



# A systematic literature review of techniques and metrics to reduce the cost of mutation testing

Alessandro Viola Pizzoleto<sup>a</sup>, Fabiano Cutigi Ferrari<sup>a,b,1,\*</sup>, Jeff Offutt<sup>b</sup>, Leo Fernandes<sup>c</sup>, Márcio Ribeiro<sup>d</sup>

<sup>a</sup> Computing Department, Federal University of São Carlos—São Carlos, SP, Brazil

<sup>b</sup> Department of Computer Science, George Mason University—Fairfax, VA, USA

<sup>c</sup> Informatics Coordination, Federal Institute of Alagoas—Maceió, AL, Brazil

<sup>d</sup> Computing Institute, Federal University of Alagoas—Maceió, AL, Brazil



## ARTICLE INFO

### Article history:

Received 5 January 2019

Revised 24 July 2019

Accepted 25 July 2019

Available online 26 July 2019

### Keywords:

Mutation analysis

Mutation testing

Cost reduction

Systematic review

## ABSTRACT

Historically, researchers have proposed and applied many techniques to reduce the cost of mutation testing. It has become difficult to find all techniques and to understand the cost-benefit tradeoffs among them, which is critical to transitioning this technology to practice. This paper extends a prior workshop paper to summarize and analyze the current knowledge about reducing the cost of mutation testing through a systematic literature review. We selected 175 peer-reviewed studies, from which 153 present either original or updated contributions. Our analysis resulted in six main goals for cost reduction and 21 techniques. In the last decade, a growing number of studies explored techniques such as selective mutation, evolutionary algorithms, control-flow analysis, and higher-order mutation. Furthermore, we characterized 18 metrics, with particular interest in the number of mutants to be executed, test cases required, equivalent mutants generated and detected, and mutant execution speedup. We found that cost reduction for mutation is increasingly becoming interdisciplinary, often combining multiple techniques. Additionally, measurements vary even for studies that use the same techniques. Researchers can use our results to find more detailed information about particular techniques, and to design comparable and reproducible experiments.

© 2019 Elsevier Inc. All rights reserved.

## 1. Introduction

Mutation testing (Hamlet, 1977; DeMillo et al., 1978; Acree et al., 1979) (*program mutation*, *mutation analysis*, or simply *mutation*) is a testing criterion that has been extensively studied yet lightly used in the last four decades (Offutt and Untch, 2000; Jia and Harman, 2011; Papadakis et al., 2019; Ferrari et al., 2018a). Mutation creates modified versions of a program, called *mutants*, by applying *mutation operators* that make small changes to the software to either simulate faults or guide the tester to edge cases (DeMillo et al., 1978). Testers are expected to find or design tests that cause these mutants to fail, that is, behave differently from the original un-mutated program. If a test case causes a mutant to

fail, then that mutant is said to be *killed*; otherwise, the mutant remains alive. Mutation can be used to help testers design high quality tests or to evaluate and improve existing test sets.

Mutation testing has been empirically shown to be stronger than test criteria such as control-flow based testing and data-flow based testing (Frankl et al., 1997; Mathur and Wong, 1994; Mathur, 2007; Li et al., 2009). Moreover, it has been used to measure the quality of tests and other test criteria in numerous studies (Andrews et al., 2005; Do and Rothermel, 2006; Just et al., 2014b). Despite its effectiveness, several factors make it expensive and therefore difficult to use in practice: large sets of mutants that must be executed, sometimes many times; generation of test cases to kill the mutants; the number of tests needed; and equivalent mutants (Offutt and Untch, 2000; Jia and Harman, 2011; Reales and Polo, 2014; Kintis et al., 2017). For even small methods with a few dozen lines of code, hundreds of mutants may be generated. Some mutants are trivially killed (that is, killed by all tests) and some are easily killed (that is, killed by many tests), yet some are more difficult and require detailed analysis by the tester. It is also very difficult (in fact undecidable in general) (Budd and Angluin, 1982))

\* Corresponding author at: Computing Department, Federal University of São Carlos—São Carlos, SP, Brazil.

E-mail addresses: [alessandro.pizzoleto@ufscar.br](mailto:alessandro.pizzoleto@ufscar.br) (A.V. Pizzoleto), [fcferrari@ufscar.br](mailto:fcferrari@ufscar.br) (F.C. Ferrari), [offutt@gmu.edu](mailto:offutt@gmu.edu) (J. Offutt), [leonardo.oliveira@ifal.edu.br](mailto:leonardo.oliveira@ifal.edu.br) (L. Fernandes), [marcio@ic.ufal.br](mailto:marcio@ic.ufal.br) (M. Ribeiro).

<sup>1</sup> Was a visiting researcher at George Mason University, Fairfax, VA, USA, in 2017 and 2018.

to determine whether a mutant is equivalent or is simply difficult to kill. And not surprisingly, the most difficult-to-kill mutants often lead to valuable test cases.

These costs have been a major obstacle to practical adoption of mutation testing. For instance, some say the difficulty of identifying equivalent mutants is the primary reason mutation is not used more widely (Jia and Harman, 2011; Madeyski et al., 2014). Researchers have therefore invented many ways to reduce the cost, focusing on several goals. Specific goals include reducing the number of mutants (Acree et al., 1979; Mathur and Wong, 1993; Offutt et al., 1993), not creating certain mutants (Marshall et al., 1990), speeding up mutant execution (Krauser et al., 1988; Untch, 1992; Weiss and Fleysgakker, 1993), automating test set generation (DeMillo and Offutt, 1991; 1993), minimizing or prioritizing test sets (Sahinoglu and Spafford, 1990), and automatically identifying equivalent mutants (Offutt and Craft, 1994).

Historically, many techniques to reduce mutation testing costs have been proposed, developed, and studied. The sheer volume of papers and the fact that they are published in dozens of different conferences and journals makes it hard to find them all, and even more difficult to understand the cost-benefit tradeoffs among them (Kurtz et al., 2016; Just et al., 2017; Gopinath et al., 2017). To help current and future researchers understand what has been done and discover new directions, as well as to organize the knowledge to support technology transition, we have performed a systematic literature review (SLR) to characterize the history and the state-of-the-art of mutation testing cost reduction. We previously presented preliminary results from our SLR in a workshop paper (Ferrari et al., 2018a), with general classifications of primary studies. This paper greatly extends and updates the prior short paper, and makes four key contributions:

- It presents comprehensive classifications for cost reduction goals and cost reduction techniques, including lists of references to selected studies.
- It provides an overview of all selected studies grouped by cost reduction goal, which is complemented by extensive online material (Ferrari et al., 2018b).
- It characterizes metrics used to measure cost reduction, including lists of references to selected studies.
- It summarizes, analyzes, and discusses the cost savings and test effectiveness as reported in the selected studies.

We analyzed 175 peer-reviewed, published, primary studies<sup>2</sup>. We down-selected these to 153 studies by removing studies that were either extended or updated by later studies (*subsumed*). We then characterized 6 main goals related to cost reduction, and then identified 21 cost reduction techniques. We also identified and characterized 18 metrics that have been used to measure cost reduction. In addition to the fact that cost reduction research has been increasing, we found that cost reduction techniques are becoming interdisciplinary and are more frequently combined. Moreover, we found out that cost reduction techniques have been measured in many ways, and the results vary even for studies that apply the same cost reduction technique.

The remainder of this paper is organized as follows. Section 2 provides an overview of mutation testing and the need to reduce its costs. Section 3 provides details of our SLR protocol, including goals, research questions (RQ1, RQ2 and RQ3), and the criteria and procedures we used to select and analyze the studies. Section 4 summarizes results from our search. Section 5 answers research question RQ1, by analyzing the results with a focus on the goals of the studies and the cost reduction

techniques. Section 6 provides answers to research questions RQ2 and RQ3. It first lists metrics used to measure cost reduction, then analyzes the results of the selected studies regarding cost reduction measurements and achieved mutation scores. We summarize threats to the validity of this SLR in Section 7, and describe related research in Section 8. Finally, conclusions, implications, and future work are presented in Section 9. Appendix A brings tables with selected studies and subsumed studies.

## 2. Background

Recent literature reviews (Jia and Harman, 2011; Silva et al., 2017; Ferrari et al., 2018a; Papadakis et al., 2019) have shown that research in mutation testing has been steadily increasing through the last four decades. Mutation is a fault-based testing criterion that creates modified versions of the program, called *mutants* (DeMillo et al., 1978). Testers use mutation to design tests and to evaluate externally created tests (Andrews et al., 2005; Do and Rothermel, 2006; Just et al., 2014b). Mutation has also been used to assess other testing criteria. Papadakis et al. recently summarized 217 papers published from 2008 and 2017 that used mutation testing to perform test assessment (Papadakis et al., 2019).

Mutation is based on two fundamental hypotheses: the *Competent Programmer Hypothesis* and the *Coupling Effect Hypothesis* (DeMillo et al., 1978). The *Competent Programmer Hypothesis* says that although programmers do not write correct programs, they write programs that are close to being correct. The coupling effect hypothesizes that test data that find faults that are very close to the original program (that is, mutants) can also find more complex faults. That is, complex faults are *coupled* to simple faults (DeMillo et al., 1978). Both hypotheses, taken together, mean that tests that kill all or most simple mutants will probably kill more complex mutants.

*Mutations operators* modify the program under test to create mutants. For example, an arithmetic operator would change the expression  $(a + b)$  to  $(a^*b)$ ,  $(a - b)$ , and  $(a/b)$ . Mutation operators use fault taxonomies that are usually based on studies of faults in real programs. Mutation operators are applied to a program  $P$  to create a set of mutants  $M$ . Each test  $t$  in a test set  $T$  is run against each mutant  $m$ ,  $m \in M$ . If  $m(t) \neq P(t)$  for some  $t$ , then we say that  $t$  has killed  $m$ . If not, the tester should find or design a test that kills  $m$ . If  $m$  and  $P$  are equivalent, then  $P(t) = m(t)$  for all possible test cases.

As originally conceived (DeMillo et al., 1978), mutation testing was applied in four steps: (1) execution of the original program; (2) generation of mutants; (3) execution of the mutants; and (4) analysis of the mutants. After each cycle, the *mutation score* is calculated with Eq. 1. In the mutation score,  $K$  is the number of mutants of  $P$  killed by  $T$ ;  $M$  is the total number of mutants created from  $P$ ; and  $E$  is the number of mutants of  $P$  that are equivalent to  $P$ . The mutation score is a value in the interval [0,1] that reflects the quality of the test set with respect to the mutants. The closer to 1 the mutation score is, the better the test set is.

$$MS = \frac{K}{M - E} \quad (1)$$

Steps 3 (execution) and 4 (analysis) are very costly, primarily due to two factors: the large number of executions and the need to hand-analyze live mutants. This paper characterizes research into reducing the cost of mutation testing. Step 3 has been highly automated, whereas step 4 is usually done by hand. The cost of step 4 is considered by some to be the main impediment for the practical adoption of mutation testing. Part of this analysis involves detecting equivalent mutants, which is a generally undecidable problem (Budd and Angluin, 1982). According to Madeyski et al. (2014), several partial automated solutions to the *equivalent mutant*

<sup>2</sup> Primary study is a term used in evidence-based research (Kitchenham et al., 2004) to describe research results from well-founded experimental procedures or from incipient research approaches. SLRs, on the other hand, are secondary studies.

problem have been developed. However, these solutions are limited and thus equivalent mutants are usually found by hand.

### 3. Study setup

A systematic literature review (SLR) is a rigorous approach to identify, evaluate, and interpret all available evidence about a particular topic of interest (Kitchenham, 2004). A basic requirement for SLRs is to document the process to ensure auditability and reproducibility (Kitchenham, 2004).

The following subsections summarize the protocol we used. The complete protocol is provided on a companion website (Ferrari et al., 2018b).

#### 3.1. Goals and research questions

The general goal of this SLR is to summarize and analyze the history and state-of-the-art of efforts to reduce the cost of mutation testing. We define three research questions:

- (RQ1) Which techniques support cost reduction of mutation testing?
- (RQ2) Which metrics have been used to measure the cost reduction of mutation testing?
- (RQ3) What are the savings (benefit) and loss of effectiveness (as proxied by mutation score) for the techniques?

These research questions are addressed in Sections 5 (RQ1) and 6 (RQ2 and RQ3).

#### 3.2. Control group and study selection criteria

For this SLR, the baseline dataset to assess the study selection process (the *control group* (Biolchini et al., 2005)) includes cost reduction-related studies described in Jia and Harman's survey (Jia and Harman, 2011). We applied the following inclusion (I) and exclusion (E) selection criteria to the studies:

- (I1): The study was included if it proposes an approach to reduce the cost of mutation testing.
- (I2): The study was included if it applies an approach to reduce the cost of mutation testing.
- (I3): The study was included if it is a primary study, peer-reviewed, and published either in a conference or a scientific journal.
- (E1): The study was excluded if it mentions or uses a technique that can be used to reduce the cost of mutation testing, but cost reduction is not a main goal of that study.
- (E2): The study was excluded if it uses mutation testing information to assess test quality (e.g. the use of mutation score to assess the quality of test sets generated based on testing criteria other than mutation testing).

We selected studies that first passed I3, and then either I1 or I2 ( $I3 \wedge (I1 \vee I2)$ ). We discarded studies that passed E1 or E2 ( $E1 \vee E2$ ). Notice that we use the term *approach* to refer to research that applies either classical cost reduction techniques (e.g. selective mutation and random mutation) or more contemporary techniques such as evolutionary algorithms and minimal mutation. As explained in Section 5.3, we found that several studies applied more than one technique in combination. Moreover, we selected studies that had a clear goal of reducing costs, either: (i) by explicitly or implicitly stating this goal in the approach description; (ii) while defining experimental goals; and/or (iii) by experimental measurements and analyses. We discarded studies that simply employed a particular mutation tool, or applied a subset of mutation operators, but did not measure cost reduction either theoretically or empirically.

#### 3.3. Search steps, search string, repositories, and selection procedures

**Step 1—Automatic search:** This step relied on search strings and search engines. The main terms in the string were “mutation testing” and “cost reduction,” which were used to compose the full search string as follows:

(“mutation testing” or “mutation analysis” or “mutant analysis”) **and** (“cost reduction” or “sufficient operator” or “sufficient mutation” or “constrained mutation” or “selective mutation” or “weak mutation” or “random selection” or “random mutation” or “random mutants” or “equivalent mutant” or “equivalent mutation”)

The search string, sometimes customized for specific databases, was submitted to the search engines of the following databases:<sup>3</sup> IEEE Xplore,<sup>4</sup> ACM Digital Library,<sup>5</sup> Elsevier ScienceDirect,<sup>6</sup> Springer SpringerLink,<sup>7</sup> and Wiley Online Library.<sup>8</sup> Additional databases are also sources of primary studies in the snowballing and author survey steps, as described next.

**Step 2—Snowballing:** The backward snowballing technique (Wohlin, 2014) was applied to the set of studies selected in Step 1. For each selected study, we analyzed the list of references to identify other studies of interest, according to the inclusion and exclusion criteria.

**Step 3—Author survey:** We contacted every author of a selected study from the previous steps and asked for suggestions of additional studies. We analyzed their suggestions according to the inclusion and exclusion criteria.

**Step 4—Post-publication update:** Our set of selected studies was updated after our preliminary publication (Ferrari et al., 2018a). Particularly, we analyzed studies published in the 13<sup>th</sup> International Workshop on Mutation Analysis (Mutation'18), 11<sup>th</sup> International Conference on Software Testing, Verification and Validation (ICST'18), and 40<sup>th</sup> International Conference on Software Engineering (ICSE'18). The updated dataset also includes studies suggested by authors previously contacted in our author survey in step 3.

