the correctness for the outputs of SINE is difficult, if not impossible, to be verified. Some properties, however, do exist for the sine function, including that for any angle  $x$  (in degrees),  $\sin(x) = \sin(x + 360^\circ)$ . This property then forms an MR in MT: "If  $y = x + 360^\circ$ , then  $\sin(x) = \sin(y)$ ". Using this MR, the implemented SINE should be executed twice: firstly with any angle  $x$  as a *source* test case; and then with the angle  $y$ , such that  $y = x + 360^\circ$ , as a *follow-up* test case. The rationale is that, although it may be very difficult to check the correctness of  $\text{SINE}(x)$  in one single execution (unless  $x$  is a special angle such as  $90^\circ$ ), executing SINE again with another angle  $y = x + 360^\circ$  and (automatically) checking whether  $\text{SINE}(y) = \text{SINE}(x)$  can be easily done.

For example, suppose SINE is executed with a test case  $x = 69^\circ$ , giving an output of 0.9336 (i.e.,  $\text{SINE}(69^\circ) = 0.9336$ ). In MT, with respect to the above MR corresponding to the property  $\sin(x) = \sin(x + 360^\circ)$ , SINE should next be executed with another test case  $y = 69^\circ + 360^\circ = 429^\circ$ . The output of the second execution should be compared with that of the first (after allowing for rounding error): Does  $\text{SINE}(429^\circ) = 0.9336$ ? If the equality does not hold, then it can be concluded that SINE is faulty. Note that the inequality only indicates that the MR "If  $y = x + 360^\circ$ , then  $\sin(x) = \sin(y)$ " has been violated and, hence, SINE is faulty. The exact values of  $\sin(69^\circ)$  and  $\sin(429^\circ)$  are still unknown. This explains that MRs and oracles are not the same.

Example 1 shows that, rather than focusing on the correctness of output from one single execution, MT checks whether a relevant MR holds among multiple executions. It is this characteristic that makes

MT very useful in testing, even when the oracle problem occurs. Formal definitions for a metamorphic relation (MR) and for metamorphic testing (MT) are as follows:

**Definition 1 (Metamorphic relation).** Let  $f$  be a target algorithm and  $P$  be its corresponding implementation. A **metamorphic relation (MR)** is a necessary property of  $f$  over a series of two or more inputs  $(x_1, x_2, \dots, x_n)$ , where  $n \geq 2$  and their corresponding outputs  $(f(x_1), f(x_2), \dots, f(x_n))$ , that is, a relation  $\mathbb{R}(x_1, x_2, \dots, x_n, f(x_1), f(x_2), \dots, f(x_n))$  involving multiple inputs and their corresponding outputs.

Without loss of generality, assume there exist  $k$  ( $1 \leq k < n$ ) inputs such that, for each  $x_j$  ( $(k+1) \leq j \leq n$ ), we have a mapping  $f_j$  to compute  $x_j$  such that  $x_j = f_j(x_1, x_2, \dots, x_k, f(x_1), f(x_2), \dots, f(x_k))$ . For  $1 \leq i \leq k$ ,  $x_i$  is called a **source test case**; and for  $(k+1) \leq j \leq n$ ,  $x_j$  is called a **follow-up test case**. In other words, if all  $x_i$  ( $1 \leq i \leq k$ ) are specified, then  $x_j$  ( $(k+1) \leq j \leq n$ ) can be determined from the source test cases and their corresponding outputs.

**Definition 2 (Metamorphic testing).** Let  $P$  be an implementation of a target algorithm  $f$ . Consider an MR:  $\mathbb{R}(x_1, x_2, \dots, x_n, f(x_1), f(x_2), \dots, f(x_n))$ . **Metamorphic testing** of this MR for  $P$  involves the following steps: (1) Given a series of source test cases  $\langle x_1, x_2, \dots, x_k \rangle$  and their respective outputs  $\langle P(x_1), P(x_2), \dots, P(x_k) \rangle$ , generate a series of follow-up test cases  $\langle x_{k+1}, x_{k+2}, \dots, x_n \rangle$  according to the mapping  $x_j = f_j(x_1, x_2, \dots, x_k, f(x_1), f(x_2), \dots, f(x_k))$ , but with  $f(x_i)$  being replaced by  $P(x_i)$ , where  $1 \leq i \leq k$  and  $(k+1) \leq j \leq n$ . (2) Check the relation  $\mathbb{R}(x_1, x_2, \dots, x_n, f(x_1), f(x_2), \dots, f(x_n))$ , but with  $f(x_i)$  being replaced by  $P(x_i)$ , where  $1 \leq i \leq n$ . If  $\mathbb{R}$  is not satisfied, then metamorphic testing of this MR reveals that  $P$  is faulty.

In MT, for any given test case  $x_j$  ( $1 \leq j \leq n$ ), there is no need to investigate whether  $P(x_j) = f(x_j)$  by means of a test oracle. This therefore alleviates the oracle problem in testing.

### 3. Category-choice framework

As mentioned in Section 1, METRIC is built within the category-choice framework, which involves categories, choices, and complete test frames. Before presenting the details of METRIC, CHOC'LATE (a category-choice framework) is first used to illustrate the concepts of categories, choices, and complete test frames.

CHOC'LATE provides a systematic framework for generating test cases from specifications using categories, choices, and complete

test frames (Chen et al., 2012; 2003a). By partitioning the input domain according to choices, CHOC'LATE generates a set of complete test frames based on the constraints among choices, formally known as **choice relations**. The complete test frames thus generated, which are valid combinations of choices, then form the basis for test case generation. If an instance is selected from each choice in a complete test frame, a test case is formed. In this paper, the notations  $\mathbb{B}$  and  $B$  are used to denote a “generic” complete test frame and a “specific” complete test frame, respectively.

**Example 2 (Parking fee system).** Consider a parking fee calculation system FEE, which accepts the parking details of each vehicle, including type of vehicle, type of car, day of week, discount coupon, and hours of parking. FEE first rounds up the parking duration to the next full hour, and then calculates the parking fee for a vehicle according to the hourly rates in Table 1. If a discount coupon is presented, a 50% discount off the parking fee will be given.

To facilitate better parking management, at the time of parking, customers can optionally provide an estimation of parking duration in terms of three different time ranges, namely (0.0, 2.0], (2.0, 4.0], and (4.0, 24.0]. Suppose a customer provides an estimation. If the estimated hours of parking and the actual hours of parking fall into the same time range, then the customer will receive a 40% discount. If, however, the estimated hours and the actual hours are in different time ranges, a 20% markup will be added. A customer can choose to either use a discount coupon, or provide an estimation of parking duration, but not both. Obviously, a customer may also choose to neither provide an estimation, nor use a discount coupon. Note that no vehicles are allowed to park across two consecutive days on a continuous basis. Suppose, for instance, that a vehicle parks at 23:30 on a Monday, then it must vacate the car park not later than 23:59 on that same day.

**Step 1.** Categories and choices are identified from the FEE specification. A **category** is a major property or characteristic of an environment condition, or a parameter of the system that affects its execution behavior. The possible values associated with each category are partitioned into disjoint subsets known as **choices**, with the assumption that all values in the same choice are similar either in their effect on the system's behavior, or in the type of output that they produce. Given a category  $P$ , the notation  $P_x$  is used to denote a choice of  $P$ . Table 2 shows the possible categories and their associated choices for FEE.

**Table 1**  
Hourly parking rates.

Actual hours of parking	Weekday			Saturday and Sunday		
	Motorcycle (\$)	Car: 2-door coupe (\$)	Car: others (\$)	Motorcycle (\$)	Car: 2-door coupe (\$)	Car: others (\$)
(0.0, 2.0]	4.00	4.50	5.00	5.00	6.00	7.00
(2.0, 4.0]	5.00	5.50	6.00	6.50	7.50	8.50
(4.0, 24.0]	6.00	6.50	7.00	8.00	9.00	10.00

**Table 2**  
Categories and choices for FEE.

Categories	Associated choices
type of vehicle	type of vehicle <sub>motorcycle</sub> , type of vehicle <sub>car</sub>
type of car	type of car <sub>2-door coupe</sub> , type of car <sub>others</sub>
day of week	day of week <sub>weekday</sub> , day of week <sub>Saturday or Sunday</sub>
discount coupon	discount coupon <sub>yes</sub> , discount coupon <sub>no</sub>
estimated hours of parking	estimated hours of parking <sub>(0.0, 2.0]</sub> , estimated hours of parking <sub>(2.0, 4.0]</sub> , estimated hours of parking <sub>(4.0, 24.0]</sub>
actual hours of parking	actual hours of parking <sub>(0.0, 2.0]</sub> , actual hours of parking <sub>(2.0, 4.0]</sub> , actual hours of parking <sub>(4.0, 24.0]</sub>



Step 2. Constraints are identified for each pair of choices, from which valid combinations are selected as complete test frames. For FEE, given the categories and choices defined in Table 2, exactly 90 complete test frames will be generated, an example of which is:  $B_1 = \{\text{type of vehicle}_{\text{motorcycle}}, \text{day of week}_{\text{weekday}}, \text{discount coupon}_{\text{no}}, \text{actual hours of parking}_{(0.0, 2.0)}\}$ .

Step 3. Test cases can be formed from each complete test frame  $\mathbb{B}$  generated in step 2, by selecting and combining an instance from each choice in  $\mathbb{B}$ . For example, taking  $B_1$  in step 2, a test case  $tc = \{\text{type of vehicle} = \text{motorcycle}, \text{day of week} = \text{Monday}, \text{discount coupon} = \text{no}, \text{actual hours of parking} = 1.4\}$  can be formed.

3. The identification exercise is repeated by going back to step 1, involving the selection of another concrete input as a new source test case.
4. If the  $R$  defined in step 2 can be generalized to a class of source test cases, then an MR is identified. Otherwise, repeat the entire process by starting from step 1, involving the selection of another concrete input as a new source test case.

#### 4.2. Methodology

The category-choice framework (including categories, choices, and complete test frames) was used to develop a new methodology to partly automate the manual MR identification pattern in Section 4.1. Reasons for this include:

- a. As explained in Section 1, MRs should have already been explicitly or implicitly included as part of the specification, from which it should be possible to extract the relevant information for MR identification. Because the category-choice framework is specification-based, it is a good candidate for the proposed methodology.
- b. In MT, each MR represents the effect on the output when an input is altered in some controlled way, therefore MRs focus on the input domain of the system. The category-choice framework also focuses on the input domain by dividing it into several disjoint partitions (known as complete test frames), from which test cases can be constructed.
- c. The category-choice framework aims at generating a set of complete test frames from which concrete test cases can be constructed. Each complete test frame is a valid combination of choices, and corresponds to an input (test) scenario rather than a single concrete input (a test case). Complete test frames are therefore more abstract than concrete test cases, and as such, an MR derived from a complete test frame is applicable to *multiple* source test cases and their corresponding follow-up test cases. This makes complete test frames a natural basis for MR identification: a relation derived from concrete test cases corresponds to an  $R$ , and one derived from complete test frames corresponds to an MR (see step 4 of Section 4.1). Example 5 (later in this section) will illustrate the advantages of using complete test frames to identify MRs from which multiple source and follow-up test cases can be generated.