**General study selection procedures:** For each search step, one researcher performed the preliminary selection by reading specific parts of the retrieved studies (initially, title and abstract, and, if necessary, introduction and conclusion) to check whether they should be selected for full reading. In the final selection step, two researchers fully read each study and made the final selection decision, according to the inclusion and exclusion criteria. All conflicts in the study selection were resolved by the two researchers in order to reduce potential threats.

#### 3.4. Data extraction, synthesis, and analysis

Data of interest were extracted and stored in customized forms. General data such as authors, study title, and year of publication were used to compose references. For this SLR, more specific pieces of information were necessary for the synthesis and analysis steps. The first consists of characteristics of the cost reduction techniques and the metrics used for cost reduction measurement. The second consists of the values obtained with respect to cost reduction savings and the achieved mutation scores.

<sup>3</sup> The original repositories of the selected studies can be found in the references at the end of this paper as well in the companion website (Ferrari et al., 2018b). It is important to note that IEEE, Elsevier, and Springer can be considered “pure” (or private) study repositories since they only store and retrieve studies from their own database. ACM and Wiley, on the other hand, have “hybrid” characteristics, since their search engines retrieve studies from a variety of databases, thus overlapping results from other searches.

<sup>4</sup> <http://ieeexplore.ieee.org> – last accessed on July, 2019.

<sup>5</sup> <https://dl.acm.org/> – last accessed on July, 2019.

<sup>6</sup> <http://www.sciencedirect.com/> – last accessed on July, 2019.

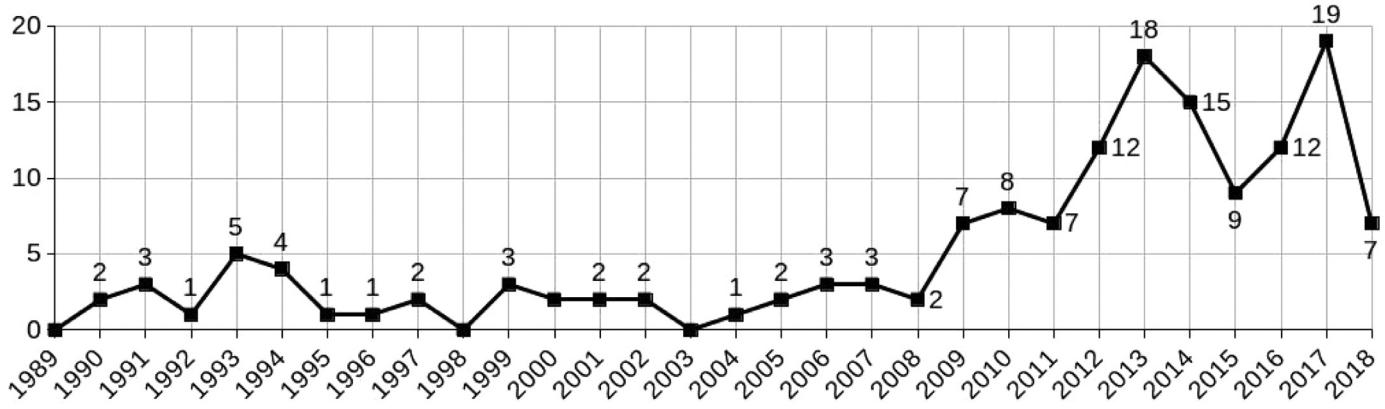
<sup>7</sup> <http://link.springer.com/> – last accessed on July, 2019.

<sup>8</sup> <http://onlinelibrary.wiley.com/> – last accessed on July, 2019.

**Table 1**

Number of studies per search round. Columns labeled with “Sel.” include non-duplicated results for final selection. Columns labeled with “Retr.” include duplicated search results. Note that the numbers of retrieved studies did not evolve consistently for Rounds 1 and 2 due to changes made in the search engines, and due to specifics of the grammar used to combine search terms in both rounds.

Search round	Database/Search engine												Total	
	IEEE		ACM		Elsevier		Springer		Wiley		Experts		..	..
	Retr.	Sel.	Retr.	Sel.	Retr.	Sel.	Retr.	Sel.	Retr.	Sel.	Retr.	Sel.	Retr.	Sel.
1 - Autom (Apr, 2016)	467 8.88%	27 17.65%	612 11.64%	19 12.42%	606 11.52%	12 7.84%	464 8.82%	0 0.00%	771 14.66%	13 8.50%	19 0.36%	13 8.50%	2939 55.89%	84 54.90%
2 - Autom (Oct, 2017)	303 5.76%	8 5.23%	642 12.21%	10 6.54%	706 13.42%	1 0.65%	78 1.48%	0 0.00%	376 7.15%	1 0.65%	0 0.00%	0 0.00%	2105 40.03%	20 13.07%
3 - Snowb (Oct, 2017)	-	-	-	-	-	-	-	-	-	-	-	-	72 1.37%	17 11.11%
4 - Survey (Oct-Dec, 2017)	-	-	-	-	-	-	-	-	-	-	-	-	135 2.57%	25 16.34%
5 - Update (Aug, 2018)	-	-	-	-	-	-	-	-	-	-	-	-	8 0.15%	7 4.58%
TOTAL	770 14.64%	35 22.88%	1254 23.84%	29 18.95%	1312 24.95%	13 8.50%	542 10.31%	0 0.00%	1147 21.81%	14 9.15%	19 0.36%	13 8.50%	5259 -	153 -

**Fig. 1.** Distribution of studies per year. Only part of 2018 is included.

#### 4. Study selection results

This section summarizes our results at each study selection step and the final set of selected studies. **Table 1** presents overall numbers for each step. **Fig. 1** shows the number of selected cost reduction studies per year from 1989 through (part of) 2018. **Table A.1** in [Appendix A](#) summarizes the studies in tabular form.

**Table 1** presents the numbers of studies retrieved and selected from each repository and search round. The term *search round* (or *round*) is used to refer to an automatic search round, the snowballing step, the author survey, and the last update performed after our preliminary publication of partial results ([Ferrari et al., 2018a](#)). Notice that **Table 1** displays only total numbers for the third, fourth, and fifth rounds since the studies were not retrieved from a particular database via automatic search.

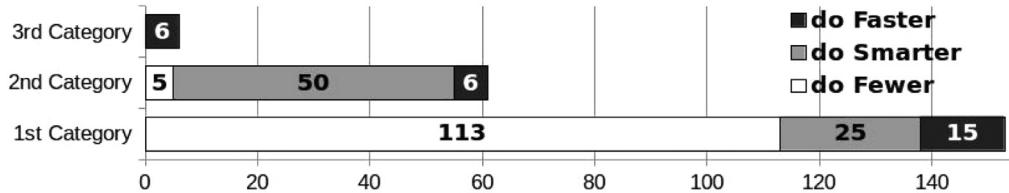
All rounds were performed with the support of the StArt tool ([Hernandes et al., 2012](#)), which supports major steps of the SLR process, with particular help for the preselection and final selection phases. The first round was performed in April 2016 and resulted in 84 selected studies. The second round was performed in October 2017 and resulted in 20 additional studies. Snowballing (third round) was performed over the set of studies selected in the first search round, and resulted in 17 additional studies. The author survey (fourth round) targeted authors of the three previous rounds. It started on 24 October 2017, and responses arrived until 18 December 2017 (54 days). In total, 42 authors replied and

provided 135 suggestions of additional studies. 25 additional, non-duplicate, studies were selected from this group. As described in [Section 3.3](#), our dataset was updated based on the analysis of studies published in three events (Mutation'18, ICST'18, and ICSE'18), and with studies suggested by authors previously contacted in our author survey. This update added 7 new studies (*Search Round 5* in **Table 1**).

It is important to note that **Table 1** only includes total numbers of non-subsumed studies, that is, it does not account for 22 studies that passed our inclusion criteria, but were subsumed (updated or extended) by more recent studies. References to the 22 subsumed studies can be found in [Appendix A](#) (**Table A.2**). Furthermore, as explained in [Section 3.3](#), due to their “hybrid” nature, ACM’s and Wiley’s search engines retrieved results that overlap results from other queried databases. Therefore, **Table 2** provides more precise information regarding the number of studies selected per database than **Table 1**.

**Table 2**  
Number of studies per database.

Database	# of Studies	Database	# of Studies		
IEEE	58	37.91%	Elsevier	17	11.11%
ACM	33	21.57%	Wiley	20	13.07%
Springer	15	9.80%	Others	10	6.54%
Total			153	100.00%	



**Fig. 2.** Number of studies per Offutt and Untch (2000)'s categories.

**Fig. 1** shows the distribution of selected studies per year. Note that only part of 2018 is included. It is clear the number of studies has been growing since 2009. Analyses of selected studies are presented in the next sections.

## 5. An overall characterization

This section presents and discusses an overall characterization of research into cost reduction of mutation testing. It starts in [Section 5.1](#) by presenting the results of an initial classification based on the traditional “*do fewer*,” “*do smarter*,” and “*do faster*” categories ([Offutt and Untch, 2000](#)). In [Section 5.2](#), we refine this classification into a set of primary goals associated with cost reduction. Then, [Section 5.3](#) describes techniques that have been developed to achieve those primary goals. Quantitative data are also presented in the following sections.

### 5.1. Traditional classification: *do fewer*, *do smarter*, *do faster*

[Offutt and Untch \(2000\)](#) grouped cost reduction techniques into three categories: *do fewer*, *do smarter*, and *do faster*. In their words:

*“The ‘do fewer’ approaches seek ways of running fewer mutant programs without incurring intolerable information loss. The ‘do smarter’ approaches seek to distribute the computational expense over several machines or factor the expense over several executions by retaining state information between runs or seek to avoid complete execution. The ‘do faster’ approaches focus on ways of generating and running each mutant program as quickly as possible.”*

We classified the 153selected studies according to [Offutt and Untch](#)'s categories, as summarized in [Fig. 2](#). The bar labelled “1st Category” shows that 113 studies used a *do fewer* approach, 25 used a *do smarter* approach, and 15 used *do faster*. This bar reflects the primary goal of the 153studies. However, some studies combined more than one goal. The bar labelled “2nd Category” shows that five studies used a second approach of *do fewer*, 50 had a second approach of *do smarter*, and six used *do faster*. For example, to automatically generate test cases, [Papadakis and Malevris \(2010b\)](#) used *control-flow analysis* supported by the *meta-mutants* technique. We classified this study as *do smarter* followed by *do faster*. To obtain the total number of studies that, for example, tried to *do fewer*, one can sum 113 (shown in the “1st Category” bar) with 5 (shown in the “2nd Category” bar).

[Fig. 3](#) shows the distribution of studies by year and category. We found more *do fewer* and *do smarter* approaches. For example, 17 out of 18 studies published in 2013 tried to run fewer mutants without significant effectiveness loss. In the same year, 11 out of 18 studies tried to reduce costs in a “smarter” way, usually by distributing computational expense, processing intermediate execution information from both the original program and its mutants, or by analyzing mutant execution statistics.

### 5.2. Primary goals for cost reduction

The *do fewer*, *do smarter*, and *do faster* categories have been used for almost 20 years. However, as cost reduction techniques

have become more complicated, this categorization has become harder to apply. Some techniques fit into more than one category and others do not really fit into any of the three. Thus, a major contribution of this paper is a new, and modern, characterization of primary goals for cost reduction techniques.

Our analysis resulted in the six primary goals described as follows. Each is labeled with a symbol that is used as a marker in subsequent figures. We also provide a complete list of references to studies that pursued each primary goal.

PG-1: (▽) *Reducing the number of mutants*: The objective is to reduce the number of mutants that will be executed, preferably without reducing effectiveness ([Petrović and Ivanković, 2018](#); [Zhu et al., 2018](#); [McMinn et al., 2018](#); [Abuljadayel and Wedyan, 2018](#); [Sun et al., 2017a](#); [Derezińska and Rudnik, 2017](#); [Delgado-Pérez et al., 2017](#); [Praphamontripang and Offutt, 2017](#); [Gopinath et al., 2017](#); [Delgado-Pérez et al., 2017c](#); [Just et al., 2017](#); [Al-Hajjaji et al., 2017](#); [Delgado-Pérez et al., 2017b](#); [Sun et al., 2017b](#); [Derezińska, 2016](#); [Quyen et al., 2016](#); [Parsai et al., 2016](#); [Kurtz et al., 2016](#); [Lima et al., 2016](#); [Derezińska and Hałas, 2015](#); [Reuling et al., 2015](#); [Bluemke and Kulesza, 2014b](#); [Harman et al., 2014](#); [Lacerda and Ferrari, 2014](#); [Ammann et al., 2014](#); [Delamaro et al., 2014a](#); [2014c](#); [Madeyski et al., 2014](#); [Delamaro et al., 2014b](#); [Bluemke and Kulesza, 2014a](#); [2013](#); [Oliveira et al., 2013](#); [Deng et al., 2013](#); [Cachia et al., 2013](#); [Inozemtseva et al., 2013](#); [Reales et al., 2013](#); [Zhang et al., 2013a](#); [Gligoric et al., 2013](#); [Derezińska and Rudnik, 2012](#); [Nobre et al., 2012](#); [Patrick et al., 2012](#); [Omar and Ghosh, 2012](#); [Wedyan and Ghosh, 2012](#); [Kaminski et al., 2011b](#); [Domínguez-Jiménez et al., 2011](#); [Papadakis and Malevris, 2010a](#); [Sridharan and Siami-Namin, 2010](#); [Kintis et al., 2010](#); [Zhang et al., 2010a](#); [Ji et al., 2009](#); [Kaminski and Ammann, 2009](#); [Polo et al., 2009](#); [Untch, 2009](#); [Siami-Namin et al., 2008](#); [Tuya et al., 2007](#); [Siami-Namin and Andrews, 2006](#); [Adamopoulos et al., 2004](#); [Vincenzi et al., 2001](#); [Barbosa et al., 2001](#); [Mresa and Bottaci, 1999](#); [Offutt et al., 1996](#); [Wong and Mathur, 1995](#); [Wong et al., 1994](#); [Offutt et al., 1993](#); [Marcozzi et al., 2018](#); [Zhu et al., 2017](#)).

PG-2: (≡) *Automatically detecting equivalent mutants*: The objective is to automatically identify which mutants are equivalent to the original program, and then eliminate them from consideration ([Marcozzi et al., 2018](#); [Devroey et al., 2017](#); [Kintis et al., 2017](#); [Holling et al., 2016](#); [Patel and Hierons, 2016](#); [Kintis et al., 2015](#); [Kintis and Malevris, 2015](#); [Papadakis et al., 2014](#); [Schuler and Zeller, 2013](#); [Ferrari et al., 2013](#); [Kintis and Malevris, 2013](#); [Papadakis and Le Traon, 2013](#); [Schuler et al., 2009](#); [Anbalagan and Xie, 2008](#); [Offutt et al., 2006](#); [Vincenzi et al., 2002](#); [Harman et al., 2000](#); [Hierons et al., 1999](#); [Offutt and Pan, 1997](#); [Offutt and Craft, 1994](#); [McMinn et al., 2018](#); [Delgado-Pérez et al., 2017](#); [Offutt and Lee, 1994](#)).

PG-3: (▷▷) *Executing faster*: The objective is to reduce execution time by using novel algorithms, tool improvements, or special-purpose hardware. Some techniques analyze each mutant to decide whether it can be partially executed or

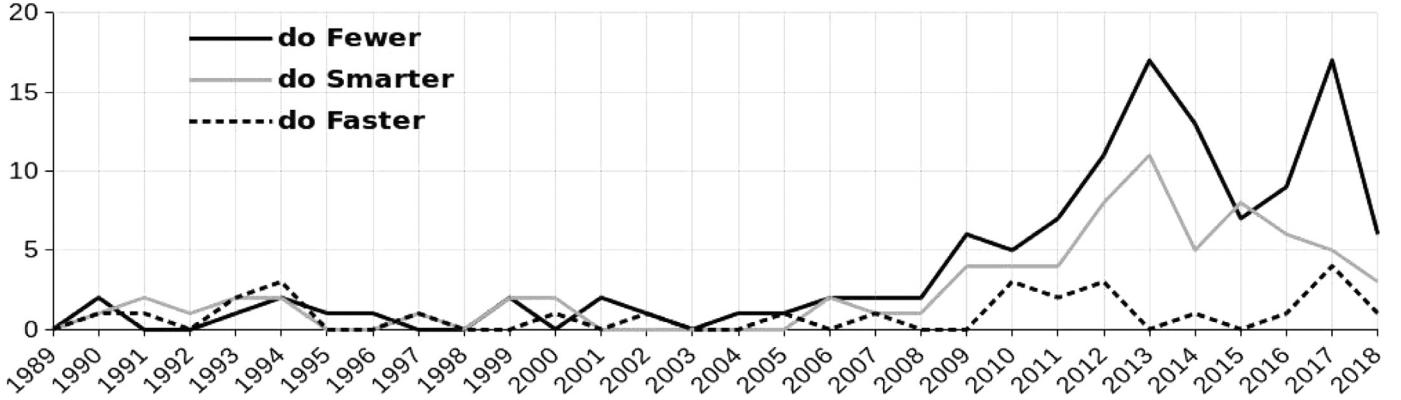


Fig. 3. Distribution of studies per year for Offutt and Untch (2000)'s categories.

even discarded without being executed (Chen and Zhang, 2018; Denisov and Pankevich, 2018; Wang et al., 2017; McMinn et al., 2016; Devroey et al., 2016; Gopinath et al., 2016b; Zhang et al., 2016; Li et al., 2015; Reales and Polo, 2014; Just et al., 2014a; Reales and Polo, 2013; Wright et al., 2013; Zhang et al., 2012; Gligoric et al., 2012; Kim et al., 2012; Durelli et al., 2012; Papadakis and Malevris, 2011b; Just et al., 2011; Gligoric et al., 2010; Bogacki and Walter, 2006a; Ma et al., 2005; Alexander et al., 2002; Jackson and Woodward, 2000; Untch et al., 1997; Fleyshgakker and Weiss, 1994; Offutt and Lee, 1994; Choi and Mathur, 1993; Untch et al., 1993; Weiss and Fleyshgakker, 1993; Offutt et al., 1992; Krauser et al., 1991; DeMillo et al., 1991; Devroey et al., 2017).

PG-4: (ε) Reducing the number of test cases or the number of executions: The objective is either to find smaller test sets that are still as effective at killing mutants, or to identify groups of similar mutants to reduce the number of test runs (Gopinath et al., 2018; Zhu et al., 2017; Ma and Kim, 2016; Fraser and Arcuri, 2015; Derezińska, 2013; Zhang et al., 2013b; Papadakis and Malevris, 2012; Kaminski and Ammann, 2011; Ayari et al., 2007; Sahinoğlu and Spafford, 1990; Zhu et al., 2018; Devroey et al., 2016; Derezińska and Hałas, 2015; Harman et al., 2014; Oliveira et al., 2013; Just et al., 2012b; Tuya et al., 2007; Adamopoulos et al., 2004).

PG-5: (⊖) Avoiding the creation of certain mutants: The objective is to define mutation operators or mutation generation algorithms to generate fewer mutants. The general idea is to generate only non-trivial (and non-equivalent) mutants (Delgado-Pérez et al., 2017a; Iida and Takada, 2017; Fernandes et al., 2017; Just and Schweiggert, 2015; Belli and Beyazit, 2015; Kintis and Malevris, 2014; Kaminski et al., 2013; Just et al., 2012b; Hu et al., 2011; Steimann and Thies, 2010; Domínguez-Jiménez et al., 2009a; Marshall et al., 1990).

PG-6: (τ) Automatically generating test cases: The objective is to generate test cases automatically, to kill as many mutants as possible. Test case generation is typically guided by characteristics of the mutants, and can substantially reduce the effort required to create tests, which is usually done by hand (Bashir and Nadeem, 2017; Matnei Filho and Vergilio, 2015; Aichernig et al., 2014; Henard et al., 2014; Aichernig et al., 2013; Fraser and Zeller, 2012; Papadakis and Malevris, 2011a; Zhang et al., 2010b; Papadakis and Malevris, 2010b; 2009; Liu et al., 2006; Baudry et al., 2005; DeMillo and Offutt, 1993; Papadakis and Malevris, 2011b; Harman et al., 2000; DeMillo and Offutt, 1991; Offutt et al., 1999).

Approaches that try to avoid the creation of certain mutants (PG-5) typically characterize circumstances during mutant generation that would lead to the creation of equivalent, trivial, and redundant mutants. These rules are then embedded into the code that implements the mutation operators to prevent such mutants from being created. Other approaches try to reduce the number of mutants (PG-1) post-creation. For these, the mutation operators are applied as originally specified; and the reduction of the mutant set is achieved either by constraining the set of operators or by strategically selecting mutants from the generated set (e.g. based on mutant characteristics or randomly). Clearly, achieving PG-5 contributes to achieving PG-1. In other words, some of these goals can be seen as “side-effects” of others. Another example involves PG-1, PG-4 and PG-6. Reducing the number of mutants (PG-1) implicitly reduces the number of test case executions (PG-4) and (possibly) the number of automatically generated test cases (PG-6). Fig. 4 summarizes the relationships among the primary goals. Each edge indicates that the primary goal in the target node is indirectly achieved when the goal in the source node is achieved. For example, if PG-1 is achieved, then PG-6 is indirectly achieved.

We also emphasize that there is no clear disjoint, one-to-one, mapping between our new primary goals and Offutt and Untch (2000)'s categories. For instance, PG-3 (Executing faster) can be achieved by (1) by using novel algorithms (e.g., meta-mutants (Untch et al., 1993)), (2) improved tools (e.g., compiler-integrated mutant generator and executor (DeMillo et al., 1991)), or (3) running mutants on special-purpose hardware (e.g. multi-processor machines (Krauser et al., 1991; Offutt et al., 1992; Wang et al., 2017)). These methods were classified by Offutt and Untch as do smarter (1 and 2) and do faster (3).

Fig. 5 shows the number of studies that used each primary goal. The chart reveals that reducing the number of mutants was the most common, followed by executing faster, then automatically detecting equivalent mutants.

### 5.3. Cost reduction techniques

This section presents a broad classification of cost reduction-related research according to 21 categories, or *techniques*. Some of the categories rely on mutation-specific research (such as *weak mutation* and *higher order mutation*), and others on classical software execution and analysis techniques (such as *parallel execution* and *control-flow analysis*). The timeline in Fig. 6 shows the year each technique was introduced.

Five techniques appear twice in Fig. 6: *random mutation* (Acree et al., 1979; Mathur and Wong, 1993); *higher order mutation* (Acree et al., 1979; Polo et al., 2009); *weak mutation* (Howden,

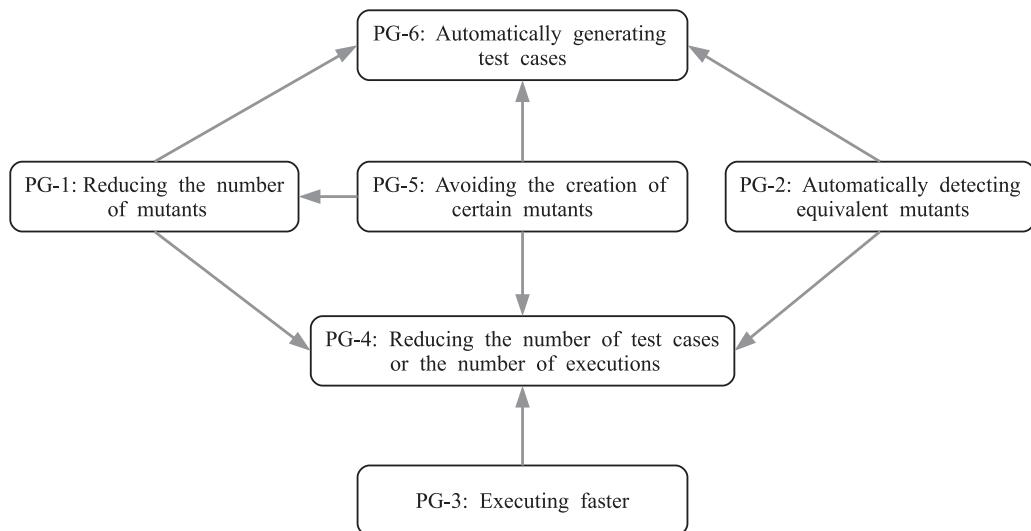


Fig. 4. Relationships among primary goals.

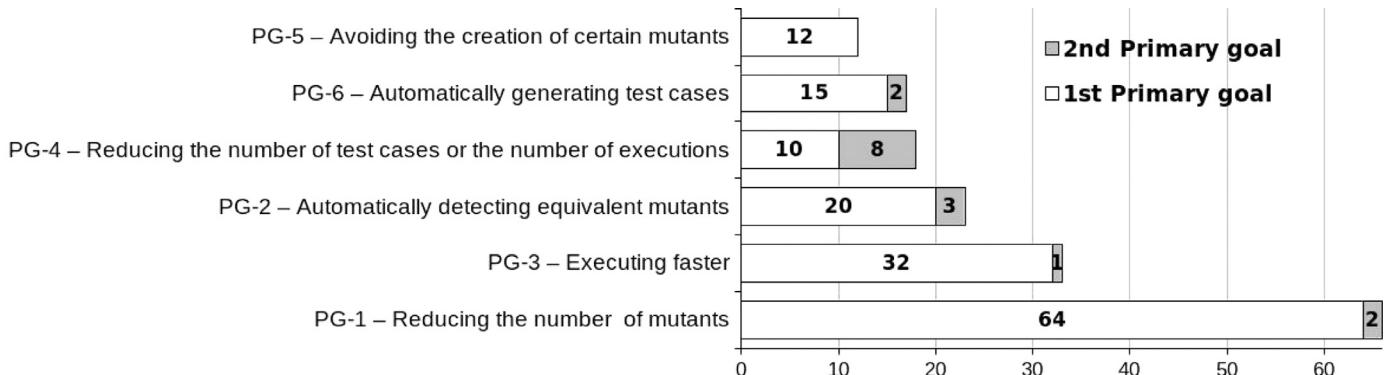


Fig. 5. Number of studies per primary goal.

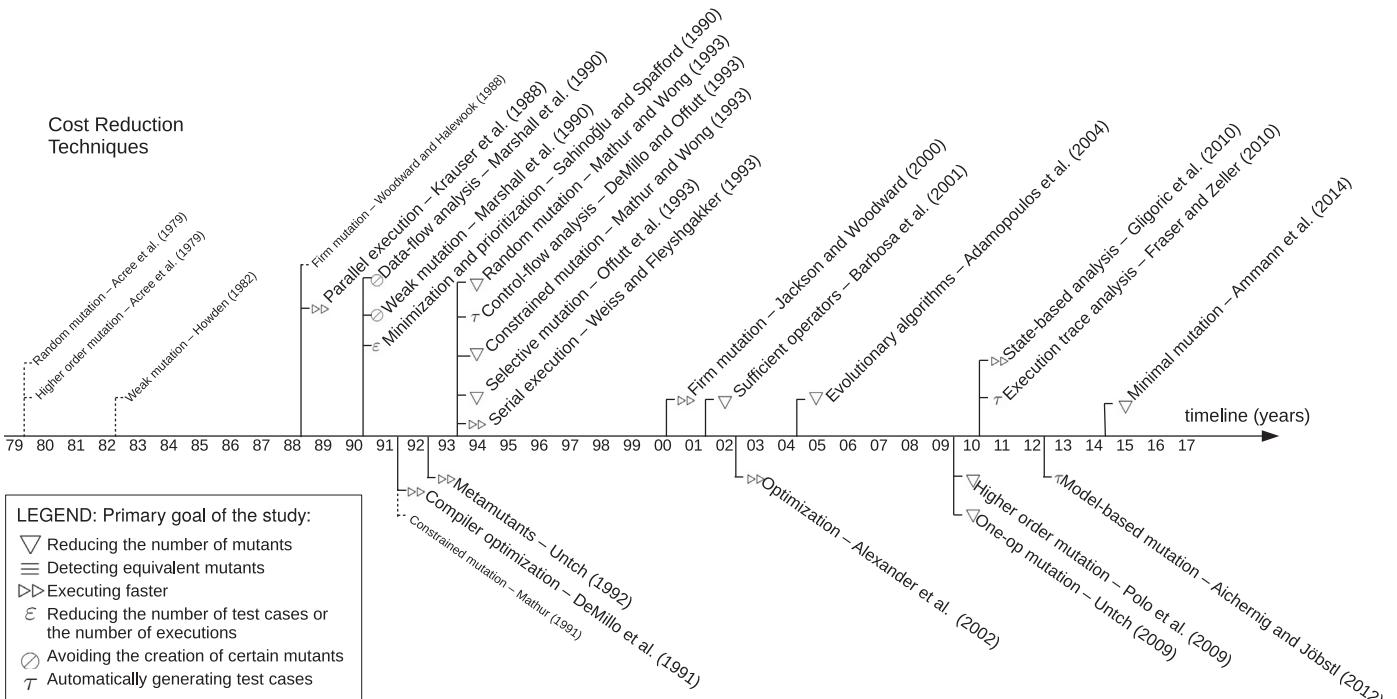


Fig. 6. Timeline for introduction of cost reduction techniques in peer-reviewed studies.

1982; Marshall et al., 1990); *firm mutation* (Woodward and Halewood, 1988; Jackson and Woodward, 2000); and *constrained mutation* (Mathur, 1991; Mathur and Wong, 1993). These techniques were first described either without a focus on cost reduction (Howden, 1982; Woodward and Halewood, 1988), in a high level way (Mathur, 1991), or in a non-peer-reviewed publication (Acree et al., 1979), and were later published with more details in peer-reviewed studies (Marshall et al., 1990; Jackson and Woodward, 2000; Mathur and Wong, 1993; Polo et al., 2009). These initial references (Acree et al., 1979; Howden, 1982; Woodward and Halewood, 1988; Mathur, 1991) appear with different notation (dotted lines and smaller font size). Also, the techniques in Fig. 6 are annotated with the symbols from the list of primary goals. For instance, *data-flow analysis* and *weak mutation* were first explored to avoid the creation of certain mutants ( $\emptyset$ ) (Marshall et al., 1990), whereas *metamutants* was first explored for executing faster ( $\gg$ ) (Untch, 1992).

The following list describes the 21 techniques. The descriptions are brief and omit details that were published in later years. Some categories have some overlap (including *selective mutation*, *sufficient operators*, *one-op mutation*, and *minimal mutation*), as discussed in Section 5.4.