The METRIC (METamorphic Relation Identification based on the Category-choice framework) methodology focuses on MRs such that each MR involves two distinct inputs, with the associated follow-up input depending only on the corresponding source input  $I$  (not on the output of  $I$ ). Consequently, only two distinct complete test frames need be considered by the tester to generate an MR, leading to the concept of a candidate pair, which serves as a cornerstone of METRIC.

In general, METRIC involves the following steps:

- Step 1. Select two relevant and distinct complete test frames as a *candidate pair* for user's consideration.
- Step 2. Enable the user to effectively and efficiently determine whether or not the selected candidate pair is useful for MR identification, and if it is, then allow the user to provide the corresponding MR description.
- Step 3. Restart from step 1, and repeat until all candidate pairs are exhausted, or the predefined number of MRs to be generated is reached.

For the remainder of this paper, candidate pairs that have associated MRs are called *usable*, and those that do not have are *unusable*.

In a typical MR identification process, testers have to *simultaneously* consider both the input and output aspects of the system (Section 4.1), which has been reported by software researchers and practitioners as a difficult task of MR identification. In view of this

## 4. METRIC methodology

### 4.1. A motivating example

When testing the SINE function in Example 1, MRs can be easily identified because the properties of this function, such as " $\sin(x) = \sin(x + 360^\circ)$ ", are well known in mathematics. However, knowledge of specific properties of many applications, especially proprietary software, may not always be explicitly available, thus making ad hoc MR identification for these applications very difficult (Groce et al., 2014). This problem has represented an obstacle to the applicability and effectiveness of MT.

Previous work and experience with MT have revealed that many software practitioners adopt a similar approach for manual MR identification, as illustrated in the following example.

**Example 3 (Online hotel availability inquiry system).** Consider a Web-based hotel availability inquiry system INQ, which outputs a list of hotels satisfying certain search criteria. Each inquiry includes input parameters such as city name, the maximum distance of the hotel from the center of the city, check-in date, check-out date, the minimum hotel category, and room type. For instance, a possible input  $I$  is: "Melbourne/10 km/Jan 3, 2014/Jan 8, 2014/5-star/single". Suppose, given  $I$ , INQ returns the output  $O = \{\text{Hotel-A}, \text{Hotel-B}\}$ . Although verifying whether or not Hotel-A and Hotel-B are correctly identified by INQ may be easily done (by calling both hotels, for example), because of the potentially large number of hotels involved, it could be very difficult to check whether any other eligible hotels were omitted by INQ. In other words, it is not practically feasible to check the correctness of  $O$ , resulting in an occurrence of the oracle problem for testing INQ.

When identifying an MR, a common approach is to consider how to *adjust* a given input in order to get a *predictable change* in output. For example, if the minimum hotel category is changed from "5-star" to "4-star", but all other input parameters remain unchanged, then the expected output should still include Hotel-A and Hotel-B, but possibly with some other eligible hotels. Thus, an MR can be generalized as follows: "If the minimum hotel category is lowered, the output for the old input will be a subset of that for the new input".

The identification pattern described in Example 3 can be summarized as the following "bottom-up" stepwise approach:

1. With respect to the implemented algorithm, select a concrete input  $I$  as a source test case.
2. Determine how to change  $I$  to another concrete input  $I'$  (a follow-up test case), such that a relation  $R$  can be defined among  $I$ ,  $I'$ , and their corresponding outputs. This step may be iterated until  $R$  is found (in which case the tester goes to step 4), or, after several failed attempts, the tester judges that  $R$  does not exist with  $I$  as the source test case (in which case the tester goes to step 3).

problem, METRIC provides testers with one pair ( $\mathbb{B}_1$  and  $\mathbb{B}_2$ , which correspond to the inputs only) of complete test frames at a time, thereby enabling them to concentrate on evaluating whether a relation exists between the outputs associated with  $\mathbb{B}_1$  and  $\mathbb{B}_2$ . If such a relation exists, then an MR can be formed.

By definition, any pair ( $\mathbb{B}_1$  and  $\mathbb{B}_2$ ) of distinct complete test frames must contain some different choices. Either (a)  $\mathbb{B}_1$  and  $\mathbb{B}_2$  contain different choices of the same category, or (b) one or more categories associated with  $\mathbb{B}_1$  are not associated with  $\mathbb{B}_2$  (or vice versa). The different choices (due to reason (a)) and extra choices (due to reason (b)) are referred to as *differentiating choices*, and the categories associated with these differentiating choices are referred to as *differentiating categories*. A pair of differentiating choices associated with the same category but which appears separately in  $\mathbb{B}_1$  and  $\mathbb{B}_2$  is called a *differentiating choice pair*.

A candidate pair ( $\mathbb{B}_1$  and  $\mathbb{B}_2$ ) can be classified into one of three types: a **choice-only pair (choice pair)**; a **category-only pair (category pair)**; or a **category-choice pair**. For a choice-only pair, the difference in choices between  $\mathbb{B}_1$  and  $\mathbb{B}_2$  is due only to reason (a); for a category-only pair, the difference is due only to reason (b); and a category-choice pair is a pair which is neither choice-only nor category-only. A choice, category, or category-choice pair ( $\mathbb{B}_1$  and  $\mathbb{B}_2$ ) can be written as a choice( $k$ ), category( $j$ ), or category-choice( $j, k$ ) pair, respectively, where  $j$  is the number of differentiating categories that are associated with either  $\mathbb{B}_1$  or  $\mathbb{B}_2$  (but not both), and  $k$  is the number of differentiating categories with different choices in both  $\mathbb{B}_1$  and  $\mathbb{B}_2$ .

**Example 4 (Types of candidate pairs).** Referring again to  $B_1$  in Example 2 (reproduced below), consider the following complete test frames for FEE:

$$\begin{aligned} B_1 &= \{\text{type of vehicle}_{\text{motorcycle}}, \text{day of week}_{\text{weekday}}, \text{discount coupon}_{\text{no}}, \text{actual hours of parking}_{(0.0, 2.0]}\}; \\ B_2 &= \{\text{type of vehicle}_{\text{motorcycle}}, \text{day of week}_{\text{weekday}}, \text{discount coupon}_{\text{no}}, \text{actual hours of parking}_{(2.0, 4.0]}\}; \\ B_3 &= \{\text{type of vehicle}_{\text{motorcycle}}, \text{day of week}_{\text{weekday}}, \text{discount coupon}_{\text{no}}, \text{estimated hours of parking}_{(2.0, 4.0]}, \\ &\quad \text{actual hours of parking}_{(2.0, 4.0]}\}; \text{ and} \\ B_4 &= \{\text{type of vehicle}_{\text{motorcycle}}, \text{day of week}_{\text{weekday}}, \text{discount coupon}_{\text{no}}, \text{actual hours of parking}_{(4.0, 24.0]}\}. \end{aligned}$$

Consider first  $B_1$  and  $B_2$ , which are identical other than having different choices in the category “actual hours of parking” (“actual hours of parking<sub>(0.0, 2.0]</sub>”  $\in B_1$  and “actual hours of parking<sub>(2.0, 4.0]</sub>”  $\in B_2$ ). Thus, “actual hours of parking” is the only differentiating category; and “actual hours of parking<sub>(0.0, 2.0]</sub>” and “actual hours of parking<sub>(2.0, 4.0]</sub>” are the differentiating choice pair. Therefore,  $B_1$  and  $B_2$  form a choice(1) pair.

Next, consider  $B_2$  and  $B_3$ . Because the choice “estimated hours of parking<sub>(2.0, 4.0]</sub>” is contained in  $B_3$  but not in  $B_2$ , the category “estimated hours of parking” is associated with  $B_3$  but not with  $B_2$ . Thus,  $B_2$  and  $B_3$  form a category (1) pair, where “estimated hours of parking” is the differentiating category, and “estimated hours of parking<sub>(2.0, 4.0]</sub>” is the differentiating choice.

Finally, consider  $B_3$  and  $B_4$ , which form a category-choice(1, 1) pair, where “estimated hours of parking” and “actual hours of parking” are the differentiating categories; and “estimated hours of parking<sub>(2.0, 4.0]</sub>”, “actual hours of parking<sub>(2.0, 4.0]</sub>”, and “actual hours of parking<sub>(4.0, 24.0]</sub>” are the differentiating choices. Furthermore, “actual hours of parking<sub>(2.0, 4.0]</sub>” and “actual hours of parking<sub>(4.0, 24.0]</sub>” are the differentiating choice pair.

**Example 5 (identification of MRs from candidate pairs).** Consider  $B_1$  and  $B_2$  in Example 4, which are identical except in one differenti-

ating choice pair, namely “actual hours of parking<sub>(0.0, 2.0]</sub>” and “actual hours of parking<sub>(2.0, 4.0]</sub>”. Given the information that a longer time range of actual parking duration (related to the system input) will incur a higher parking fee (corresponding to the system output) for the same type of vehicle and day of week, if the customer does not possess a discount coupon and does not provide an estimation of parking duration, then, when comparing  $B_1$  and  $B_2$ , the parking fee for  $B_2$  must be higher than that for  $B_1$ . Therefore,  $B_1$  and  $B_2$  are usable for defining the following MR:

(MR<sub>1</sub>) Assume that a motorcycle parks on a weekday, and the customer neither possesses a discount coupon nor provides an estimation of parking duration upfront. If the time range of actual parking duration increases from (0.0, 2.0] to (2.0, 4.0], the parking fee will increase.

Next, consider  $B_2$  and  $B_3$  in Example 4.  $B_3$  is associated with the differentiating category “estimated hours of parking”, but  $B_2$  is not. Because the scenario corresponding to  $B_3$  will receive a 40% discount (an estimation of parking duration was provided upfront and it falls into the same time range as the actual hours of parking), the final parking fee for  $B_3$  must be lower than that for  $B_2$ . In this regard,  $B_2$  and  $B_3$  are usable, and lead to the following MR:

(MR<sub>2</sub>) Assume that a motorcycle parks on a weekday with an actual parking duration of (2.0, 4.0] and the customer does not possess a discount coupon. The provision of an estimation of parking duration of (2.0, 4.0] upfront will result in a lower parking fee than not providing such an estimation.

Finally, consider  $B_3$  and  $B_4$  in Example 4. The parking fee for  $B_3$  must be lower than  $B_4$ , because: (a)  $B_3$  has a correct range of estimation of parking duration, resulting in a 40% discount off the parking fee; and (b)  $B_3$  has a shorter actual parking duration than  $B_4$ , resulting in a lower hourly rate and a lower parking fee (even before applying the discount). Thus,  $B_3$  and  $B_4$  are usable for defining the following MR:

(MR<sub>3</sub>) Assume that a motorcycle parks on a weekday and the customer neither possesses a discount coupon nor provides an estimation of parking duration upfront. If the time range of actual parking duration changes from (4.0, 24.0] to (2.0, 4.0], and an estimation of parking duration of (2.0, 4.0] is provided, the parking fee will be lowered.

Note that not all candidate pairs are usable. Consider  $B_3$  in Example 4 and  $B_5 = \{\text{type of vehicle}_{\text{car}}, \text{type of car}_{\text{others}}, \text{day of week}_{\text{weekday}}, \text{discount coupon}_{\text{no}}, \text{actual hours of parking}_{(0.0, 2.0]}\}$ . They form a category-choice(2, 2) pair, but no relationship exists between their corresponding outputs (the computed parking fees). The following three test cases illustrate this situation, where  $tc_1$  is generated from  $B_3$ , and  $tc_2$  and  $tc_3$  are generated from  $B_5$ .

$$\begin{aligned} tc_1 &= \{\text{type of vehicle} = \text{motorcycle}, \text{day of week} = \text{Monday}, \\ &\quad \text{discount coupon} = \text{no}, \text{estimated hours of parking} = 3.0, \\ &\quad \text{actual hours of parking} = 2.5\}; \\ tc_2 &= \{\text{type of vehicle} = \text{car}, \text{type of car} = 4 - \text{door sedan}, \\ &\quad \text{actual hours of parking} = 0.5\}; \text{ and} \\ tc_3 &= \{\text{type of vehicle} = \text{car}, \text{type of car} = 4 - \text{door sedan}, \\ &\quad \text{day of week} = \text{Monday}, \text{discount coupon} = \text{no}, \\ &\quad \text{actual hours of parking} = 2.0\}. \end{aligned}$$

According to the hourly parking rates in Table 1, the parking fees for  $tc_1$ ,  $tc_2$ , and  $tc_3$  are \$9.00 (= \$5.00/h  $\times$  3.0 h  $\times$  60%), \$5.00 (= \$5.00/h  $\times$  1.0 h), and \$10.00 (= \$5.00/h  $\times$  2.0 h), respectively. These parking fees are calculated based on the following points, as stated earlier in the paper: (a) parking durations will be rounded up to the

next full hour and (b) a 40% discount will be given if the estimated hours of parking and the actual hours of parking fall into the same time range. Because the parking fee for  $tc_1$  ( $= \$9.00$ ) is higher than  $tc_2$  ( $= \$5.00$ ), but lower than  $tc_3$  ( $= \$10.00$ ), there is no relationship between the computed parking fees for  $B_3$  and  $B_5$ , and therefore  $B_3$  and  $B_5$  are unusable for defining any MR.

For any candidate pair  $B_1$  and  $B_2$ , their differentiating categories indicate the difference between their concrete source and follow-up test cases, and their non-differentiating categories indicate the commonality between them. Thus, to facilitate MR identification, when a candidate pair has some common choices, the *same* value should be used for each common choice. For instance, in test cases  $tc_1$ ,  $tc_2$ , and  $tc_3$  of Example 5, “Monday” is selected as the value from the choice “day of week<sub>weekday</sub>” in  $B_3$  and  $B_5$ .

An advantage of using complete test frames for MR identification (reason (c) in Section 4.2) is illustrated by  $MR_1$  in Example 5, which allows the tester to generate a number of source and follow-up test cases for MT. Example 5 also clearly and convincingly shows that candidate pairs provide a natural and central basis upon which possible MRs can be systematically and effectively identified.

## 5. An associated generator tool

In order to support and partly automate the three steps involved in METRIC (Section 4.2), a Java-based tool called **MR-GENerator** (MR-GEN) has been developed. MR-GEN does the following: (a) takes a spreadsheet file containing a set of complete test frames (TF), which is generated by a tool such as the one described in Chen et al. (2003a); (b) allows the user to specify the type of candidate pairs for evaluation; (c) displays selected candidate pairs for the user to determine if they are usable; (d) records the user’s MR description for usable candidate pairs; and (e) outputs an HTML file containing the identified MRs. Other facilitating features of MR-GEN include, for example, postponing evaluation of certain candidate pairs, and allowing the user to modify the currently displayed candidate pairs to control the selection of other candidate pairs. These steps and features are next illustrated using the FEE system from Example 2.

Using the categories and choices in Table 2, a set of 90 complete test frames for FEE, denoted  $TF_{FEE}$ , were generated and stored in a spreadsheet file, and then input to MR-GEN. MR-GEN prompts the user to specify the type of candidate pairs and other associated details as shown in Fig. 1. Often, resource constraints impose a limit on the number of test cases that can be executed, thus, MR-GEN provides a parameter  $M$ , to specify the *maximum* number of MRs to be identified from the TF. MR-GEN also allows the user to specify a set of critical categories  $\{C_1, C_2, \dots, C_n\}$  (where  $C_i$  is a category,  $1 \leq i \leq n$ ) such that each  $C_i$  must be associated with at least one of the two complete test frames in any selected candidate pair (Fig. 2). A *critical* category is one judged by the user to play a significant role in the functionality of the software under test. As a further option, the user may specify

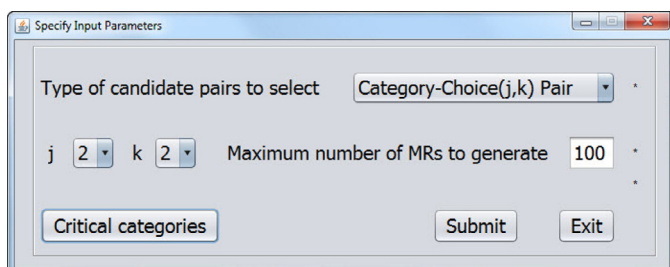


Fig. 1. Input screen for specifying parameters.

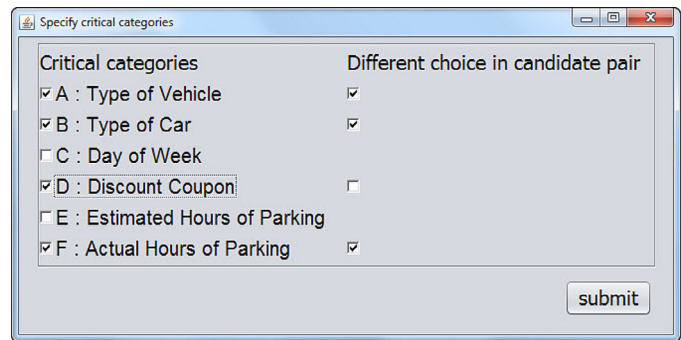


Fig. 2. Selection screen for critical categories.

whether each selected candidate pair must have different choices in a particular  $C_i$  (Fig. 2).

After the user has specified the above parameters and critical categories, MR-GEN displays the relevant candidate pairs, one at a time, to be judged usable or not (the upper half of Fig. 3). A combo box allows the user to indicate whether or not the pair is usable for identifying an MR. If usable, the MR details are recorded. If the user is not certain whether or not the candidate pair is usable, the pair can be labeled as “yet-to-be-decided”, and the pair can be revisited later. After evaluating other candidate pairs, the user may accumulate sufficient knowledge to judge whether or not a “yet-to-be-decided” pair is usable.

Because MR-GEN presents candidate pairs based on the parameters entered in Figs. 1 and 2, the user only needs to determine if a relation exists between the corresponding *outputs*. The user does not need to select complete test frames (corresponding to the *input* aspects) for consideration, and thus a significant portion of the MR identification burden is alleviated (see the discussion on the problem of simultaneous consideration of both the input and output aspects of the system in Section 4.2). Furthermore, to facilitate comparison, MR-GEN presents the complete test frames, side-by-side, for each candidate pair, highlighting differentiating choices in different colors (Fig. 3). This arrangement is obviously more useful than relying only on the user’s own ability to select some related inputs from a large pool for possible MR identification.

After identifying some MRs with MR-GEN, the user may develop a better acquaintance with the software under test, and may know what changes to the current candidate pair’s choices would lead to a relation between outputs of a new pair. To cater for this, MR-GEN provides an “on-demand selection” option to modify the current candidate pairs, with a view to obtaining a new, but usable candidate pair (Fig. 4). Consider an example: Although  $B_3$  and  $B_5$  in Example 5 are not usable, the user may know that by replacing the choice “actual hours of parking<sub>{0.0,2.0}</sub>” in  $B_5$  with another choice “actual hours of parking<sub>{2.0,4.0}</sub>” to form a new complete test frame  $B_6$ , a relationship exists between the outputs associated with  $B_3$  and  $B_6$ —the parking fee for  $B_6$  must be higher than  $B_3$ —and therefore  $B_3$  and  $B_6$  are usable. Note that the initial candidate pair which serves as the basis for on-demand selection can be usable or unusable.

## 6. Experimental setting

### 6.1. Research questions, specifications, and participants

Although METRIC is intuitive, its viability and effectiveness should be verified from the tester’s perspective, leading to the following primary research question: (RQ1) **Can METRIC and its associated tool MR-GEN be practically and feasibly used for MR identification?** This research question can be decomposed into four sub-research



Define metamorphic relations

	Complete Test Frame 41	Complete Test Frame 67
Category		
A : Type of Vehicle	a1: Motorcycle	a2: Car
B : Type of Car		b2: Others
C : Day of Week	c1: Weekday	c1: Weekday
D : Discount Coupon	d2: No	d2: No
E : Estimated Hours of Parking	e2: $2.0 < f \leq 4.0$	
F : Actual Hours of Parking	f2: $2.0 < f \leq 4.0$	f1: $0.0 < f \leq 2.0$

Is the above candidate pair usable?

Name of this MR

Definition of this MR

MR ID: -; Created time: 2014-06-21-13:59:22; Last modified: 2014-06-21-13:59:22

Intelligent selection of categories/choices

Preferred maximal number of MRs: 100

Total number of MRs defined: 0

Progress: 1 / 144 pairs

Fig. 3. Screen for evaluating candidate pairs and providing MR descriptions.