- T-1: *Random mutation*: This technique selects randomly from the complete set of mutants according to a predefined probability distribution. Random mutation has been interpreted as “choose X% of all mutants,” “for each mutant, generate it with X% probability,” and as “choose X% of mutants generated by each operator” (Acree et al., 1979; Mathur and Wong, 1993; Derezińska and Rudnik, 2017; Gopinath et al., 2017; Parsai et al., 2016; Bluemke and Kulesza, 2013; Papadakis and Malevris, 2010a; Zhang et al., 2010a; Wong and Mathur, 1995; Zhang et al., 2013a; Petrović and Ivanković, 2018; Kurtz et al., 2016).
- T-2: *Higher order mutation*: This technique combines two or more simple mutations to create a single complex mutant (Acree et al., 1979; Gopinath et al., 2018; Lima et al., 2016; Reuling et al., 2015; Harman et al., 2014; Madeyski et al., 2014; Reales et al., 2013; Omar and Ghosh, 2012; Kaminski et al., 2011b; Kintis et al., 2010; Polo et al., 2009; Abuljadayel and Wedyan, 2018; Derezińska, 2016; Kaminski and Ammann, 2011; Papadakis and Malevris, 2010a; Devroey et al., 2016; Kintis et al., 2015).
- T-3: *Weak mutation*: This execution technique checks whether the state of the mutant is infected shortly after the mutated location has been executed, rather than checking the output after execution ends. If the state is infected, the mutant is killed immediately (Howden, 1982; Zhu et al., 2017; Papadakis and Malevris, 2011b; Offutt and Lee, 1994; Zhu et al., 2018; Devroey et al., 2017; Ma and Kim, 2016; Fraser and Arcuri, 2015; Zhang et al., 2010b; Kintis et al., 2010; Marshall et al., 1990; Kim et al., 2012).
- T-4: *Firm mutation*: This is a variant of weak mutation where execution is allowed to proceed for some pre-defined duration after the state is infected (Woodward and Halewood, 1988; Jackson and Woodward, 2000; Offutt and Lee, 1994).
- T-5: *Parallel execution*: This technique executes mutants in parallel processors, reducing the total time needed to perform mutation analysis (Li et al., 2015; Reales and Polo, 2013; Choi and Mathur, 1993; Offutt et al., 1992; Krauser et al., 1991; Wang et al., 2017; Jackson and Woodward, 2000).
- T-6: *Data-flow analysis*: This technique uses program data flow-related information to decide which mutants to generate and to analyze mutants. It considers whether variables that are more prone to failure during execution are reached and referenced (Al-Hajjaji et al., 2017; Kintis and Malevris, 2015;

Papadakis et al., 2014; Kintis and Malevris, 2014; Schuler and Zeller, 2013; Papadakis and Le Traon, 2013; Just et al., 2012b; Patrick et al., 2012; Schuler et al., 2009; Offutt et al., 2006; Hierons et al., 1999; Offutt and Craft, 1994; Marshall et al., 1990; Delgado-Pérez et al., 2017a; Kintis et al., 2015; Harman et al., 2000).

- T-7: *Minimization and prioritization of test sets*: This technique analyzes the test suite to score test cases based on their effectiveness at killing mutants, then either eliminates test cases that are ineffective or runs the most effective test cases before the less effective test cases (Sahinoğlu and Spafford, 1990; Derezińska, 2013; Tuya et al., 2007).
- T-8: *Compiler optimization*: This technique uses compiler-related techniques to optimize mutant execution and analysis (for instance, to automatically detect equivalent mutants) (Denisov and Pankevich, 2018; Kintis et al., 2017; DeMillo et al., 1991; Delgado-Pérez et al., 2017; Just et al., 2011; Offutt and Craft, 1994).
- T-9: *Constrained mutation*: This technique chooses a subset of mutation operators to use. The choice relies on testers' intuition regarding the significance of particular groups of mutants (Mathur, 1991; Mathur and Wong, 1993; Petrović and Ivanković, 2018; Wong et al., 1994; Gopinath et al., 2017; Wong and Mathur, 1995).
- T-10: *Metamutants (or mutant schemata)*: This technique generates and executes mutants by embedding all mutants in one parameterized program called a *metamutant*. The metamutant is then compiled for fast execution. When run, the metamutant takes a parameter that tells it which mutant to run (Untch, 1992; Gopinath et al., 2016b; Reales and Polo, 2014; Wright et al., 2013; Ma et al., 2005; Untch et al., 1997; Kim et al., 2012; Papadakis and Malevris, 2010b; Weiss and Fleysgakker, 1993; Ma and Kim, 2016; Papadakis and Malevris, 2011b).
- T-11: *Control-flow analysis*: This technique uses program control flow-related information, focusing on execution characteristics to identify branches and commands that help determine which structures are most relevant to the generation and execution of mutants (Marcozzi et al., 2018; Chen and Zhang, 2018; McMinn et al., 2018; Delgado-Pérez et al., 2017a; Just et al., 2017; Sun et al., 2017b; Holling et al., 2016; Patel and Hierons, 2016; Kintis et al., 2015; Zhang et al., 2012; Papadakis and Malevris, 2012; Kaminski and Ammann, 2011; Zhang et al., 2010b; Papadakis and Malevris, 2010b; 2009; Ji et al., 2009; Liu et al., 2006; Harman et al., 2000; Offutt and Pan, 1997; DeMillo and Offutt, 1993; Petrović and Ivanković, 2018; Papadakis et al., 2014; Schuler and Zeller, 2013; Papadakis and Le Traon, 2013; Just et al., 2012b; Patrick et al., 2012; Papadakis and Malevris, 2011b; Schuler et al., 2009; Offutt et al., 2006; Bashir and Nadeem, 2017; Fraser and Arcuri, 2015; DeMillo and Offutt, 1991; Offutt et al., 1999).
- T-12: *Selective mutation*: This technique tries to avoid the application of mutation operators that are responsible for the most mutants or to select mutation operators that result in mutants that are killed by tests that also kill lots of mutants created by other operators. The idea is that if a test set,  $T_{op}$ , that is adequate for a subset of mutation operators  $op$ , also kills a very high percentage of all mutants, then we can select only the operators in  $op$  (Sun et al., 2017a; Praphamontripong and Offutt, 2017; Delgado-Pérez et al., 2017c; Bluemke and Kulesza, 2014b; 2014a; Zhang et al., 2013a; Gligoric et al., 2013; Derezińska and Rudnik, 2012; Sridharan and Siami-Namin, 2010; Mresa and Bottaci, 1999; Offutt et al., 1996; 1993; Kurtz et al., 2016; Zhang et al., 2010a; Tuya et al., 2007; Delgado-Pérez et al., 2017; Gopinath et al., 2017).

- T-13: *Serial execution*: This technique dynamically determines classes of mutants that behave similarly, thus decreasing the number of mutants to be executed. The overhead incurred needs to be small relative to the time saved (Ma and Kim, 2016; Kim et al., 2012; Fleyshgakker and Weiss, 1994; Weiss and Fleyshgakker, 1993).
- T-14: *Sufficient operators*: This technique tries to determine an essential set of mutation operators by applying customized procedures (Just and Schweiggert, 2015; Lacerda and Ferrari, 2014; Delamaro et al., 2014b; Siami-Namin et al., 2008; Siami-Namin and Andrews, 2006; Vincenzi et al., 2001; Barbosa et al., 2001; Zhang et al., 2010a). It can be seen as a special case of *selective mutation*, where the main difference is the complexity of the procedures applied to identify the final subset of mutation operators, which may be based, for example, on heuristics or statistics).
- T-15: *Optimization of generation, execution, and analysis of mutants*: This technique groups approaches that reduce the cost of mutation testing by exploring strategies that did not fit other categories on this list. This category appears as *Optimization* in Fig. 6 for brevity. More details are given later in this section (Zhu et al., 2018; Iida and Takada, 2017; Fernandes et al., 2017; McMinn et al., 2016; Zhang et al., 2016; Derezińska and Hałas, 2015; Belli and Beyazit, 2015; Just et al., 2014a; Derezińska, 2013; Kaminski et al., 2013; Ferrari et al., 2013; Kintis and Malevris, 2013; Cachia et al., 2013; Inozemtseva et al., 2013; Zhang et al., 2013b; Wedyan and Ghosh, 2012; Gligoric et al., 2012; Durelli et al., 2012; Hu et al., 2011; Just et al., 2011; Steimann and Thies, 2010; Kaminski and Ammann, 2009; Anbalagan and Xie, 2008; Bogacki and Walter, 2006a; Vincenzi et al., 2002; Alexander et al., 2002; Chen and Zhang, 2018; Zhu et al., 2017; Devroey et al., 2016; Gopinath et al., 2016b; Just and Schweiggert, 2015; Reuling et al., 2015; Wright et al., 2013; Zhang et al., 2012).
- T-16: *Evolutionary algorithms*: This technique uses evolutionary algorithms to reduce the number of mutants, to reduce the number of test cases, or to identify equivalent mutants (Abuljadayel and Wedyan, 2018; Delgado-Pérez et al., 2017; Bashir and Nadeem, 2017; Delgado-Pérez et al., 2017b; Quyen et al., 2016; Matnei Filho and Vergilio, 2015; Fraser and Arcuri, 2015; Henard et al., 2014; Oliveira et al., 2013; Nobre et al., 2012; Fraser and Zeller, 2012; Papadakis and Malevris, 2011a; Domínguez-Jiménez et al., 2011; Domínguez-Jiménez et al., 2009a; Ayari et al., 2007; Baudry et al., 2005; Adamopoulos et al., 2004; Harman et al., 2014; Sridharan and Siami-Namin, 2010; Ji et al., 2009; Papadakis and Malevris, 2011b).
- T-17: *One-op mutation*: This technique uses only a single mutation operator, which has the advantage of producing few mutants but also providing comprehensive coverage of the program (Derezińska, 2016; Delamaro et al., 2014a; 2014c; Deng et al., 2013; Untch, 2009; Kurtz et al., 2016).
- T-18: *Execution trace analysis*: This technique uses traces of execution of the original program and some mutants to decide which of the remaining mutants should be executed (Reales and Polo, 2014; Fraser and Zeller, 2012).
- T-19: *Model-based mutation*: This technique mutates formal or informal models of the program, and then uses the mutants to automatically generate test cases that are later used to kill mutants of the program (Devroey et al., 2017; 2016; Aichernig et al., 2014; 2013).
- T-20: *State-based analysis*: This technique compares states of different mutant executions. When two mutants lead to the same mutation state, that is, when the same execution path is observed, only one needs to be executed and the result of

the other can be inferred. Similarly, the technique creates groups of classes that define certain transition sequences such that only one needs to be verified (Wang et al., 2017; Gligoric et al., 2010; Bashir and Nadeem, 2017; Just et al., 2014a; Gligoric et al., 2012).

- T-21: *Minimal mutation*: This technique identifies and eliminates redundant mutants by applying the concepts of mutant subsumption and dominator mutants (Kurtz et al., 2016; Ammann et al., 2014; Just et al., 2017).

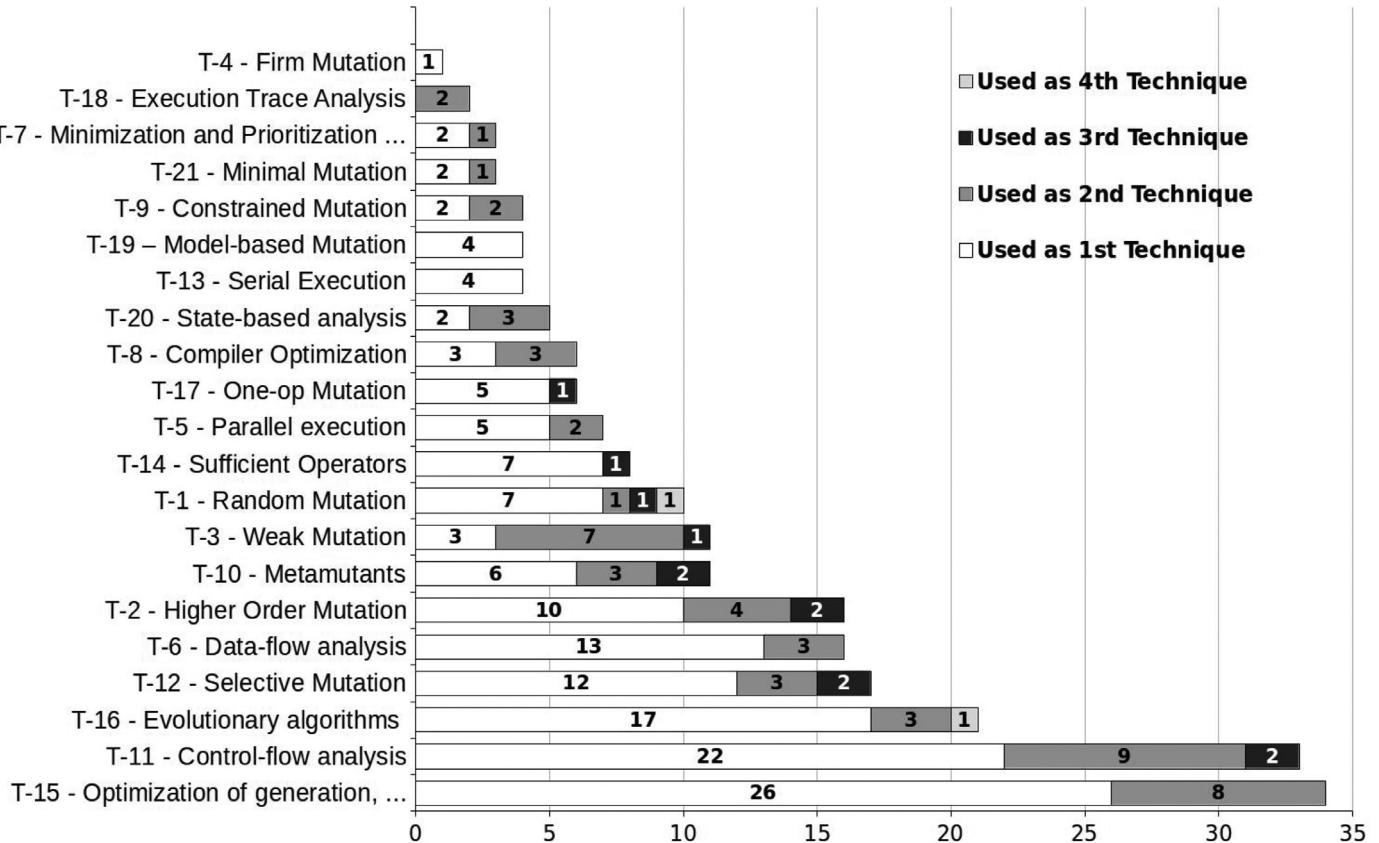
Note that T-15 (*optimization of generation, execution, and analysis of mutants*) classifies studies that do not fit into any other category. For instance, Alexander et al. (2002) implemented a mutation tool that uses Java reflection to mutate objects at run-time, thus avoiding the need to compile separate mutated programs. Durelli et al. (2012) evolved the Java Virtual Machine to embed native support to speed up the execution of the original program and its mutants. Derezińska (2013) proposed a clustering algorithm that uses results of mutant execution on subsets of test cases. A more recent example is Zhang et al. (2016)'s work, which used machine learning to devise predictive models to avoid the execution of some mutants. We classified these and other similar techniques as *optimization-related*. These are described in Section 5.5.