Intelligent selection of categories/choices

	Complete Test Frame 41	Complete Test Frame 67
Category		
A : Type of Vehicle	a1: Motorcycle	a2: Car
B : Type of Car	N/A	b2: Others
C : Day of Week	c1: Weekday	c1: Weekday
D : Discount Coupon	d2: No	d2: No
E : Estimated Hours of Parking	e2: $2.0 < f \leq 4.0$	N/A
F : Actual Hours of Parking	f2: $2.0 < e \leq 4.0$	f1: $0.0 < e \leq 2.0$

Is the above candidate pair usable?

Name of this MR

Definition of this MR

Fig. 4. Screen for on-demand selection of subsequent candidate pairs.

questions: (RQ1.1) Does use of METRIC result in a significant difference in MR identification compared with the ad hoc approach? (RQ1.2) Is METRIC easy to use for users with only some primitive knowledge about MT? (RQ1.3) Does previous experience with MT have a significant impact on the user's MR identification when using METRIC? (RQ1.4) Does increased familiarity with METRIC and MR-GEN result in a significant improvement in the user's MR identification?

Intuitively, METRIC should ease the difficulties encountered by testers in MR identification. Furthermore, the users, especially those familiar with the specifications and MR-GEN, can use the tool's on-demand feature to select particular candidate pairs for later evaluation, leading to the second research question: (RQ2) Does use of MR-GEN's on-demand selection feature result in a significant improvement in MR identification?

Two commercial specifications were used in the experiment:  $SCAR$ , which is related to a company car and expense claim system (CAR) for a multinational trading firm (Chen et al., 2012); and  $S_{MOS}$  which is related to a meal ordering system (MOS) for an international company providing catering services for many different airlines (Chen et al., 2003a). The experiment involved 19 participants (denoted  $P_1, P_2, \dots, P_{19}$ ), who were postgraduate IT students or who had earned a PhD degree in IT. The participants were classified into two broad groups: the *inexperienced* (I-group), who had no knowledge of MT before the experiment (ten participants,  $P_1$  to  $P_{10}$ ); and the *experienced* (E-group) who have acquired prior MT knowledge through previous research work related to MT (nine participants,  $P_{11}$  to  $P_{19}$ ). Furthermore, the I-group and the E-group were randomly and evenly divided into two (I-group1 and I-group2) and three (E-group1, E-group2, and E-group3) subgroups, respectively.

**Table 3**  
MR identification arrangements for different parts of the experiment.

Participant groups	Participants	Part I:	Part II: using MR-GEN without on-demand selection		Part III: using MR-GEN with on-demand selection	
		Ad hoc Within 90 min	Ila For 40 MRs	Ilb Within 90 min Specification(s) used	IIla For 40 MRs	IIlb within 90 min
I-group1	P <sub>1</sub> to P <sub>5</sub>	—	$S_{CAR}$	$S_{CAR}$	—	—
I-group2	P <sub>6</sub> to P <sub>10</sub>	—	$S_{MOS}$	$S_{MOS}$	—	—
E-group1	P <sub>11</sub> to P <sub>13</sub>	$S_{MOS}$	$S_{CAR}$	$S_{CAR}$	—	—
E-group2	P <sub>14</sub> to P <sub>16</sub>	$S_{CAR}$	$S_{MOS}$	$S_{MOS}$	—	—
E-group3	P <sub>17</sub> to P <sub>19</sub>	—	—	—	$S_{CAR}$ & $S_{MOS}$	$S_{CAR}$ & $S_{MOS}$

## 6.2. Experimental preparation and components

Three tutorials were organized to prepare the participants for the experiment. **Tutorial 1 (Specifications)** introduced the  $S_{CAR}$  and  $S_{MOS}$  specifications to all participants. **Tutorial 2 (MT)** was provided to the inexperienced participants only. It introduced the concept of MT and included two hands-on exercises to reinforce understanding. **Tutorial 3 (Category-choice framework and MR-GEN)** was provided to all participants and had two purposes, the first of which was to introduce the concepts of categories, choices, and complete test frames. This was done by revisiting  $S_{CAR}$  and  $S_{MOS}$ , and examining some complete test frames for CAR and MOS, generated by a tool developed for CHOC'LATE (Chen et al., 2003a). The second purpose of this tutorial was to teach the functionality and operation of MR-GEN, including the evaluation of the usability of candidate pairs for MR identification, illustrated by using the FEE system described in Example 2. Tutorials 1 and 2 each lasted 1 h, and Tutorial 3 lasted for 90 min.

The experiment was divided into three parts as follows:

**Part I (Ad hoc MR identification)** took place immediately after Tutorial 1, and involved only groups E-group1 and E-group2, who were asked to independently identify as many MRs as possible for  $S_{MOS}$  and  $S_{CAR}$ , respectively, within a time limit of 90 min, using an ad hoc approach. The duration of 90 min was set to avoid participants becoming exhausted from performing manual MR identification for too long a period, but in the actual experiment it was observed that, even after about 1 h, some participants could not identify any more MRs.

**Part II (MR identification supported by MR-GEN without on-demand selection)** was conducted after Tutorial 3, and involved all groups except E-group3. This part was further divided into two phases, with a 1-h break between them. In the first phase (Part IIa), each participant independently identified 40 MRs from either  $S_{CAR}$  or  $S_{MOS}$ , without any time constraint: because this was their first time to use MR-GEN, the participants would not yet be very familiar with the tool, and therefore by not setting a time limit, the participants were expected to focus more on correctness than on speed of performing the task. Note that E-group1 and E-group2 used different specifications in Parts I and II. After this phase, the participants should have become used to the tool, and were therefore given a time limit of 90 min in the second phase (Part IIb), in which they continued to identify as many MRs as possible, with the same specifications used by them as in Part IIa, but with different candidate pairs.

**Part III (MR identification supported by MR-GEN with on-demand selection)** took place after Tutorial 3, and involved only E-group3. This part has two rounds, the first using  $S_{CAR}$ , and the second using  $S_{MOS}$ . Similar to Part II, each round was divided into two phases with a 1-h break between them, with each participant identifying MRs for the relevant specification, firstly to get 40 MRs with no time constraint (Part IIIa), and then to get as many MRs as possible within a period of 90 min (Part IIIb).

Table 3 summarizes the details of these arrangements.

## 6.3. Specification of input parameters

The input/selection screens (Figs. 1 and 2) allow users to specify the candidate pairs' types and other associated details, the maximum number of MRs to be identified, and the critical categories. Provision of such parameters allows the users to customize the selection process according to their own needs and preferences, but may also render impossible a fair examination of the different components in the research questions (Section 6.1). The parameter specification feature was therefore disabled for the experiment, and the participants were instead provided with a standard setup, where candidate pairs of any type were randomly selected by MR-GEN for evaluation. This setup made it possible to filter out the effects due to a particular sequence of candidate pairs, and thereby measure the viability and effectiveness of METRIC. However, the random order of candidate pairs displayed for user evaluation was not followed in Part III, in which the impact of users' choice of candidate pairs using the on-demand selection feature was examined.

## 6.4. Measurement

MR-GEN timestamps and logs all major events, such as the determination of whether or not candidate pairs are usable, making it possible to examine the time taken and correctness of the evaluations after the experiment.

If an MR identified by the participants satisfied Definition 1, then it was called a *genuine* MR, otherwise it was *non-genuine*. The researchers examined all experiment logs to categorize the participant-identified MRs as genuine or non-genuine. From here until the end of Section 7, the term “metamorphic relations (MRs)” refers to those identified by the participants, which may or may not be genuine. The experiment determined the soundness and completeness of participants' evaluations. **Soundness** ( $R_g$ ) was calculated as the ratio of genuine MRs identified by the participants to either (a) the candidate pairs evaluated as usable by the participants using MR-GEN; or (b) the MRs identified by the participants in an ad hoc manner. In this way,  $R_g$  is a measure of the true-positive rate. **Completeness** ( $R_n$ ) was calculated as the proportion of the *actually* unusable candidate pairs among the candidate pairs evaluated as unusable by the participants, and is therefore a measure of the true-negative rate. Obviously, higher values of  $R_g$  and  $R_n$  mean better results. In MT, soundness is more important than completeness because: (i) every MR used in the construction of source and follow-up test cases must be genuine, otherwise, interpretation of the testing results may be meaningless; and (ii) if the user has already identified some genuine MRs, MT can still be conducted even if some *actually* usable candidate pairs have been *incorrectly* identified as unusable.

Table 4 shows some measures from which important statistics for each participant were collected. Other measures will be introduced in the next section.



**Table 4**  
Measures for participants' performance.

Notation	Meaning
$i_{MR}$	Number of <b>MRs</b> identified by a participant in the ad hoc approach
$g_{MR}$	Number of <b>genuine MRs</b> identified by a participant
$e_u$	Number of candidate pairs evaluated as <b>usable</b> by a participant <sup>1</sup>
$e_n$	Number of candidate pairs evaluated as <b>not usable</b> (that is, unusable) by a participant
$c_n$	Number of candidate pairs correctly evaluated as <b>not usable</b> (that is, unusable) by a participant
$R_g$	$R_g = \frac{g_{MR}}{i_{MR}}$ if using the ad hoc approach, or $R_g = \frac{g_{MR}}{e_u}$ if using MR-GEN; $R_g$ serves as the <b>soundness</b> measure of MR identification <sup>2</sup>
$R_n$	$R_n = \frac{c_n}{e_n}$ if $e_n > 0$ , or $R_n = 1$ if $e_n = 0$ ; $R_n$ serves as the <b>completeness</b> measure of MR identification <sup>2</sup>
$T$	Time taken (in minutes) for identifying a certain number of MRs by a participant <sup>3</sup>
$T_g$	Average time taken (in minutes) for identifying a <b>genuine MR</b> by a participant <sup>4</sup>

<sup>1</sup> In Part I,  $e_u$  was inapplicable because  $i_{MR}$  was used instead. In Parts IIa and IIIa,  $e_u = 40$ . In Parts IIb and IIIb,  $e_u$  was not fixed.

<sup>2</sup> Higher values of  $R_g$  and  $R_n$  mean better results.

<sup>3</sup> In Parts I, IIb, and IIIb,  $T = 90$  min. In Parts IIa and IIIa,  $T$  was not fixed.