Fig. 7 shows the distribution of studies per technique. Some studies used two or more techniques (discussed further in Section 5.4). Studies that used more than one technique are listed in Fig. 7 with counts. For example, *minimal mutation* (T-21) has two units for *Used as 1st Technique* and one unit for *Used as 2nd Technique*. The amount of emphasis of each technique is based on the focus authors gave to that particular technique. As an example, Marshall et al. (1990) reduced mutation costs by applying *data-flow analysis* to support *weak mutation*. Thus, *data-flow analysis* was the main technique while *weak mutation* was secondary. In another example, Wong and Mathur (1995) compared *random mutation* with *constrained mutation*, without any clear preference for either. In this case, assigning higher participation for a particular technique was arbitrary. Therefore, even though Fig. 7 shows varied levels of participation for most techniques, both techniques might have had similar importance.

Fig. 8 shows the distribution of the nine most investigated techniques per year, split into three charts to improve readability. Even without statistical tests for time series, it is clear that researchers are studying some techniques more in recent years, such as *evolutionary algorithms*, *control-flow analysis*, *higher order mutation*, and *selective mutation*. Fig. 9 illustrates this by showing the distribution of the six techniques most investigated over the five last years (2014–2018).

#### 5.4. Additional notes about the cost reduction techniques and primary goals

As mentioned in Section 5.3, some of the techniques are similar. For instance, *constrained mutation* (Mathur, 1991), *selective mutation* (Offutt et al., 1993), *sufficient operators* (Barbosa et al., 2001), *one-op mutation* (Untch, 2009), and the recent *minimal mutation* technique (Ammann et al., 2014) all reduce the number of mutants in some way, as far as possible without reducing effectiveness. After considering putting them all into one large category, we decided that a fine-grained categorization would be more helpful, and that small categories would be simpler to understand. Thus, we divided them on fairly specific characteristics. For example, *constrained mutation* relies on the testers' intuition, while *one-op mutation* relies on individual mutation operators. Along the same lines, we chose to separate *weak mutation* (Howden, 1982) and *firm mutation* (Woodward and Halewood, 1988), even though



**Fig. 7.** Number of studies per cost reduction technique.

weak mutation could be viewed as a special case of firm mutation.

Another interesting observation is that the frequency of studies that combine techniques has increased over time. Fig. 10 shows that techniques have been combined since research on cost reduction started, but this combined usage has increased since 2010. In 2009, two studies combined two techniques (Schuler et al., 2009; Ji et al., 2009), but in 2010 five studies combined two techniques (Kintis et al., 2010; Papadakis and Malevris, 2010a; Sridharan and Siami-Namin, 2010; Papadakis and Malevris, 2010b; Zhang et al., 2010b), and one combined three (Zhang et al., 2010a). By 2017, six studies combined two techniques (Wang et al., 2017; Just et al., 2017; Zhu et al., 2017; Devroey et al., 2017; Delgado-Pérez et al., 2017a; Abuljadayel and Wedyan, 2018), and three combined three (Bashir and Nadeem, 2017; Delgado-Pérez et al., 2017; Gopinath et al., 2017). We also found a few studies that combined four techniques starting in 2011 (Papadakis and Malevris, 2011b; Kurtz et al., 2016). In total, of the 107 studies selected since 2010, 45 combined at least two techniques. Table 3 lists the studies that combined two or more techniques per year (the techniques are shown between parentheses).

Table 4 relates primary goals, cost reduction techniques, and studies. For each primary goal, the table lists all studies and the applied cost reduction techniques. The number of times a study is listed for a given goal corresponds to the number of applied techniques applied in that study, irrespective of which technique was applied to achieve a particular goal. For example, for reducing the number of mutants (PG-1), Gopinath et al. (2017) applied three techniques: t-1 (random mutation), T-9 (constrained mutation) and T-12 (selective mutation). Therefore, Gopinath et al.'s study is listed three times for PG-1.

Note that some studies tried to achieve two primary goals and applied two or more techniques. An example is by Just et al. (2012b). That study applied techniques T-6 (Data-flow analysis) and T-11 (Control-flow analysis) to reduce the number of test cases or the number of executions (PG-4), and to avoid the creation of certain mutants (PG-5). In cases like this, the combination of techniques may allow one to achieve a combination of primary goals, therefore we do not distinguish between particular techniques applied to achieve particular goals.

Table 5 relates cost reduction techniques, primary goals, and studies. For each technique, the table lists all studies that addressed a particular goal. Similarly to Table 4, the number of times a study is listed corresponds to the number of primary goals it pursued times the number of techniques it applied (e.g. Just et al. (2012b)'s study is listed two times for T-6 and two times for T-11).

The reader should notice that providing a complete list of tools that implement these cost reduction techniques is outside the scope of this paper. Readers could use information in this paper to identify the tools used in the experiments. Examples of useful information are the lists of references to studies with respect to primary goals (Section 5.2), techniques (Section 5.3), and the relationships between them (Tables 4 and 5).

### 5.5. Overview of selected studies

We describe each of the 153 selected studies in detail on the companion website to this paper (Ferrari et al., 2018b).<sup>9</sup> We

<sup>9</sup> <http://goo.gl/edyF1n> –last accessed November 2018.

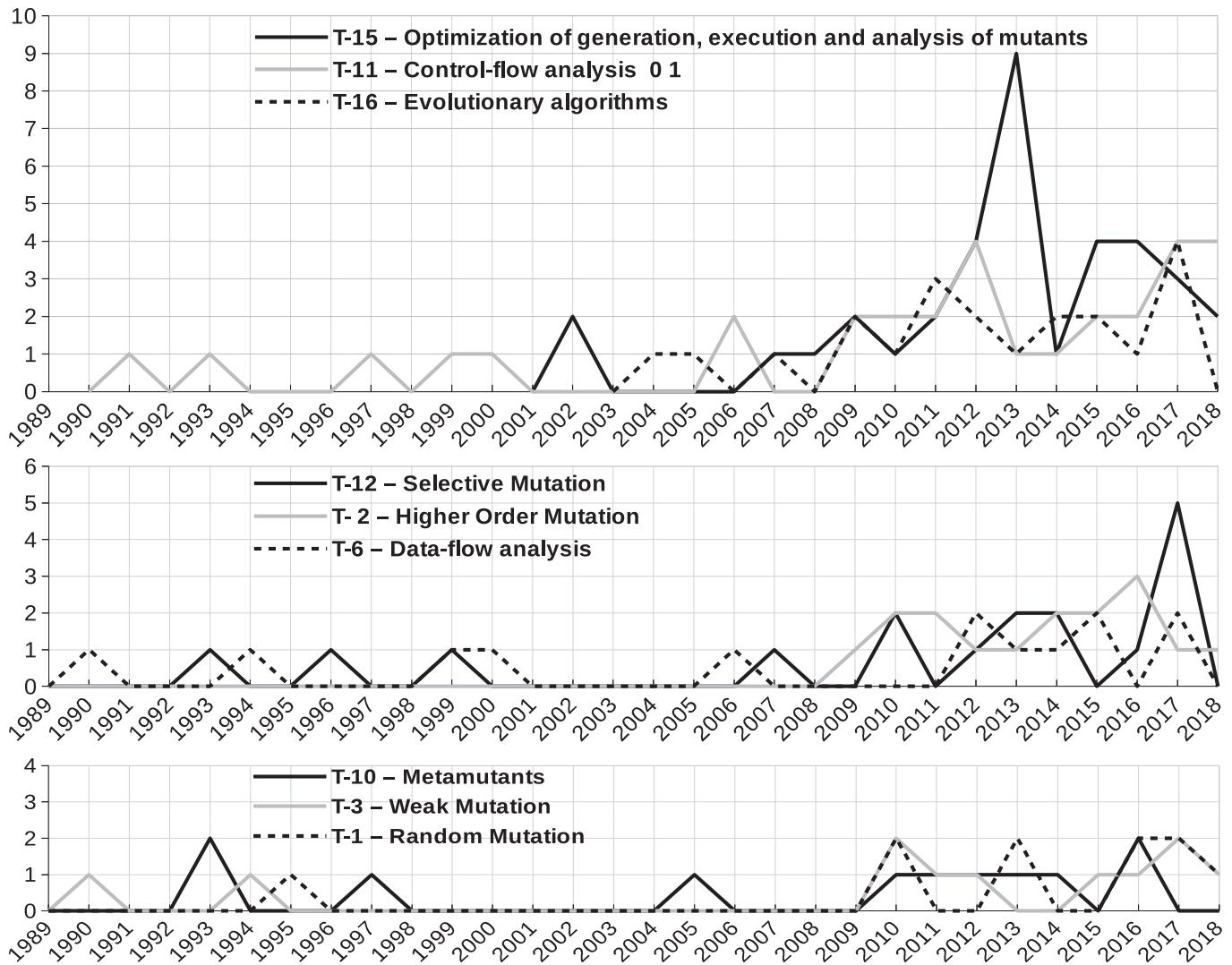


Fig. 8. Distribution of studies of the top nine techniques per year.

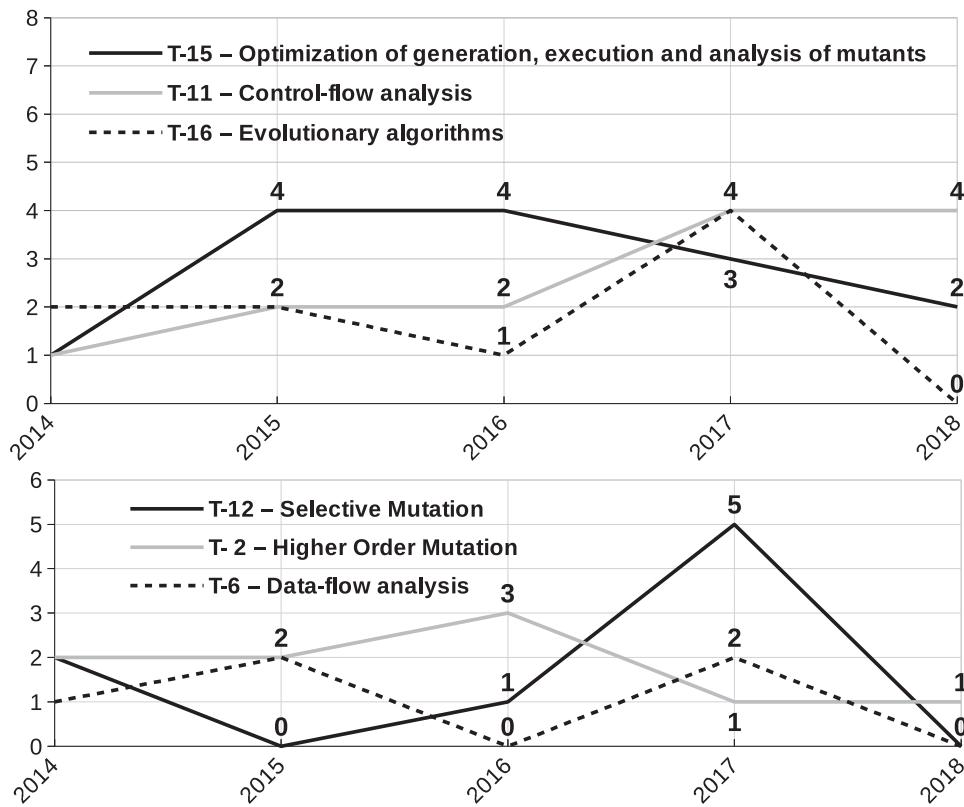
include a BibTeX file with an entry for every study, and a description field with more details than in this paper. The descriptions include the cost reduction goals of the studies (either explicit or implicit), and the achieved results. Every description also includes the cost reduction technique as described in Section 5.3, and the technology and artifacts addressed (if provided by the original authors). As an example, the study of Offutt and Craft (1994) is described in the complementary material as: ‘‘Offutt and Craft (1994) had the goal of detecting equivalent mutants. They used six techniques developed for compiler optimization. They analyzed the flow of data in mutants to identify equivalent mutants. In their experiment, they used 15 small Fortran programs and 14 mutation operators. The study automatically detected an average of 19.80% of the equivalent mutants.’’

The remainder of this section provides a more concise description of the selected studies, grouped by cost reduction goal and applied cost reduction technique. The groups of studies are organized chronologically, with a few exceptions when very interrelated studies are described within a given group (e.g. the studies by Domínguez-Jiménez et al. (2011) and Delgado-Pérez et al. (2017b), which explored Evolutionary Mutation Testing to reduce the number of mutants to be executed). When possible, we established re-

lationships between studies that extended prior studies. For the sake of completeness, studies that pursued two primary goals appear twice, each one in the subsection for each goal.

#### 5.5.1. Studies that tried to reduce the number of mutants (PG-1)

Offutt et al. (1993) extended and used the Mothra (DeMillo and Offutt, 1991) tool to experiment with several selective mutation options on 10 Fortran programs. Offutt et al. (1993) by establishing a reduced set of operators (ABS, UOI, LCR, AOR, and ROR). Mresa and Bottaci (1999) proposed a new type of selective mutation that takes into account the cost of detecting equivalent mutants. Tuya et al. (2007) proposed mutation operators for SQL. They used selective mutation and minimization and prioritization of test sets to reduce the number of mutants and the number of test cases. Zhang et al. (2010a) empirically compared random mutation with selective mutation and sufficient operators. Results of experiments with 7 medium-sized C programs provided hints that 7%-random selection would achieve similar test effectiveness. Sridharan and Siami-Namin (2010) used the same programs used by Zhang et al. to experiment with a Bayesian (evolutionary) approach to select the operators that



**Fig. 9.** Distribution of studies of the top six techniques in the last five years.

**Table 3**

List of studies that applied combined techniques per year (years with no studies are omitted).