<sup>4</sup> In Parts I, IIb, and IIIb,  $T_g = \frac{90min}{g_{MR}}$ . In Parts IIa and IIIa,  $T_g = \frac{T}{g_{MR}}$ . The smaller the value of  $T_g$ , the better the result.

## 7. Experimental results and analyses

### 7.1. MR identification feasibility of METRIC and MR-GEN (RQ1)

#### 7.1.1. Comparison of ad hoc approach and METRIC (RQ1.1)

This question was addressed by involving E-group1 and E-group2 in MR identification in an ad hoc manner (Part I) and comparing this with MR identification supported by MR-GEN (Part II). As a reminder, for each group, the specifications used in Parts I and II were different, which prevented the participants from gaining any advantage through MR knowledge for a particular specification from the previous part, and thus facilitated a fair comparison of the effectiveness of the ad hoc approach and METRIC (supported by MR-GEN). Because inexperienced participants might not be able to identify MRs in an ad hoc manner, the analyses here only involved the experienced participants. Table 5 shows the relevant statistics.

It was observed that the mean values of  $R_g$  in Part IIa were 2.1% (for  $S_{CAR}$ ;  $= \frac{0.95-0.97}{0.95} \times 100\%$ ) and 5.7% (for  $S_{MOS}$ ;  $= \frac{0.88-0.93}{0.88} \times 100\%$ ) higher than those in Part I. The results indicate that less mistakes were made in identifying genuine MRs in Part IIa than in Part I. Thus, it can be concluded that MR-GEN improved the identification of genuine MRs in terms of soundness (even though the soundness for ad hoc identification in Part I was already high). As additional information, it was observed that most of the MRs identified in an ad hoc manner were also identified by using METRIC.

Readers are reminded that the main purpose of developing METRIC is to provide a systematic methodology for supporting MR identification for those testers who are *inexperienced* in MT. Thus, the relative importance of RQ1.1 (which investigates the difference in performance in MR identification between the ad hoc approach and the use of METRIC for *experienced* MT testers) is lower than other research questions.

**Table 5**  
Comparison of MR identification (ad hoc versus tool support).

Ad Hoc: Part I			Supported by MR-GEN: Part IIa		
$S_{CAR}$					
<b>E-group2</b>	$g_{MR}/i_{MR}$	$R_g$	<b>E-group1</b>	$g_{MR}/e_u$	$R_g$
$P_{14}$	25/26	0.96	$P_{11}$	40/40	1.00
$P_{15}$	14/14	1.00	$P_{12}$	38/40	0.95
$P_{16}$	9/10	0.90	$P_{13}$	38/40	0.95
<b>Mean</b>	<b>16.0/16.7</b>	<b>0.95</b>	<b>Mean</b>	<b>38.7/40.0</b>	<b>0.97</b>
$S_{MOS}$					
<b>E-group1</b>	$g_{MR}/i_{MR}$	$R_g$	<b>E-group2</b>	$g_{MR}/e_u$	$R_g$
$P_{11}$	5/5	1.00	$P_{14}$	40/40	1.00
$P_{12}$	10/10	1.00	$P_{15}$	32/40	0.80
$P_{13}$	10/16	0.63	$P_{16}$	40/40	1.00
<b>Mean</b>	<b>8.3/10.3</b>	<b>0.88</b>	<b>Mean</b>	<b>37.3/40.0</b>	<b>0.93</b>

#### 7.1.2. Practicality and feasibility of METRIC for users with primitive MT background (RQ1.2)

RQ1.2 was addressed by examining the performance of inexperienced participants in Part IIa of the experiment, in which they identified 40 MRs without time constraint, supported by MR-GEN (with the on-demand selection feature disabled).

In Table 6, the columns with the headings “I-group1” and “I-group2” show the relevant statistics. Looking first at  $S_{CAR}$ . For I-group1, although the five inexperienced participants P<sub>1</sub> to P<sub>5</sub> received only 1 h of MT training in Tutorial 2 and another 90 min training on MR-GEN in Tutorial 3, on average, they were able to identify MRs with a very high soundness measure ( $R_g$ : mean = 0.98, range = [0.95,1.00]). Regarding the completeness measure, on average, the result was also encouraging ( $R_n$ : mean = 0.93, range = [0.85,1.00]). Similar results were observed for  $S_{MOS}$ . These promising results clearly confirm that using METRIC and its associated tool MR-GEN are practically feasible from the perspective of those users with primitive MT knowledge.

#### 7.1.3. Impact of user's MT experience on MR identification when using METRIC (RQ1.3)

RQ1.3 was addressed by comparing the performance of experienced and inexperienced participants in Part IIa of the experiment. Refer to Table 6 again. For  $S_{CAR}$ , the soundness and completeness measures for I-group1 (see Section 7.1.2 for the means and ranges of these two measures) are quite comparable to the E-group1 statistics ( $R_g$ : mean = 0.97, range = [0.95,1.00];  $R_n$ : mean = 0.96, range = [0.92,1.00]). Similar comparable results were observed for  $S_{MOS}$ . Furthermore, considering  $S_{CAR}$  and  $S_{MOS}$  together, the individual values of  $R_g$  and  $R_n$  were quite similar across all participants. In other words, the impact of previous MT experience on the soundness and completeness aspects of MR identification (using METRIC) appears to be negligible.

Consider next how quickly genuine MRs were identified, in terms of  $T_g$ . Here, mixed results were obtained: For  $S_{CAR}$ , the mean value of  $T_g$  was lower in I-group1 than in E-group1; but for  $S_{MOS}$ , the opposite was observed. The mixed results suggest that more experimental studies are needed to investigate whether there is a general relationship between user's MT experience on MR identification and  $T_g$ . (Analyses were not conducted on  $R_g$ ,  $R_n$ , or  $T_g$  in Part IIb of the experiment. This is because, after identifying 40 MRs using MR-GEN in Part IIa, different participants might have different proficiency with the tool and, hence, might make the analyses related to RQ1.3 invalid. Note that the learning aspect of MR-GEN by the participants will be addressed by RQ1.4.)

At first, it may appear somewhat surprising that, in some cases, the  $R_g$ ,  $R_n$ , and  $T_g$  values for inexperienced participants were slightly better than those for the experienced participants. On reflection,

**Table 6**

Performance of MR identification in Part IIa (inexperienced versus experienced participants).

$S_{CAR}$										
Statistics	I-group1						E-group1			
	$P_1$	$P_2$	$P_3$	$P_4$	$P_5$	Mean	$P_{11}$	$P_{12}$	$P_{13}$	Mean
$R_g$	0.98	0.95	0.98	0.98	1.00	<b>0.98</b>	1.00	0.95	0.95	<b>0.97</b>
(= $g_{MR}/e_u$ )	(39/40)	(38/40)	(39/40)	(39/40)	(40/40)		(40/40)	(38/40)	(38/40)	
$R_n$	0.92	1.00	0.86	0.85	1.00	<b>0.93</b>	0.95	0.92	1.00	<b>0.96</b>
(= $c_n/e_n$ )	(23/25)	(20/20)	(12/14)	(22/26)	(21/21)		(21/22)	(22/24)	(11/11)	
$T_g(\text{min})$	1.92	1.97	1.79	0.90	1.88	<b>1.69</b>	1.88	2.37	1.97	<b>2.07</b>
$S_{MOS}$										
Statistics	I-group2						E-group2			
	$P_6$	$P_7$	$P_8$	$P_9$	$P_{10}$	Mean	$P_{14}$	$P_{15}$	$P_{16}$	Mean
$R_g$	0.95	1.00	0.98	1.00	0.98	<b>0.98</b>	1.00	0.80	1.00	<b>0.93</b>
(= $g_{MR}/e_u$ )	(38/40)	(40/40)	(39/40)	(40/40)	(39/40)		(40/40)	(32/40)	(40/40)	
$R_n^1$	1.00	1.00	1.00	1.00	1.00	<b>1.00</b>	1.00	1.00	1.00	<b>1.00</b>
$T_g(\text{min})$	3.82	2.75	2.31	1.75	2.82	<b>2.69</b>	2.50	3.28	1.63	<b>2.47</b>

<sup>1</sup> For participants  $P_6$  to  $P_{10}$  and  $P_{14}$  to  $P_{16}$ ,  $e_n = 0$ . According to the scheme for determining the value of  $R_n$  in Table 4, if  $e_n = 0$ , then  $R_n = 1$ .

however, some plausible reasons for this can be identified as follows. The determination of “experienced participants” (in E-group1 and E-group2) was partly based on having previous experience in ad hoc MR identification in some other application domains, which was assumed to be somehow advantageous. The apparent impact of this “advantage”, however, was absent in the experiment because of three plausible reasons. Firstly, although MT has a sound concept, its technical content is simple to understand, and offering Tutorial 2 to the inexperienced participants largely reduced or eliminated the MT knowledge gap. Secondly, since METRIC is highly intuitive, simple (but elegant), and supported by its generator tool MR-GEN, both experienced and inexperienced participants were able to easily apply the methodology, which reduced the possible performance gap caused by previous MR identification experience. Thirdly, although the experienced participants had MR identification experience in other application domains, such experience might not have been directly applicable to  $S_{CAR}$  and  $S_{MOS}$ . On the contrary, when compared with the inexperienced participants, the “historical baggage” of previously using an ad hoc approach might have made the experienced participants less able to change their mindset to use METRIC appropriately shortly after learning the methodology.

One may argue that, by using METRIC (and supported by MR-GEN), experienced and inexperienced participants performed comparably well, possibly because of the simplicity of the two specifications rather than the contribution of METRIC. Table 5 provides an observation to address this issue. In Part I (the ad hoc approach) of the experiment, experienced participants in E-group2 (for  $S_{CAR}$ ) and E-group1 (for  $S_{MOS}$ ) obtained the mean values 0.95 and 0.88 for  $R_g$ , respectively. Because the two mean values were less than 1.00, it indicated that even experienced participants identified some invalid MRs using the ad hoc approach. This observation rules out that the two specifications used in the experiment are simple, causing the comparable performance between the experienced and inexperienced participants.

#### 7.1.4. MR identification improvement after some use of METRIC and MR-GEN (RQ1.4)

In I-group1, I-group2, E-group1, and E-group2, participants used MR-GEN to identify MRs for the same assigned specifications in Parts IIa and IIb. This arrangement made it possible to analyze the changes in MR identification after using MR-GEN for a certain period of time, which was done by, for each relevant participant, analyzing and comparing the statistics  $R_g$ ,  $R_n$ , and  $T_g$  in Parts IIa and IIb. From Parts

IIa to IIb, the percentage decrease in the time taken to identify a genuine MR (denoted  $R_{T_g}$ ) was also calculated, using the formula  $R_{T_g} = \frac{T_g \text{ in Part IIa} - T_g \text{ in Part IIb}}{T_g \text{ in Part IIa}}$ . Obviously, a positive value of  $R_{T_g}$  indicates an increase in speed for identifying genuine MRs.