Year	Studies and Respective Cost Reduction Techniques
1990	Marshall et al. (1990) (T-6,T-3)
1993	Weiss and Fleyshgakker (1993) (T-13,T-10)
1994	Offutt and Craft (1994) (T-6,T-8)
1995	Wong and Mathur (1995) (T-1,T-9)
2000	Harman et al. (2000) (T-11,T-6); Jackson and Woodward (2000) (T-4,T-5)
2006	Offutt et al. (2006) (T-6,T-11)
2007	Tuya et al. (2007) (T-7,T-12)
2009	Schuler et al. (2009) (T-6,T-11); Ji et al. (2009) (T-11,T-16)
2010	Zhang et al. (2010a) (T-1,T-12,T-14); Kintis et al. (2010) (T-2,T-3); Papadakis and Malevris (2010a) (T-1,T-2); Sridharan and Siami-Namin (2010) (T-12,T-16); Papadakis and Malevris (2010b) (T-11,T-10); Zhang et al. (2010b) (T-11,T-3)
2011	Just et al. (2011) (T-15,T-8); Kaminski and Ammann (2011) (T-11,T-2); Papadakis and Malevris (2011b) (T-3,T-11,T-10,T-16)
2012	Kim et al. (2012) (T-13,T-10,T-3); Gligorac et al. (2012) (T-15,T-20); Zhang et al. (2012) (T-11,T-15); Just et al. (2012b) (T-6,T-11); Patrick et al. (2012) (T-6,T-11); Fraser and Zeller (2012) (T-16,T-18)
2013	Papadakis and Le Traon (2013) (T-6,T-11); Zhang et al. (2013a) (T-12,T-1); Schuler and Zeller (2013) (T-6,T-11); Wright et al. (2013) (T-10,T-15); Derezińska (2013) (T-15,T-7)
2014	Just et al. (2014a) (T-15,T-20); Papadakis et al. (2014) (T-6,T-11); Harman et al. (2014) (T-2,T-16); Reales and Polo (2014) (T-10,T-18)
2015	Reuling et al. (2015) (T-2,T-15); Kintis et al. (2015) (T-11,T-6,T-2); Fraser and Arcuri (2015) (T-16,T-3,T-11); Just and Schweigert (2015) (T-14,T-15)
2016	Ma and Kim (2016) (T-13,T-3,T-10); Kurtz et al. (2016) (T-21,T-12,T-17,T-1); Gopinath et al. (2016b) (T-10,T-15); Devroey et al. (2016) (T-19,T-15,T-2); Derezińska (2016) (T-17,T-2)
2017	Wang et al. (2017) (T-20,T-5); Just et al. (2017) (T-11,T-21); Zhu et al. (2017) (T-3,T-15); Bashir and Nadeem (2017) (T-16,T-20,T-11); Devroey et al. (2017) (T-19,T-3); Gopinath et al. (2017) (T-1,T-9,T-12); Delgado-Pérez et al. (2017a) (T-11,T-6); Delgado-Pérez et al. (2017) (T-16,T-8,T-12); Abuljadayel and Wedyan (2018) (T-16,T-2)
2018	Chen and Zhang (2018) (T-11,T-15); Zhu et al. (2018) (T-15,T-3); Petrović and Ivanković (2018) (T-9,T-11,T-1)

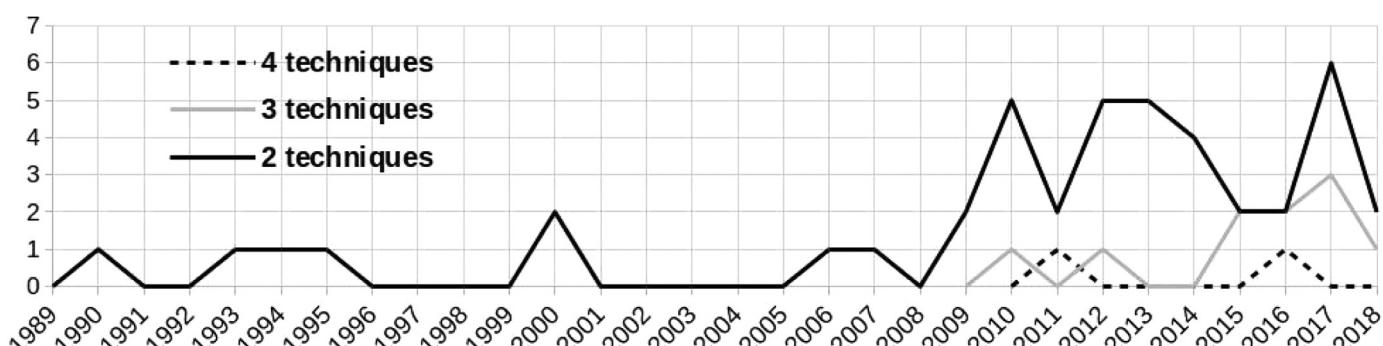
generated mutants that were hard to kill. Derezińska and Rudnik (2012) investigated selective mutation for C# programs. Gligorac et al. (2013) applied selective mutation to concurrent Java programs. Zhang et al. (2013a) applied selective mutation and subsequently applied random mutation. Bluemke and Kulesza (2014a,b) also explored selective mutation and a combination of

tools for test case generation, mutant generation, and program analysis. Kurtz et al. (2016) claimed that selective mutation must be specialized to avoid noise introduced by redundant mutants. Gopinath et al. (2017) theoretically and empirically compared several cost reduction techniques (selective mutation, constrained mutation, and random mutation), and used statistical analysis to

**Table 4**

Mapping between primary goals and cost reduction techniques.

Primary Goal	Cost Reduction Techniques and Respective Studies
PG-1 – Reducing the number of mutants	T-1 (Zhang et al., 2010a; 2013a; Wong and Mathur, 1995; Papadakis and Malevris, 2010a; Kurtz et al., 2016; Parsai et al., 2016; Gopinath et al., 2017; Derezińska and Rudnik, 2017; Petrović and Ivanković, 2018), T-2 (Kintis et al., 2010; Kaminski et al., 2011b; Polo et al., 2009; Reales et al., 2013; Reuling et al., 2015; Papadakis and Malevris, 2010a; Omar and Ghosh, 2012; Madeyski et al., 2014; Harman et al., 2014; Lima et al., 2016; Derezińska, 2016; Abuljadayel and Wedyan, 2018), T-3 (Kintis et al., 2010; Zhu et al., 2017; 2018), T-6 (Patrick et al., 2012; Al-Hajjaji et al., 2017), T-7 (Tuya et al., 2007), T-8 (Delgado-Pérez et al., 2017), T-9 (Wong and Mathur, 1995; Wong et al., 1994; Gopinath et al., 2017; Petrović and Ivanković, 2018), T-11 (Sun et al., 2017b; Patrick et al., 2012; Ji et al., 2009; Just et al., 2017; Marcozzi et al., 2018; Petrović and Ivanković, 2018; McMinn et al., 2018), T-12 (Zhang et al., 2010a; Offutt et al., 1993; Tuya et al., 2007; Mresa and Bottaci, 1999; Offutt et al., 1996; Gligoric et al., 2013; Bluemke and Kulesza, 2014a; Zhang et al., 2013a; Sridharan and Siami-Namin, 2010; Kurtz et al., 2016; Delgado-Pérez et al., 2017c; Gopinath et al., 2017; Praphamontripong and Offutt, 2017; Derezińska and Rudnik, 2012; Bluemke and Kulesza, 2014b; Delgado-Pérez et al., 2017; Sun et al., 2017a), T-14 (Zhang et al., 2010a; Siami-Namin et al., 2008; Vincenzi et al., 2001; Barbosa et al., 2001; Delamaro et al., 2014b; Siami-Namin and Andrews, 2006; Lacerda and Ferrari, 2014), T-15 (Wedyan and Ghosh, 2012; Inozemtseva et al., 2013; Reuling et al., 2015; Cachia et al., 2013; Kaminski and Ammann, 2009; Zhu et al., 2017; Derezińska and Hałas, 2015; Zhu et al., 2018), T-16 (Domínguez-Jiménez et al., 2011; Oliveira et al., 2013; Sridharan and Siami-Namin, 2010; Adamopoulos et al., 2004; Ji et al., 2009; Harman et al., 2014; Delgado-Pérez et al., 2017b; Quyen et al., 2016; Nobre et al., 2012; Delgado-Pérez et al., 2017; Abuljadayel and Wedyan, 2018), T-17 (Untch, 2009; Delamaro et al., 2014c; Deng et al., 2013; Kurtz et al., 2016; Derezińska, 2016), T-21 (Ammann et al., 2014; Kurtz et al., 2016; Just et al., 2017)
PG-2 – Automatically detecting equivalent mutants	T-2 (Kintis et al., 2015), T-3 (Offutt and Lee, 1994; Devroey et al., 2017), T-6 (Hierons et al., 1999; Papadakis and Le Traon, 2013; Schuler et al., 2009; Schuler and Zeller, 2013; Offutt et al., 2006; Kintis and Malevris, 2015; Papadakis et al., 2014; Offutt and Craft, 1994; Harman et al., 2000; Kintis et al., 2015), T-8 (Offutt and Craft, 1994; Kintis et al., 2017; Delgado-Pérez et al., 2017), T-11 (Offutt and Pan, 1997; Papadakis and Le Traon, 2013; Schuler et al., 2009; Schuler and Zeller, 2013; Offutt et al., 2006; Papadakis et al., 2014; Harman et al., 2000; Kintis et al., 2015; Patel and Hierons, 2016; Holling et al., 2016; Marcozzi et al., 2018; McMinn et al., 2018), T-12 (Delgado-Pérez et al., 2017), T-15 (Kintis and Malevris, 2013; Ferrari et al., 2013; Anbalagan and Xie, 2008; Vincenzi et al., 2002), T-16 (Delgado-Pérez et al., 2017), T-19 (Devroey et al., 2017)
PG-3 – Executing faster	T-2 (Devroey et al., 2016), T-3 (Offutt and Lee, 1994; Kim et al., 2012; Papadakis and Malevris, 2011b; Devroey et al., 2017), T-4 (Jackson and Woodward, 2000), T-5 (Choi and Mathur, 1993; Offutt et al., 1992; Reales and Polo, 2013; Wang et al., 2017; Krauser et al., 1991; Li et al., 2015; Jackson and Woodward, 2000), T-8 (Just et al., 2011; DeMillo et al., 1991; Denisov and Pankevich, 2018), T-10 (Weiss and Fleyshgakker, 1993; Ma et al., 2005; Kim et al., 2012; Untch et al., 1993; 1997; Wright et al., 2013; Papadakis and Malevris, 2011b; Gopinath et al., 2016b; Reales and Polo, 2014), T-11 (Zhang et al., 2012; Papadakis and Malevris, 2011b; Chen and Zhang, 2018), T-13 (Weiss and Fleyshgakker, 1993; Kim et al., 2012; Fleyshgakker and Weiss, 1994), T-15 (Durelli et al., 2012; Gligoric et al., 2012; Zhang et al., 2012; Just et al., 2014a; 2011; Wright et al., 2013; Zhang et al., 2016; Gopinath et al., 2016b; Devroey et al., 2016; McMinn et al., 2016; Alexander et al., 2002; Bogacki and Walter, 2006a; Chen and Zhang, 2018), T-16 (Papadakis and Malevris, 2011b), T-18 (Reales and Polo, 2014), T-19 (Devroey et al., 2016; 2017), T-20 (Gligoric et al., 2012; Just et al., 2014a; Gligoric et al., 2010; Wang et al., 2017)
PG-4 – Reducing the number of test cases or the number of executions	T-2 (Kaminski and Ammann, 2011; Harman et al., 2014; Devroey et al., 2016; Gopinath et al., 2018), T-3 (Fraser and Arcuri, 2015; Ma and Kim, 2016; Zhu et al., 2017; 2018), T-6 (Just et al., 2012b), T-7 (Tuya et al., 2007; Sahinoğlu and Spafford, 1990), T-10 (Ma and Kim, 2016), T-11 (Papadakis and Malevris, 2012; Kaminski and Ammann, 2011; Just et al., 2012b; Fraser and Arcuri, 2015), T-12 (Tuya et al., 2007), T-13 (Ma and Kim, 2016), T-15 (Zhang et al., 2013b; Derezińska, 2013; Devroey et al., 2016; Zhu et al., 2017; Derezińska and Hałas, 2015; Zhu et al., 2018), T-16 (Ayari et al., 2007; Oliveira et al., 2013; Adamopoulos et al., 2004; Fraser and Arcuri, 2015; Harman et al., 2014), T-19 (Devroey et al., 2016), T-3 (Marshall et al., 1990), T-6 (Marshall et al., 1990; Just et al., 2012b; Kintis and Malevris, 2014; Delgado-Pérez et al., 2017a), T-11 (Just et al., 2012b; Delgado-Pérez et al., 2017a), T-14 (Just and Schweiggert, 2015), T-15 (Steimann and Thies, 2010; Kaminski et al., 2013; Hu et al., 2011; Belli and Beyazit, 2015; Just and Schweiggert, 2015; Fernandes et al., 2017; Iida and Takada, 2017), T-16 (Domínguez-Jiménez et al., 2009a)
PG-5 – Avoiding the creation of certain mutants	T-3 (Papadakis and Malevris, 2011b; Zhang et al., 2010b), T-6 (Harman et al., 2000), T-10 (Papadakis and Malevris, 2010b; 2011b), T-11 (Offutt et al., 1999; Papadakis and Malevris, 2010b; Harman et al., 2000; Liu et al., 2006; Papadakis and Malevris, 2009; 2011b; DeMillo and Offutt, 1993; Zhang et al., 2010b; Bashir and Nadeem, 2017; DeMillo and Offutt, 1991), T-16 (Papadakis and Malevris, 2011a; Baudry et al., 2005; Papadakis and Malevris, 2011b; Bashir and Nadeem, 2017; Matnei Filho and Vergilio, 2015; Henard et al., 2014; Fraser and Zeller, 2012), T-18 (Fraser and Zeller, 2012), T-19 (Aichernig et al., 2013; 2014), T-20 (Bashir and Nadeem, 2017)
PG-6 – Automatically generating test cases	

**Fig. 10.** Distribution of studies that applied combined techniques per year.

**Table 5**

Mapping between cost reduction techniques and primary goals.