Table 7 shows the relevant statistics for Part II related to  $S_{CAR}$ . (The  $S_{MOS}$  data were similar to those for  $S_{CAR}$  in Table 7. Due to page limitation, the  $S_{MOS}$  data and related discussion are not included in this paper.) As shown in the table, all participants had high  $R_{T_g}$  values, with mean values of 51.5% and 35.1% for inexperienced and experienced groups, respectively. In other words, for all participants, much less time was required to identify a genuine MR after having used MR-GEN for a while. Note that, in addition to increased familiarity with MR-GEN, another factor may also have influenced the improvement in MR identification. After identifying 40 MRs in Part IIa for the assigned specifications, the participants had become more familiar with the specifications, which may have improved their performance in Part IIb. Nevertheless, because time savings for identification of a genuine MR ( $R_{T_g}$ ) were so substantial, by intuition, there were two plausible reasons for this phenomenon: increased familiarity with the specifications and increased familiarity with METRIC/MG-GEN. These observations, together with the earlier finding related to RQ1.3, reinforce the conclusion that the performance using METRIC depends more on familiarity with the methodology and knowledge of the application domains, rather than on previous experience with MT.

In addition to the above observations, Table 7 shows a slight improvement in soundness and completeness for  $S_{CAR}$  from Part IIa (grand mean  $R_g = 0.97$ ; grand mean  $R_n = 0.94$ ) to Part IIb (grand mean  $R_g = 0.99$ ; grand mean  $R_n = 0.96$ ). This phenomenon of little improvement of  $R_g$  and  $R_n$  is understandable because their values in Part IIa were already very close to the possible maximum of 1.0.

A summary of the results is as follows. Firstly, even for experienced participants, MR identification was better when using MR-GEN than without it (RQ1.1). Secondly, even for inexperienced participants (with no previous knowledge about MT prior to the experiment) who have been trained with MT knowledge for only an hour and with the functionality of MR-GEN for 90 min, their performance of MR identification by using METRIC (with the support of MR-GEN) was fairly impressive and encouraging (RQ1.2). Thirdly, when using MR-GEN, both inexperienced and experienced participants had comparable, very good performance, indicating that METRIC is effective for users with different backgrounds and knowledge of MT (RQ1.3). Fourthly, after using MR-GEN for some time, all participants had significantly better MR identification (RQ1.4). Therefore, a clear answer to RQ1 is that METRIC and MR-GEN are practical and feasible for MT identification.

**Table 7**  
Change in performance of MR identification in Part II related to  $S_{CAR}$ .

Participants	Part IIa			Part IIb			From Part IIa to IIb
	$R_g$	$R_n$	$T_g$ (min)	$R_g$	$R_n$	$T_g$ (min)	
<b>I-group1</b>							
$P_1$	0.98	0.92	1.92	1.00	0.91	1.05	45.3
$P_2$	0.95	1.00	1.97	1.00	1.00	0.86	56.3
$P_3$	0.98	0.86	1.79	0.99	0.93	1.11	38.0
$P_4$	0.98	0.85	0.90	0.98	1.00	0.45	50.0
$P_5$	1.00	1.00	1.88	0.99	1.00	0.60	68.1
<b>Mean</b>	<b>0.98</b>	<b>0.93</b>	<b>1.69</b>	<b>0.99</b>	<b>0.97</b>	<b>0.81</b>	<b>51.5</b>
<b>E-group1</b>							
$P_{11}$	1.00	0.95	1.88	1.00	1.00	1.13	39.9
$P_{12}$	0.95	0.92	2.37	1.00	0.93	2.00	15.6
$P_{13}$	0.95	1.00	1.97	0.95	0.93	0.99	49.7
<b>Mean</b>	<b>0.97</b>	<b>0.96</b>	<b>2.07</b>	<b>0.98</b>	<b>0.95</b>	<b>1.37</b>	<b>35.1</b>
<b>Grand mean</b>	<b>0.97</b>	<b>0.94</b>	<b>1.84</b>	<b>0.99</b>	<b>0.96</b>	<b>1.02</b>	<b>45.4</b>

## 7.2. Impact of on-demand selection on the performance of METRIC (RQ2)

Changes in participants' MR identification when provided with MR-GEN's on-demand selection feature were also investigated, with the data for  $S_{CAR}$  shown in Table 8. (The  $S_{MOS}$  data were similar to those for  $S_{CAR}$  in Table 8. Due to page limitation, the  $S_{MOS}$  data and related discussion are not included in this paper. Also, some statistics collected in Part II related to  $S_{CAR}$ , which are shown in Table 7, are included in Table 8 for easy comparison.)

Consider the soundness ( $R_g$ ) and completeness ( $R_n$ ) measures in Table 8: Between Parts II (without on-demand selection) and III (with on-demand selection), the mean values (0.98) for  $R_g$  were identical; but there was a small percentage decrease of 8.3% ( $\frac{0.96-0.88}{0.96} \times 100\%$ ) for  $R_n$ . For  $T_g$ , the percentage improvement was 20.3% ( $\frac{1.72-1.37}{1.72} \times 100\%$ ) from Parts II to III, which was a large change.

Recall that soundness ( $R_g$ ) is more important than completeness ( $R_n$ ) in MT. In addition, although there was a small percentage decrease in  $R_n$  when using the on-demand selection, this decrease was more than offset by the large percentage improvement in  $T_g$  (to some extent,  $T_g$  can also be considered a soundness measure because it indicates how quickly genuine MRs are identified). In summary, with respect to RQ2, Table 8 confirms the merit of the on-demand selection feature for improving MR identification.

## 8. Discussion, limitations, and future work

### 8.1. Discussion

Based on our previous research experience on ad hoc MR identification, three problems associated with this approach are observed. Such problems, however, do not occur in METRIC. Firstly, the ad hoc

MR identification often goes through several peaks and troughs. The testers spend time discovering new knowledge about the input and output aspects of the assigned specifications, but once understood, MRs related to the specification are quickly identified (corresponding to a peak). The MR identification then gradually slows down, eventually to the point where no further MR can be identified based on this new knowledge (corresponding to a trough). When another new piece of knowledge is later discovered, the MR identification again speeds up and reaches another peak. As time passes, and the testers gradually run out of new ideas, the time required to discover new knowledge becomes increasingly longer, until a stage is reached when no further new knowledge can be discovered. Thus, if the testers are asked to identify a fairly large number of MRs, say 100, using the ad hoc approach, then they may not be able to finish the task.

On the other hand, since selection of candidate pairs for evaluation is done automatically by the tool, this problem of decreasing MR identification will not occur when using MR-GEN. With MR-GEN, the users need only focus on one candidate pair at a time, to determine whether or not it is usable. In practice, the number of usable candidate pairs is very often much larger than the maximum number of MRs that can be used for testing under a given amount of testing resources. Thus, the "unprocessed" usable candidate pairs will not be exhausted before the required number of MRs have been identified. Furthermore, in general, usable candidate pairs are uniformly distributed over the candidate pairs selected by MR-GEN (because the candidate pairs are randomly selected according to a uniform distribution), if the on-demand selection feature is disabled. Therefore, the performance of the users in MR identification over time should not decrease due to the gradual use-up of the "unprocessed" usable candidate pairs. This point is confirmed by the observation that the values of  $T_g$  are smaller in Part Ib than in Part Ia in Table 7.

**Table 8**  
With and without on-demand selection in Parts II and III related to  $S_{CAR}$ .

Without on-demand selection				With on-demand selection			
Participants	$R_g$	$R_n$	$T_g$ (min)	Participants	$R_g$	$R_n$	$T_g$ (min)
<b>Part IIa</b>				<b>Part IIIa</b>			
P <sub>11</sub> (in E-group1)	1.00	0.95	1.88	P <sub>17</sub> (in E-group3)	0.95	0.91	1.58
P <sub>12</sub> (in E-group1)	0.95	0.92	2.37	P <sub>18</sub> (in E-group3)	0.95	0.75	1.58
P <sub>13</sub> (in E-group1)	0.95	1.00	1.97	P <sub>19</sub> (in E-group3)	1.00	0.89	1.75
<b>Part IIb</b>				<b>Part IIIb</b>			
P <sub>11</sub> (in E-group1)	1.00	1.00	1.13	P <sub>17</sub> (in E-group3)	0.97	0.93	0.92
P <sub>12</sub> (in E-group1)	1.00	0.93	2.00	P <sub>18</sub> (in E-group3)	1.00	0.90	1.22
P <sub>13</sub> (in E-group1)	0.95	0.93	0.99	P <sub>19</sub> (in E-group3)	1.00	0.87	1.14
<b>Mean</b>	<b>0.98</b>	<b>0.96</b>	<b>1.72</b>	<b>Mean</b>	<b>0.98</b>	<b>0.88</b>	<b>1.37</b>



Secondly, the ad hoc approach does not have an *explicit* distinction between the two steps: (a) selecting, from a large pool of test cases, those which are potentially useful for MR identification; and (b) determining the corresponding MRs (if any) from the selected test cases. Handling steps (a) and (b) together makes the ad hoc identification task more complex.

Thirdly, in the ad hoc approach, many testers are unlikely to identify MRs involving differences in *multiple* variables. This is because such an attempt requires the testers to determine what simultaneous changes of multiple variables would give rise to MRs, which is a fairly complex cognitive task<sup>1</sup> (this problem is further complicated by handling steps (a) and (b) above at the same time). On the other hand, the tendency of not identifying MRs involving differences in multiple variables does not exist in METRIC. In the methodology supported by MR-GEN, after the users have specified the numbers of differentiating categories and differentiating choices (and other parameters), the tool will *automatically* select different types of candidate pairs for consideration during MR identification, even including those involving *multiple* differentiating categories and choices (if the users so specify). Thereafter, the relevant candidate pairs will be displayed, one at a time, for user evaluation. This arrangement makes an explicit distinction between steps (a) and (b) above, resulting in their separate processing. This advantage and the visual side-by-side presentation of candidate pairs (supported by highlighting differentiating choices with different colors) ease the task of identifying MRs from candidate pairs, even if these candidate pairs involve several differentiating categories and choices.

To apply METRIC, users may start by identifying categories and choices, and then construct complete test frames. In the experiment, however, the participants were provided with complete test frames at the start of MR identification. This was done because category-choice framework techniques are now very well developed, with a lot of work already done for identifying proper and appropriate categories and choices (Chen et al., 2005; 2012; Hierons et al., 2003; Poon et al., 2011; Singh et al., 1997), and for constructing complete test frames (Chen et al., 2003a; Grochtmann and Grimm, 1993; Hierons et al., 2003; Kruse and Luniak, 2010; Singh et al., 1997). Furthermore, several tools have been developed to support the construction of complete test frames (Cain et al., 2004; Chen et al., 2003a; Grochtmann and Grimm, 1993; Lehmann and Wegener, 2000). Therefore, it should not be a difficult task to obtain categories, choices, and complete test frames from specifications. In other words, complete test frames are easily available for the application of METRIC. Since there exist several techniques for constructing complete test frames, providing the participants with complete test frames in the experiment enabled the effectiveness of METRIC to be measured alone, by filtering out the impact on the effectiveness due to different construction techniques.

As with most other methodologies, it is unrealistic to expect that METRIC could be effectively used to identify MRs for all application domains, including those such as “ $\sin(x) = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \dots$ ”. Nevertheless, METRIC does have very broad applicability, partly because it is based on the category-choice framework, which has been confirmed by many studies to be applicable to a wide range of application domains, including non-scientific ones. Examples include airline catering systems, integrated ship management systems, adaptive cruise control systems for vehicles, rail ticketing systems, automatic mail sorting systems, online telephone inquiry systems, and inventory management systems (Chen et al., 2012; 2003a; Conrad and Krupp, 2006; Grochtmann and Grimm, 1993; Hierons et al., 2003; Singh et al., 1997). Based on previous work on MR, and feedback

obtained from software researchers and practitioners, there appears to be a general perception that, although it may not be difficult to identify MRs for scientific applications, it may be more challenging to do so for non-scientific applications. The work on METRIC provides a solution to this problem.

## 8.2. Study limitations

There are four limitations to the present studies. Firstly, the studies only involved 19 participants. Although it would be desirable to have more participants involved, it is not easy to find a large group of software testers (especially those with MT knowledge and research experience), who are willing to participate. Secondly, it is often impossible to “perfectly” classify participants into experienced and inexperienced groups. In our experiment, the experienced group included those participants who have acquired prior MT knowledge through previous research work related to MT, whereas the inexperienced group included those who had no knowledge of MT before the experiment. We consider that the classification scheme was fairly rigorous and, hence, the experimental results and comparison between the two groups should be largely valid. Thirdly, only two commercial specifications were used in the studies. As is known, obtaining real-life specifications from industry is difficult because most companies are hesitant to release them for external use. Nonetheless, the studies do provide a convincing demonstration of the viability and effectiveness of METRIC. Lastly, it would be better to compare METRIC with other similar methodologies, however, as pointed out in Section 1, we are not aware of any other systematic and generally applicable methodologies for identifying MRs directly from specifications. Thus, such a comparison is not applicable. This issue, in fact, clearly demonstrates the novelty and contribution of METRIC. It would be worth mentioning that the sequence of presenting the candidate pairs could be a threat to validity. In our experiment, these candidate pairs were presented at random, which has alleviated the problem to some extent.

## 8.3. Future work

As the first study, the current version of METRIC focuses on generating those MRs which are associated with exactly two distinct complete test frames. Obviously, generating MRs involving more than two complete test frames is also possible with METRIC, by extending the current methodology to make use of more complete test frames for MR identification. This extension will be explored further in our future work. In addition, our current experimental study has mainly targeted at the effectiveness of MR identification with the use of METRIC. We plan to conduct further experiments (possibly involving mutation analysis) to investigate the fault-detection capability of the MRs generated by METRIC. Furthermore, it would be worth investigating strategies for determining the appropriate number of differentiating categories and that of differentiating choices (see Fig. 1), and the critical categories (see Fig. 2) according to some criteria such as the cost-effectiveness of fault detection.

## 9. Related work

Kaiser and her colleagues (Mishra and Kaiser, 2012; Mishra et al., 2013) conducted two studies involving teaching students about MT at Columbia University. Their objective was to determine how well the students identified MRs in an ad hoc manner. Both studies found that the students occasionally failed to identify genuine MRs, and some of their identified MRs were non-genuine. However, their studies have not provided a systematic technique for MR identification.

Although systematic generation of MRs is essential for the automation of MT, to date, very limited work has been done in this

<sup>1</sup> Miller (1956) observed that there exists a capacity limit of human cognition. More specifically, he found that the number of objects an average human can hold in working memory is  $7 \pm 2$ . Increasing the number of objects will reduce the performance of a human. His finding is often referred to as *Miller's Law* and is one of the most highly cited papers in psychology (Baddeley, 1994).

area, with the only published investigations by Kanewala and Bieman (2013), Liu et al. (2012), and Zhang et al. (2014). The technique developed by Kanewala and Bieman (2013) uses a machine-learning approach to detect MRs from a system's control flow graph, but machine learning requires the use of a training set for creating the predictive model. In other words, Kanewala and Bieman's technique requires that some MRs have already been identified as a training set, through which new MRs can then be generated. However, obtaining a good and comprehensive training set in machine learning may not be easy (Cortes et al., 1995). (In the future, it would be interesting to investigate the effectiveness of MR identification by first generating MRs using METRIC and then augmenting the set of generated relations with Kanewala and Bieman's technique.) The technique from Liu et al. (2012) attempts to construct new MRs from existing ones through composition and, hence, can only generate certain specified types of MRs. In summary, both of these techniques (Kanewala and Bieman, 2013; Liu et al., 2012) require the presence of some existing MRs in order to identify new ones. Finally, the technique from Zhang et al. (2014) only targets at identifying polynomial MRs, and is thereby applicable to specific domains. METRIC, however, differs from the above three techniques because it does not have the requirement of existing MRs and can be applied for many different types of MRs. Consequently, METRIC is applicable to broader application domains.

## 10. Summary and conclusion

The problem of systematic identification for metamorphic relations (MRs) directly from specifications is generally considered to be a challenging, but important, prerequisite for effective metamorphic testing (MT). To date, only very limited work has been done in the identification of MRs, a situation which is far from desired. To address this problem, our paper has proposed the METRIC identification methodology, which helps users identify MRs from specifications in a systematic manner.

METRIC is built upon the category-choice framework. Through the systematic identification of usable candidate pairs, MRs can be easily and intuitively generated, even for those "complex" MRs involving multiple differentiating categories and choices. To further improve the effectiveness of METRIC, and ease the identification task, an associated generator tool MR-GEN has been developed. The tool includes several useful features such as: (a) specification of the type of candidate pairs for user evaluation; (b) specification of critical categories; (c) visual side-by-side presentation of candidate pairs and highlighting of differentiating choices to facilitate user evaluation; and (d) on-demand selection of subsequent candidate pairs.

To determine the viability and effectiveness of METRIC (supported by MR-GEN), an experiment using two commercial specifications was conducted, involving two groups of participants (one group experienced in MT, and the other not). The main purpose of the experiment was to answer two research questions: (i) Can METRIC and MR-GEN be practically and feasibly used for MR identification? (ii) Does use of MR-GEN's on-demand selection feature result in a significant improvement in MR identification? With respect to both research questions, the experimental results were largely positive and encouraging, even for those users without substantial MT knowledge. Thus, METRIC represents a significant contribution, and will, without doubt, improve the applicability, effectiveness, and automation of MT.

## Acknowledgments

We would like to express our gratitude to Prof. Huimin Lin of the Institute of Software at the Chinese Academy of Sciences for providing the laboratory facilities for the experiment. We gratefully acknowledge the contribution of all the participants from Beihang University and the University of Science and Technology Beijing who

were involved in the experiment. Furthermore, we are indebted to Dr. Dave Towey of the University of Nottingham Ningbo China for his invaluable comments and discussions.