Cost reduction technique	Primary goals and respective studies
T-1 – Random mutation	PG-1 (Zhang et al., 2010a; 2013a; Wong and Mathur, 1995; Papadakis and Malevris, 2010a; Kurtz et al., 2016; Parsai et al., 2016; Gopinath et al., 2017; Derezińska and Rudnik, 2017; Petrović and Ivanković, 2018)
T-2 – Higher order mutation	PG-1 (Kintis et al., 2010; Kaminski et al., 2011b; Polo et al., 2009; Reales et al., 2013; Reuling et al., 2015; Papadakis and Malevris, 2010a; Omar and Ghosh, 2012; Madeyski et al., 2014; Harman et al., 2014; Lima et al., 2016; Derezińska, 2016; Abuljadayel and Wedyan, 2018), PG-2 (Kintis et al., 2015), PG-3 (Devroey et al., 2016), PG-4 (Kaminski and Ammann, 2011; Harman et al., 2014; Devroey et al., 2016; Gopinath et al., 2018)
T-3 – Weak mutation	PG-1 (Kintis et al., 2010; Zhu et al., 2017; 2018), PG-2 (Offutt and Lee, 1994; Devroey et al., 2017), PG-3 (Offutt and Lee, 1994; Kim et al., 2012; Papadakis and Malevris, 2011b; Devroey et al., 2017), PG-4 (Fraser and Arcuri, 2015; Ma and Kim, 2016; Zhu et al., 2017; 2018), PG-5 (Marshall et al., 1990), PG-6 (Papadakis and Malevris, 2011b; Zhang et al., 2010b)
T-4 – Firm mutation	PG-3 (Jackson and Woodward, 2000)
T-5 – Parallel execution	PG-3 (Choi and Mathur, 1993; Offutt et al., 1992; Reales and Polo, 2013; Wang et al., 2017; Krauser et al., 1991; Li et al., 2015; Jackson and Woodward, 2000)
T-6 – Data-flow analysis	PG-1 (Patrick et al., 2012; Al-Hajjaji et al., 2017), PG-2 (Hierons et al., 1999; Papadakis and Le Traon, 2013; Schuler et al., 2009; Schuler and Zeller, 2013; Offutt et al., 2006; Kintis and Malevris, 2015; Papadakis et al., 2014; Offutt and Craft, 1994; Harman et al., 2000; Kintis et al., 2015), PG-4 (Just et al., 2012b), PG-5 (Marshall et al., 1990; Just et al., 2012b; Kintis and Malevris, 2014; Delgado-Pérez et al., 2017a), PG-6 (Harman et al., 2000)
T-7 – Minimization and prioritization of test sets	PG-1 (Tuya et al., 2007), PG-4 (Tuya et al., 2007; Sahinoğlu and Spafford, 1990)
T-8 – Compiler optimization	PG-1 (Delgado-Pérez et al., 2017), PG-2 (Offutt and Craft, 1994; Kintis et al., 2017; Delgado-Pérez et al., 2017), PG-3 (Just et al., 2011; DeMillo et al., 1991; Denisov and Pankevich, 2018)
T-9 – Constrained mutation	PG-1 (Wong and Mathur, 1995; Wong et al., 1994; Gopinath et al., 2017; Petrović and Ivanković, 2018)
T-10 – Metamutants	PG-3 (Weiss and Fleyshgakker, 1993; Ma et al., 2005; Kim et al., 2012; Untch et al., 1993; 1997; Wright et al., 2013; Papadakis and Malevris, 2011b; Gopinath et al., 2016b; Reales and Polo, 2014), PG-4 (Ma and Kim, 2016), PG-6 (Papadakis and Malevris, 2010b; 2011b)
T-11 – Control-flow analysis	PG-1 (Sun et al., 2017b; Patrick et al., 2012; Ji et al., 2009; Just et al., 2017; Marcozzi et al., 2018; Petrović and Ivanković, 2018; McMinn et al., 2018), PG-2 (Offutt and Pan, 1997; Papadakis and Le Traon, 2013; Schuler et al., 2009; Schuler and Zeller, 2013; Offutt et al., 2006; Papadakis et al., 2014; Harman et al., 2000; Kintis et al., 2015; Patel and Hierons, 2016; Holling et al., 2016; Marcozzi et al., 2018; McMinn et al., 2018), PG-3 (Zhang et al., 2012; Papadakis and Malevris, 2011b; Chen and Zhang, 2018), PG-4 (Papadakis and Malevris, 2012; Kaminski and Ammann, 2011; Just et al., 2012b; Fraser and Arcuri, 2015), PG-5 (Just et al., 2012b; Delgado-Pérez et al., 2017a), PG-6 (Offutt et al., 1999; Papadakis and Malevris, 2010b; Harman et al., 2000; Liu et al., 2006; Papadakis and Malevris, 2009; 2011b; DeMillo and Offutt, 1993; Zhang et al., 2010b; Bashir and Nadeem, 2017; DeMillo and Offutt, 1991)
T-12 – Selective mutation	PG-1 (Zhang et al., 2010a; Offutt et al., 1993; Tuya et al., 2007; Mresa and Bottaci, 1999; Offutt et al., 1996; Gligoric et al., 2013; Bluemke and Kulesza, 2014a; Zhang et al., 2013a; Sridharan and Siami-Namin, 2010; Kurtz et al., 2016; Delgado-Pérez et al., 2017c; Gopinath et al., 2017; Phramontripong and Offutt, 2017; Derezińska and Rudnik, 2012; Bluemke and Kulesza, 2014b; Delgado-Pérez et al., 2017; Sun et al., 2017a), PG-2 (Delgado-Pérez et al., 2017), PG-4 (Tuya et al., 2007)
T-13 – Serial execution	PG-3 (Weiss and Fleyshgakker, 1993; Kim et al., 2012; Fleyshgakker and Weiss, 1994), PG-4 (Ma and Kim, 2016)
T-14 – Sufficient operators	PG-1 (Zhang et al., 2010a; Siami-Namin et al., 2008; Vincenzi et al., 2001; Barbosa et al., 2001; Delamaro et al., 2014b; Siami-Namin and Andrews, 2006; Lacerda and Ferrari, 2014), PG-5 (Just and Schweiggert, 2015)
T-15 – Optimization of generation, execution and analysis of mutants	PG-1 (Wedyan and Ghosh, 2012; Inozemtseva et al., 2013; Reuling et al., 2015; Cachia et al., 2013; Kaminski and Ammann, 2009; Zhu et al., 2017; Derezińska and Hałas, 2015; Zhu et al., 2018), PG-2 (Kintis and Malevris, 2013; Ferrari et al., 2013; Anbalagan and Xie, 2008; Vincenzi et al., 2002), PG-3 (Durelli et al., 2012; Gligoric et al., 2012; Zhang et al., 2012; Just et al., 2014a; 2011; Wright et al., 2013; Zhang et al., 2016; Gopinath et al., 2016b; Devroey et al., 2016; McMinn et al., 2016; Alexander et al., 2002; Bogacki and Walter, 2006a; Chen and Zhang, 2018), PG-4 (Zhang et al., 2013b; Derezińska, 2013; Devroey et al., 2016; Zhu et al., 2017; Derezińska and Hałas, 2015; Zhu et al., 2018), PG-5 (Steimann and Thies, 2010; Kaminski et al., 2013; Hu et al., 2011; Belli and Beyazit, 2015; Just and Schweiggert, 2015; Fernandes et al., 2017; Iida and Takada, 2017), PG-1 (Dominguez-Jiménez et al., 2011; Oliveira et al., 2013; Sridharan and Siami-Namin, 2010; Adamopoulos et al., 2004; Ji et al., 2009; Harman et al., 2014; Delgado-Pérez et al., 2017b; Quyen et al., 2016; Nobre et al., 2012; Delgado-Pérez et al., 2017; Abuljadayel and Wedyan, 2018), PG-2 (Delgado-Pérez et al., 2017), PG-3 (Papadakis and Malevris, 2011b), PG-4 (Ayari et al., 2007; Oliveira et al., 2013; Adamopoulos et al., 2004; Fraser and Arcuri, 2015; Harman et al., 2014), PG-5 (Dominguez-Jiménez et al., 2009a), PG-6 (Papadakis and Malevris, 2011a; Baudry et al., 2005; Papadakis and Malevris, 2011b; Bashir and Nadeem, 2017; Matnei Filho and Vergilio, 2015; Henard et al., 2014; Fraser and Zeller, 2012)
T-16 – Evolutionary algorithms	PG-1 (Untch, 2009; Delamaro et al., 2014c; 2014a; Deng et al., 2013; Kurtz et al., 2016; Derezińska, 2016)
T-17 – One-op mutation	PG-3 (Reales and Polo, 2014), PG-6 (Fraser and Zeller, 2012)
T-18 – Execution trace analysis	PG-2 (Devroey et al., 2017), PG-3 (Devroey et al., 2016; 2017), PG-4 (Devroey et al., 2016), PG-6 (Aichernig et al., 2013; 2014)
T-19 – Model-based mutation	PG-3 (Gligoric et al., 2012; Just et al., 2014a; Gligoric et al., 2010; Wang et al., 2017), PG-6 (Bashir and Nadeem, 2017)
T-20 – State-based analysis	PG-1 (Ammann et al., 2014; Kurtz et al., 2016; Just et al., 2017)
T-21 – Minimal mutation	

evaluate operator-based mutant selection. Phramontripong and Offutt (2017) applied a set of new mutation operators for Web applications and analyzed redundancy among operators. Delgado-Pérez et al. (2017) explored the Trivial Compiler Equivalence technique (Kintis et al., 2017; Papadakis et al., 2015) in combination with selective mutation and evolutionary algorithms, and Delgado-Pérez et al. (2017c) applied selective mutation to 83 classes written in C++. Sun et al. (2017a) investigated mutation testing for WS-

BPEL programs and used selective mutation to assess applicability and efficacy of the operators.

Wong et al. (1994) used constrained mutation and explored varied strategies to compose subsets of operators, and Wong and Mathur (1995) compared constrained mutation with random mutation. Recently, Petrović and Ivanković (2018) described an approach for selecting productive mutants. These are, to some extent, similar to constrained mutation as originally proposed. It relies on the

developers' understanding of the code for selecting either killable or equivalent mutants.

Papadakis and Malevris (2010a) empirically evaluated random mutation and higher order mutation. Bluemke and Kulesza (2013) explored random mutation with eight Java classes and used various tools (MuClipse, CodePro, and some tailor-made tools). Parsai et al. (2016) proposed to apply weights to randomly selected mutants and to focus on only acceptable mutants. Derezińska and Rudnik (2017) explored random mutation together with equivalence partitioning for object-oriented (C#) programs.

Barbosa et al. (2001) defined a procedure to identify sufficient operators. Vincenzi et al. (2001) applied the same procedure (Barbosa et al., 2001) to empirically evaluate unit-level and integration-level mutation operators for C programs. Siami-Namin and Andrews (2006) and Siami-Namin et al. (2008) explored sufficient operators for C programs supported by statistical methods (linear regression analysis). Delamaro et al. (2014b) presented a new procedure for defining sufficient operators. The process iteratively includes mutation operators until 100% mutation score is achieved. Lacerda and Ferrari (2014) applied Barbosa et al. (2001)'s method in AspectJ programs.

Adamopoulos et al. (2004) applied an evolutionary algorithm that relies on a fitness function to avoid equivalent mutants during the co-evolution process. Only mutants that are hard to kill and test cases that are good at detecting mutants are selected. Ji et al. (2009) applied an evolutionary algorithm to support a domain reduction technique to identify mutant clusters, and then executed one mutant per cluster. Domínguez-Jiménez et al. (2011) proposed Evolutionary Mutation Testing (EMT) to reduce the number of mutants. They used the GAmara tool to apply the approach to WS-BPEL projects. Delgado-Pérez et al. (2017b) also explored the EMT approach for C++ programs. Nobre et al. (2012) explored three multi-objective algorithms (namely, MTabu, NSGA-II, and PACO). Oliveira et al. (2013) proposed a new genetic co-evolutionary algorithm and compared it with five other evolutionary approaches. Harman et al. (2014) investigated strongly subsuming higher order mutants (SSHOMs) and evolutionary algorithms to reduce both the number of mutants and the number of test cases. Quyen et al. (2016) explored mutation testing of Simulink models with evolutionary algorithms. Abuljadayel and Wedyan (2018) applied genetic algorithms to kill higher order mutants for Java programs.

Polo et al. (2009) explored higher order mutation for Java programs. Kintis et al. (2010) experimented with higher order mutation and weak mutation, and compared results against strong mutation for Java programs. Kaminski et al. (2011b) suggested that generating only higher order logical mutants was sufficient. Omar and Ghosh (2012) applied higher order mutation to AspectJ programs. Reales et al. (2013) explored higher order mutation for Java programs by using the Bacterio tool and different algorithms. Madeyski et al. (2014) implemented four strategies for higher order mutation of Java programs with support of the JudyDiffOP tool. Reuling et al. (2015) proposed an approach to perform higher order mutation testing on feature models of software product lines (SPLs), trying to avoid redundant and equivalent mutants. Lima et al. (2016) experimented with four strategies for higher order mutation and compared results with first order mutation. Derezińska (2016) analyzed first order and higher order mutation based on statement deletion operators for C# programs, and presented a high-level analysis focused on the benefits of the deletion operators.

Sun et al. (2017b) applied control-flow analysis and proposed four strategies that focus on the most diverse mutants. Just et al. (2017) proposed an approach for handling mutants based on abstract syntax tree analysis and on the minimal mutation ap-

proach, originally proposed by Ammann et al. (2014). Minimal mutation was also explored by Kurtz et al. (2016), together with selective mutation, one-op mutation, and random mutation.

Marcozzi et al. (2018) applied control-flow analysis to reduce the number of equivalent mutants based on proofs of logical assertions. McMinn et al. (2018) analyzed paths in SQL expressions to remove ineffective mutants, including non-compilable, impaired, equivalent, and redundant mutants.

Kaminski and Ammann (2009) reduced the mutant set related to logical expressions by optimizing mutant generation based on minimal DNF predicates. Wedyan and Ghosh (2012) proposed an preliminary analysis of source code to prevent the generation of equivalent mutants for AspectJ programs. Zhu et al. (2017) explored weak mutation and Formal Concept Analysis (FCA) to group similar mutants and test cases. Inozemtseva et al. (2013) prioritized mutating parts of the programs that are most fault-prone based on fault history information; they used mutants generated by the PIT tool. Cachia et al. (2013) proposed incremental mutation testing, which tests newly modified code only. Derezińska and Hałas (2015) tried to optimize the generation and execution of mutants by applying syntax tree analysis and exploring runtime resources of the Python interpreter. Zhu et al. (2018) evolved their previous work (Zhu et al., 2017) and explored six compression techniques based on weak mutation to optimize strong mutation.

Untch (2009) proposed the one-op mutation technique to compose a reduced set of mutants. Deng et al. (2013) applied one-op mutation to Java classes. Delamaro et al. (2014c) designed new deletion-related mutation operators for C programs, and implemented them in the Proteum tool. In another study, Delamaro et al. (2014a) explored other mutation operators for one-op mutation.

Patrick et al. (2012) used control-flow analysis and data-flow analysis to identify mutants that had little impact on the program output, and hence are harder to kill. Al-Hajjaji et al. (2017) performed static data-flow analysis to select reduced sets of mutants for configurable systems.

### 5.5.2. Studies that tried to automatically detect equivalent mutants (PG-2)

Dervoey et al. (2017) detected equivalent mutants in finite automata, i.e. behavioral models, by applying weak mutation and language equivalence concepts.

Offutt and Craft (1994) performed data-flow analysis in mutants of Fortran programs and used six compiler optimization techniques to identify equivalent mutants. Kintis et al. (2017) and Papadakis et al. (2015) also explored compiler optimization to detect equivalent mutants. In particular, they proposed the TCE (Trivial Compiler Equivalence) technique to remove equivalent and duplicated mutants. Hierons et al. (1999) used program slicing to create simple mutants for C programs. Harman et al. (2000) explored program dependence (data-flow and control-flow analysis) to detect equivalent mutants by checking weakly killed mutants, and to generate test data automatically. Offutt et al. (2006) specified heuristics, based on data-flow analysis, to avoid equivalent mutants for class-level mutation operators. Schuler et al. (2009) applied the concept of dynamic invariants (derived via control-flow and data-flow analyses) and showed that invariant-violating mutants are less likely to be equivalent. More recently, Schuler and Zeller (2013) applied the same concepts to classify mutants as killable or equivalent, and to provide hints about the most relevant mutants to be analyzed. Papadakis and Le Traon (2013) and Papadakis et al. (2014) empirically investigated Schuler and Zeller (2013)'s approach as a way to reduce the effects of equivalent mutants in mutation testing. Kintis and Malevris (2015) investigated problematic data-flow patterns, through static analysis, to automatically identify equivalent mutants. Kintis et al. (2015)

relied on runtime information to detect equivalent mutants. They explored first and *high order mutation* and applied *control-flow* and *data-flow analyses* to classify killable and equivalent mutants.

Offutt and Pan (1997) devised a constraint-based technique to detect equivalent mutants. They applied *symbolic execution* and heuristics to recognize infeasible constraints among rules to generate test cases—if the constraints cannot be satisfied, the mutant cannot be killed and thus is equivalent. Patel and Hierons (2016) presented Interlocutory Mutation Testing (IMT) as a predicate-based *control-flow analysis* technique to automatically classify equivalent and killable mutants in programs that are non-deterministic and susceptible to coincidental correctness. Holling et al. (2016) used static analysis and symbolic execution to classify mutants as being killable, equivalent, and “don’t know.” More recently, Marcozzi et al. (2018) implemented a technique to identify equivalent mutants by proving the validity of logical assertions.

Vincenzi et al. (2002) optimized the prediction of equivalent mutants by applying Bayesian-learning. The outcomes depend on the size of randomly-generated test sets. Anbalagan and Xie (2008) and Ferrari et al. (2013) implemented an approach to automatically detect equivalent mutants of pointcut expressions in *AspectJ* programs. Kintis and Malevris (2013) introduced the concept of *mirrored mutants*, which are mutants that exhibit analogous behavior, such that identifying one mirrored mutant as being equivalent could help recognize other equivalent mutants.

#### 5.5.3. Studies that tried to execute faster (PG-3)

Devroey et al. (2016) explored *higher order mutation*, *model-based mutation*, and *optimization strategies* to reduce the number of mutants of formal transition systems.

Offutt and Lee (1994) explored *weak mutation* to speed up mutant execution and to automatically detect equivalent mutants. Papadakis and Malevris (2011b) applied *control-flow analysis* and *evolutionary algorithms* to support test case generation, and performed weak mutation to check the coverage of generated *metamutants*. Kim et al. (2012) proposed optimizing test executions by avoiding redundant executions that were identified using *weak mutation*. Devroey et al. (2017) detected equivalent mutants of finite automata (behavioral) models. They applied *weak mutation* together with language equivalence concepts, to detect equivalent mutants and speed up mutation analysis.