## References

- Baddeley, A., 1994. The magical number seven: still magic after all these years? *Psychol. Rev.* 101 (2), 353–356.
- Barr, E.T., Harman, M., McMinn, P., Shahbaz, M., Yoo, S., 2014. The oracle problem in software testing: a survey. *IEEE Trans. Softw. Eng.* doi:10.1109/TSE.2014.2372785.
- Cain, A., Chen, T.Y., Grant, D., Poon, P.-L., Tang, S.-F., Tse, T.H., 2004. An automatic test data generation system based on the integrated classification-tree methodology. In: Ramamoorthy, C.V., Lee, R.Y., Lee, K.W. (Eds.), *Software Engineering Research & Applications*. In: *Lecture Notes in Computer Science*, 3026. Springer, Berlin, Germany, pp. 225–238.
- Chan, W.K., Chen, T.Y., Lu, H., Tse, T.H., Yau, S.S., 2006. Integration testing of context-sensitive middleware-based applications: a metamorphic approach. *Int. J. Softw. Eng. Knowl. Eng.* 16 (5), 677–703.
- Chan, W.K., Cheung, S.C., Leung, K.R., 2007. A metamorphic testing approach for online testing of service-oriented software applications. *Int. J. Web Serv. Res.* 4 (2), 60–80.
- Chen, T.Y., Cheung, S.C., Yiu, S.M., 1998. Metamorphic Testing: A New Approach for Generating Next Test Cases. Technical report HKUST-CS98-01. Department of Computer Science, Hong Kong University of Science & Technology, Hong Kong.
- Chen, T.Y., Ho, J.W.K., Liu, H., Xie, X., 2009. An innovative approach for testing bioinformatics programs using metamorphic testing. *BMC Bioinformatics* 10. doi:10.1186/1471-2105-10-24.
- Chen, T.Y., Kuo, F.-C., Towey, D., Zhou, Z.Q., 2015. A revisit of three studies related to random testing. *Sci. China: Inf. Sci.* 58 (5). Article 052104(9).
- Chen, T.Y., Poon, P.-L., Tang, S.-F., Tse, T.H., 2005. Identification of categories & choices in activity diagrams. In: *Proceedings of the 5th International Conference on Quality Software (QSIC 2005)*. IEEE Computer Society Press, Los Alamitos, CA, pp. 55–63.
- Chen, T.Y., Poon, P.-L., Tang, S.-F., Tse, T.H., 2012. DESSERT: a Divide-and-conquer methodology for identifying categories, choices, & choice Relations for Test case generation. *IEEE Trans. Softw. Eng.* 38 (4), 794–809.
- Chen, T.Y., Poon, P.-L., Tse, T.H., 2003a. A choice relation framework for supporting category-partition test case generation. *IEEE Trans. Softw. Eng.* 29 (7), 577–593.
- Chen, T.Y., Tse, T.H., Zhou, Z.Q., 2002. Semi-proving: an integrated method based on global symbolic evaluation & metamorphic testing. In: *Proceedings of the ACM SIGSOFT International Symposium of Software Testing & Analysis (ISSTA 2002)*. ACM Press, New York, NY, pp. 191–195.
- Chen, T.Y., Tse, T.H., Zhou, Z.Q., 2003b. Fault-based testing without the need of oracles. *Inf. Softw. Technol.* 45 (1), 1–9.
- Chen, T.Y., Tse, T.H., Zhou, Z.Q., 2011. Semi-proving: an integrated method for program proving, testing, & debugging. *IEEE Trans. Softw. Eng.* 37 (1), 109–125.
- Conrad, M., Krupp, A., 2006. An extension of the classification-tree method for embedded systems for the description of events. *Electronic Notes in Theor. Comput. Sci.* 164 (4), 3–11.
- Cortes, C., Jackel, L.D., Chiang, W.-P., 1995. Limits on learning machine accuracy imposed by data quality. In: *Proceedings of the 1st International Conference on Knowledge Discovery & Data Mining (KDD 1995)*, 95. AAAI Press, Palo Alto, CA, pp. 57–62.
- Gotlieb, A., Botella, B., 2003. Automated metamorphic testing. In: *Proceedings of the 27th Annual International Computer Software & Applications Conference (COMPSAC 2003)*. IEEE Computer Society Press, Los Alamitos, CA, pp. 34–40.
- Groce, A., Kulesza, T., Zhang, C., Shamasunder, S., Burnett, M., Wong, W.-K., Stumpf, S., Das, S., Shinsell, A., Bice, F., McIntosh, K., 2014. You are the only possible oracle: effective test selection for end users of interactive machine learning systems. *IEEE Trans. Softw. Eng.* 40 (3), 307–323.
- Grochtmann, M., Grimm, K., 1993. Classification trees for partition testing. *Softw. Test. Verif. Reliab.* 3 (2), 63–82.
- Hamlet, R., 2002. Random testing. In: Marciniak, J.J. (Ed.), *Encyclopedia of Software Engineering*. Wiley, New York, NY. doi:10.1002/0471028959.sof268.
- Hierons, R.M., 2012. Oracles for distributed testing. *IEEE Trans. Softw. Eng.* 38 (3), 629–641.
- Hierons, R.M., Harman, M., Singh, H., 2003. Automatically generating information from a Z specification to support the classification tree method. In: *Proceedings of the 3rd International Conference on Formal Specification & Development in Z & B (ZB 2003)*. Springer, Berlin, Germany, pp. 388–407.
- Hui, Z., Huang, S., Ren, Z., Yao, Y., 2013. Metamorphic testing integer overflow faults of mission critical program: a case study. *Math. Probl. Eng.* 2013. Article ID 381389. <http://dx.doi.org/10.1155/2013/381389> (accessed January 2015).
- Jones, J.A., Harrold, M.J., 2005. Empirical evaluation of the Tarantula automatic fault-localization technique. In: *Proceedings of the 20th IEEE/ACM International Conference on Automated Software Engineering (ASE '05)*. ACM Press, New York, NY, pp. 273–282.
- Kanewala, U., Bieman, J.M., 2013. Using machine learning techniques to detect metamorphic relations for programs without test oracles. In: *Proceedings of the 24th IEEE International Symposium on Software Reliability Engineering (ISSRE 2013)*. IEEE Computer Society Press, Los Alamitos, CA, pp. 1–10.
- Kruse, P.M., Luniak, M., 2010. Automated test case generation using classification trees. *Softw. Qual. Prof. Mag.* 13 (1), 4–12.
- Le, V., Afshari, M., Su, Z., 2014. Compiler validation via equivalence modulo inputs. In: *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design & Implementation (PLDI'14)*. ACM Press, New York, NY, pp. 216–226.

- Lehmann, E., Wegener, J., 2000. Test case design by means of the CTE XL. In: Proceedings of the 8th European International Conference on Software Testing, Analysis & Review (EuroSTAR 2000), Copenhagen, Denmark. <http://systematic-testing.com/documents/eurostar2000.pdf> (accessed January 2015).
- Liu, H., Kuo, F.-C., Towey, D., Chen, T.Y., 2014a. How effectively does metamorphic testing alleviate the oracle problem? *IEEE Trans. Softw. Eng.* 40 (1), 4–22.
- Liu, H., Liu, X., Chen, T.Y., 2012. A new method for constructing metamorphic relations. In: Proceedings of the 12th International Conference on Quality Software (QSIQ 2012). IEEE Computer Society Press, Los Alamitos, CA, pp. 59–68.
- Liu, H., Yusuf, I.I., Schmidt, H.W., Chen, T.Y., 2014b. Metamorphic fault tolerance: an automated & systematic methodology for fault tolerance in the absence of test oracle. In: Companion Proceedings of the 36th International Conference on Software Engineering (ICSE Companion 2014). ACM Press, New York, NY, pp. 420–423.
- Lyu, M.R., Horgan, J., London, S., 1994. A coverage analysis tool for the effectiveness of software testing. *IEEE Trans. Reliab.* 43 (4), 527–535.
- Mayer, J., Guderlei, R., 2006. An empirical study on the selection of good metamorphic relations. In: Proceedings of the 30th Annual International Computer Software & Applications Conference (COMPSAC 2006), 1. IEEE Computer Society Press, Los Alamitos, CA, pp. 475–484.
- Miller, G.A., 1956. The magical number seven, plus or minus two: some limits on our capacity for processing information. *Psychol. Rev.* 63 (2), 81–97.
- Mishra, K.S., Kaiser, G., 2012. Effectiveness of Teaching Metamorphic Testing. Technical report CUCS-020-12. Department of Computer Science, Columbia University, New York, NY.
- Mishra, K.S., Kaiser, G.E., Sheth, S.K., 2013. Effectiveness of Teaching Metamorphic Testing, part II. Technical report CUCS-022-13. Department of Computer Science, Columbia University, New York, NY.
- Morell, L.J., 1990. A theory of fault-based testing. *IEEE Trans. Softw. Eng.* 16 (8), 844–857.
- Murphy, C., Raunak, M., King, A., Chen, S., Imbriano, C., Kaiser, G., Lee, I., Sokolsky, O., Clarke, L., Osterweil, L., 2011. On effective testing of health care simulation software. In: Proceedings of the 3rd Workshop on Software Engineering in Health Care (SEHC 2011). ACM Press, New York, NY, pp. 40–47.
- Myers, G.J., Sandler, C., Badgett, T., 2004. *The Art of Software Testing*, 3rd ed. Wiley, Hoboken, NJ.
- Nie, C., Leung, H., 2011. A survey of combinatorial testing. *ACM Comput. Surv.* 43 (2), 11:1–11:29.
- Peters, D.K., Parnas, D.L., 1998. Using test oracles generated from program documentation. *IEEE Trans. Softw. Eng.* 24 (3), 161–173.
- Poon, P.-L., Kuo, F.-C., Liu, H., Chen, T.Y., 2014. How can non-technical end users effectively test their spreadsheets? *Inf. Technol. People* 27 (4), 440–462.
- Poon, P.-L., Tse, T.H., Tang, S.-F., Kuo, F.-C., 2011. Contributions of tester experience & a checklist guideline to the identification of categories & choices for software testing. *Softw. Qual. J.* 19 (1), 141–163.
- Segura, S., Hierons, R.M., Benavides, D., Ruiz-Cortés, A., 2011. Automated metamorphic testing on the analyses of feature models. *Inf. Softw. Technol.* 53 (3), 245–258.
- Singh, H., Conrad, M., Sadeghipour, S., 1997. Test case design based on Z & the classification-tree method. In: Proceedings of the 1st IEEE International Conference on Formal Engineering Methods (ICFEM 1997). IEEE Computer Society Press, Los Alamitos, CA, pp. 81–90.
- Tao, Q., Wu, W., Zhao, C., Shen, W., 2010. An automatic testing approach for compiler based on metamorphic testing technique. In: Proceedings of the 17th Asia Pacific Software Engineering Conference (APSEC 2010). IEEE Computer Society Press, Los Alamitos, CA, pp. 270–279.
- Weyuker, E.J., 1982. On testing non-testable programs. *The Comput. J.* 25 (4), 465–470.
- White, L.J., Cohen, E.L., 1980. A domain strategy for computer program testing. *IEEE Trans. Softw. Eng.* 6 (3), 247–257.
- Xie, X., Ho, J.W.K., Murphy, C., Kaiser, G., Xu, B., Chen, T.Y., 2011. Testing & validating machine learning classifiers by metamorphic testing. *J. Syst. Softw.* 84 (4), 544–558.
- Xie, X., Wong, W.E., Chen, T.Y., Xu, B., 2013. Metamorphic slice: an application in spectrum-based fault localization. *Inf. Softw. Technol.* 55 (5), 866–879.
- Zhang, J., Chen, J., Hao, D., Xiong, Y., Xie, B., Zhang, L., Mei, H., 2014. Search-based inference of polynomial metamorphic relations. In: Proceedings of the 29th ACM/IEEE International Conference on Automated Software Engineering (ASE'14). ACM Press, New York, NY, pp. 701–712.
- Zhou, Z.Q., Zhang, S., Hagenbuchner, M., Tse, T.H., Kuo, F.-C., Chen, T.Y., 2012. Automated functional testing of online search services. *Softw. Test. Verif. Reliab.* 22 (4), 221–243.

**Tsong Yueh Chen** received the BSc and MPhil degrees from The University of Hong Kong, the MSc and DIC degrees from the Imperial College of Science and Technology, and the PhD degree from The University of Melbourne. He is a Professor of Software Engineering at Swinburne University of Technology, Australia. He is the principal co-inventor of metamorphic testing and adaptive random testing. His current research interests include software testing, debugging, and program repair.

**Pak-Lok Poon** received his PhD degree in Software Engineering from The University of Melbourne. Prior to joining RMIT University in Australia, he has taught at The Hong Kong Polytechnic University. His research interests include software testing, requirements engineering and inspection, information systems audit and control, electronic commerce, and computers in education. He served on the Editorial Board of the *Information Systems Control Journal*, and is currently an Editor of the *International Journal of Information Processing and Management*.

**Xiaoyuan Xie** received her BEng and MEng degrees from Southeast University (China), and her PhD degree from Swinburne University of Technology (Australia). Her research interests include software testing, debugging and fault localization, search-based software engineering, and evolutionary computing. She is currently a Professor at the State Key Laboratory of Software Engineering, Wuhan University, China.