Krauser et al. (1991) unified mutants for *parallel execution* on Single Instruction Multiple Data (SIMD) machines. Offutt et al. (1992) tried similar ideas with Multiple Instruction Multiple Data (MIMD) machines. Similarly to Offutt et al. (1992), Choi and Mathur (1993) developed PMothra, which enabled *parallel execution* of mutants using a hypercube computer. Jackson and Woodward (2000) introduced the parallel *firm mutation* technique for Java programs, using Java threads to increase execution speed. Reales and Polo (2013) reported results from five algorithms for *parallel execution of mutants* and showed that the mutant execution cost is reduced proportionally to the number of nodes being used. Li et al. (2015) explored mutation testing for Ruby programs with *parallel execution* in a cloud infrastructure.

DeMillo et al. (1991) investigated a *compiler optimization* technique that created mutants by applying small patches at compile time. Just et al. (2011) explored conditional mutation to reduce the time to generate and execute the mutants. Denisov and Pankevich (2018) introduced a new tool for mutation testing based on the Low-Level Virtual Machine (LLVM), which compiles only the code fragments that are mutated.

Untch (1992) and Untch et al. (1993) introduced *metamutants* to compile and execute mutants faster. Untch et al. (1997) later implemented a metamutant-based tool (TUMs) using Mutant Schema Generation (MSG). Weiss and Fleyshgakker (1993) proposed a

new algorithm for *serial mutation of metamutants* for fast mutant execution. They later invented a new algorithm for *serial execution* and analysis of mutants called Lazy Mutation Analysis (LMA). Ma et al. (2005) applied the *metamutant* concept in the MuJava tool to generate mutants for compiled Java code. Wright et al. (2013) explored four different approaches that combine *metamutants* and *optimization techniques* to speed up mutation testing of database schemas. Reales and Polo (2014) proposed the MUSIC technique, which implements *metamutants* and *execution trace analysis* to speed up mutation. Gopinath et al. (2016b) explored the concepts of *metamutants* and *execution optimization techniques* to speed up mutant execution by avoiding multiple mutant compilation and redundant, partial test case executions.

Zhang et al. (2012) used *control-flow analysis* and *optimization techniques* to speed up mutation testing during regression testing. Chen and Zhang (2018) also applied regression testing ideas, reporting fewer test cases required, fewer mutants executed, and fewer test executions.

Alexander et al. (2002) used *optimization techniques* to mutate Java objects. The study describes a multi-threaded tool that sped up mutation. Bogacki and Walter (2006a) used aspect-oriented programming concepts to perform mutation while executing the original program to try to reduce mutant compilation and execution time. Durelli et al. (2012) modified the Java Virtual Machine to embed native support to speed up execution. Gligoric et al. (2010) presented a method to efficiently explore the states of multithreaded programs. The approach was implemented by Gligoric et al. (2012). Just et al. (2014a) implemented *optimizations* in the Major mutation tool, focusing on infected states and propagation to output. Wang et al. (2017) extended Just et al. (2014a)'s approach by using meta-functions to fork new processes to reduce redundant executions. Zhang et al. (2016) predicted whether mutants would be killed before execution by using a model that relies on features of mutants, tests, and coverage measures. McMinn et al. (2016) virtualized mutation testing of databases, thus reducing communication with the database.

#### 5.5.4. Studies that tried to reduce the number of test cases or the number of executions (PG-4)

Kaminski and Ammann (2011) introduced the minimal-MUMCUT logic-based testing criterion based on *higher order mutation*, *control-flow analysis*, and fault hierarchies. Harman et al. (2014) investigated strongly subsuming *higher order mutants* (SSHOMs) and *evolutionary algorithms* to reduce both the number of mutants and the number of test cases. Devroey et al. (2016) explored *higher order mutation*, *model-based mutation*, and *optimization strategies* to reduce the number of mutants in formal transition systems. Gopinath et al. (2018) reduced the number of test case executions by creating “supermutants,” that combine several first order mutants (i.e., *metamutants*).

Fraser and Arcuri (2015) investigated *weak mutation*, *evolutionary algorithms* and *control-flow analysis* for automatic generation of reduced test sets. Ma and Kim (2016) devised a mutant clustering approach to execute only one mutant from a group of mutants that are weakly killed by a given test case. The approach generates *metamutants* and *serial mutants* to speed up execution. Zhu et al. (2017) explored *weak mutation* (more precisely, state infection analysis) and Formal Concept Analysis (FCA) to group similar mutants and test cases. Zhu et al. (2018) evolved their previous work (Zhu et al., 2017) and explored six compression techniques to improve the efficiency of strong mutation based on *weak mutation*.

Tuya et al. (2007) proposed mutation operators for SQL. They used *selective mutation* and *minimization and prioritization of test sets* to reduce the number of mutants and required test cases. Sahinoğlu and Spaford (1990) addressed *minimization and prioritization of test sets* by applying a sequential statistical procedure

based on prespecified thresholds. Derezińska (2013) defined a mutant clustering approach for C# programs to reduce the numbers of mutants and tests.

Papadakis and Malevris (2012) used control-flow analysis to select test cases. Just et al. (2012b) also used control-flow analysis to prioritize test cases.

Zhang et al. (2013b) optimized mutant execution by predicting, without running, which mutants can be killed based on historical test coverage.

Adamopoulos et al. (2004) applied an evolutionary algorithm that used a fitness function to avoid equivalent mutants during co-evolution. Only mutants that were hard to kill and test cases that were good at detecting mutants were selected. Ayari et al. (2007) proposed an evolutionary approach based on ant colony optimization to generate tests automatically. Oliveira et al. (2013) explored genetic co-evolutionary algorithms to generate reduced sets of test cases that achieve higher mutation scores.

### 5.5.5. Studies that tried to avoid creation of certain mutants (PG-5)

Marshall et al. (1990)'s approach predicts the impact of mutations of program variables from the strong mutation and weak mutation perspectives, thus avoiding the creation of many mutants. Just et al. (2012b) used data-flow analysis to avoid creating redundant relational and conditional operator mutants. They also prioritized test cases based on control-flow analysis. Kintis and Malevris (2014) prevented the generation of equivalent mutants by defining data flow patterns that reveal code locations that should not be mutated. Delgado-Pérez et al. (2017a) defined class-level mutation operators for C++ programs and a set of restrictions (based on data-flow analysis and control-flow analysis) to avoid the creation of unproductive mutants (equivalent, duplicate, invalid, and trivial).

Just and Schweigert (2015) defined rules to avoid the creation of unnecessary operator mutants (conditional, unary, and relational) mutants by using sufficient operators, optimization, and prior results (Just et al., 2012b).

Steimann and Thies (2010) explored behavior-preserving constraints (based on refactoring rules) to optimize the generation of killable mutants. Hu et al. (2011) evaluated static and dynamic nature of class-level mutation operators. They proposed new rules to avoid the creation of equivalent mutants. Kaminski et al. (2013) theoretically analyzed and updated fault hierarchies for relational mutation operators (ROR), which eliminate redundancy among mutants. Their analysis showed that only three out of seven possible ROR mutants need to be generated. Inspired by Kaminski et al. (2013)'s, Iida and Takada (2017) introduced the notion of mutant killable preconditions to identify redundant mutants in control-flow statements. Belli and Beyazit (2015) proposed an event-based approach named *k*-Reg-based that avoids creating equivalent and redundant mutants for models written in a regular grammar. Fernandes et al. (2017) proposed a strategy to create rules to avoid useless (equivalent and duplicated) mutants. The rules discard those mutants right before they are generated.

Domínguez-Jiménez et al. (2009a) explored the use of genetic algorithms to avoid creating mutants for WS-BPEL compositions.

### 5.5.6. Studies that tried to automatically generate test cases (PG-6)

Zhang et al. (2010b) used control-flow analysis to introduce mutant-killing constraints into the program under test, and explored weak mutation concepts to guide the generation of test inputs. Papadakis and Malevris (2011b) used control-flow analysis and evolutionary algorithms to generate test cases, and performed weak mutation to evaluate coverage on the generated metamutants.

Harman et al. (2000) explored program dependence (data-flow analysis and control-flow analysis) to detect equivalent mutants by checking weakly killed mutants, as well as to automatically generate test cases.

Papadakis and Malevris (2010b) performed control-flow analysis of metamutants to generate better test cases.

DeMillo and Offutt (1991, 1993) used symbolic execution (control-flow analysis) to develop Constraint-Based Testing (CBT) to automatically generate test cases. CBT generated tests that satisfy the conditions needed to reach mutants and infect the program state after the mutant, and discarded tests that did not add to the mutation score. Later, Offutt et al. (1999) evolved the CBT technique by replacing symbolic execution algorithms with dynamic symbolic execution in a technique they called dynamic domain reduction (DDR). Liu et al. (2006) explored path-wise test data generation, claiming it was more efficient than symbolic evaluation. Papadakis and Malevris (2009) applied path selection based on control-flow analysis to generate tests.

Bashir and Nadeem (2017) proposed a combined fitness function for evolutionary algorithms that uses both control-flow and state-based properties of the mutants to generate test data automatically. Baudry et al. (2005) explored the concept of ant colonies in evolutionary algorithms to generate test cases. Papadakis and Malevris (2011a) used dynamic information from mutant executions to generate test data automatically. Fraser and Zeller (2012) presented *uTest*, which generates unit tests for Java classes based on genetic algorithms and execution trace analysis. Henard et al. (2014) explored search-based optimization methods to minimize the number of selected configurations in Software Product Line testing and to maximize the number of mutants killed. Matnei Filho and Vergilio (2015) also tested Software Product Lines using a multi-objective evolutionary approach.

Aichernig et al. (2013) combined model-based testing and mutation testing to automatically generate test cases. Their ideas can be used for several languages (e.g. Prolog, Java and C). They also developed a tool named MoMut (Aichernig et al., 2014) to generate test cases from UML models.

**Answer for RQ1:** Regarding RQ1 (Which techniques support cost reduction of mutation testing?), we noticed that mutation-related costs can be reduced by applying a wide range of techniques, sometimes individually and sometimes in combination. Some of the most common techniques are traditional software analysis methods such as control-flow analysis and data-flow analysis. Others, such as selective mutation and higher-order mutation, are common only within the mutation testing field, whereas still others are widely used in CS and Math (e.g. compiler optimization, evolutionary algorithms, and optimization-related techniques). These observations lead us to conclude that mutation-related cost reduction research is more interdisciplinary than in the past. That implies that collaboration with researchers from other areas can be very productive and should be encouraged.

## 6. Cost reduction metrics and results

Another contribution of this paper is the analysis of the benefits (in terms of cost savings) versus the loss of effectiveness (in terms of the mutation score or test quality). The analysis concerns research questions RQ2 and RQ3 (defined in Section 3.1) and starts with the characterization of metrics that have been used in the selected studies (Section 6.1). Then, Section 6.2 summarizes the main findings regarding results collected with those metrics.

### 6.1. Metrics and their timeline

Table 6 lists metrics and studies that applied those metrics. The third column of the table, “Intent,” shows the intended goal

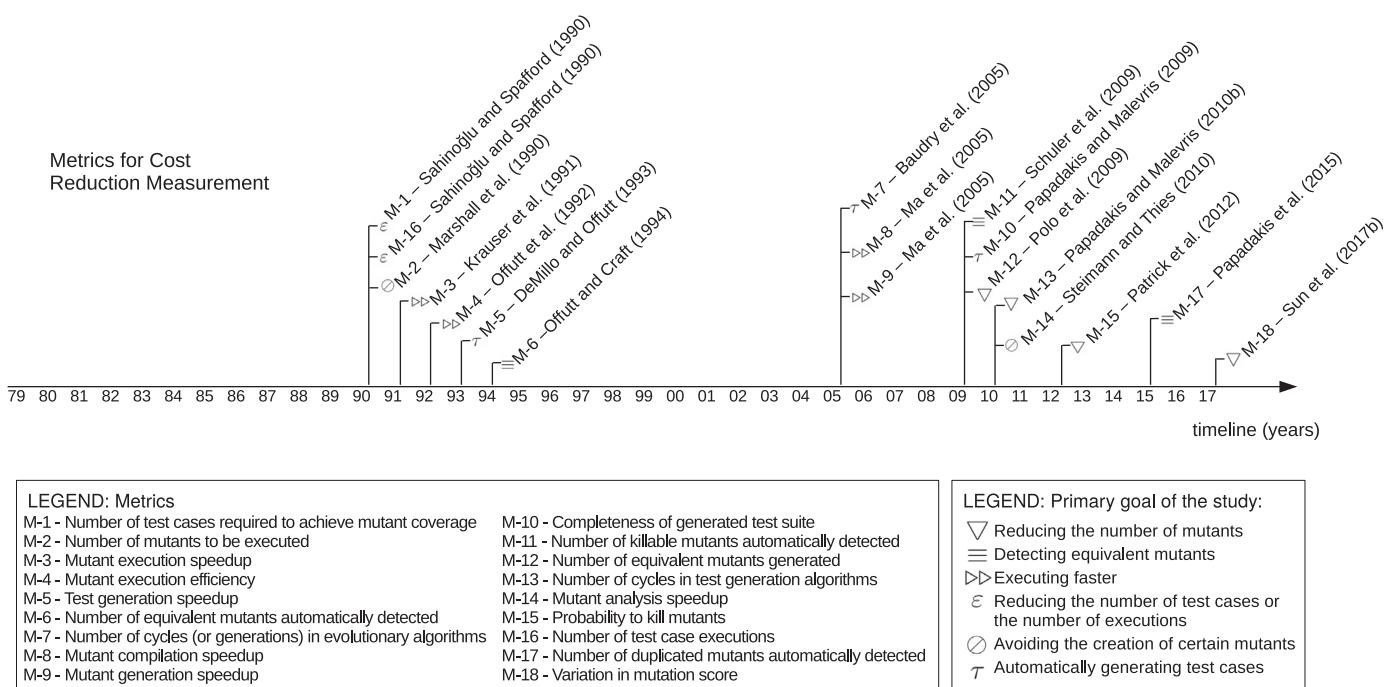
**Table 6**

Metrics to measure the cost reduction of mutation testing.

ID	Metric name	Intent	Studies that collected the metric
M-1	<i>Number of test cases required to achieve mutant coverage</i>	Estimate the number of test cases that need to be created (automatically or manually) or executed.	Wong and Mathur (1995), Liu et al. (2006), Tuya et al. (2007), Papadakis and Malevris (2010a), Kaminski and Ammann (2011), Papadakis and Malevris (2012), Derezińska (2013), Harman et al. (2014), Delamaro et al. (2014c,b), Bluemke and Kulesza (2014a), Lima et al. (2016), Reales et al. (2013), Delgado-Pérez et al. (2017c), Chen and Zhang (2018), Delamaro et al. (2014a), Zhu et al. (2017), Wong et al. (1995), Ji et al. (2009), Cachia et al. (2013), Derezińska and Rudnik (2017), Derezińska and Rudnik (2012), Sahinoglu and Spafford (1990), Oliveira et al. (2013), Wong et al. (1994), Offutt et al. (1993), Wong and Mathur (1995), Offutt et al. (1996), Barbosa et al. (2001), Vincenzi et al. (2001), Tuya et al. (2007), Siami-Namin et al. (2008), Kaminski and Ammann (2009), Polo et al. (2009), Untch (2009), Sridharan and Siami-Namin (2010), Papadakis and Malevris (2010a), Steimann and Thies (2010), Kaminski et al. (2011b), Domínguez-Jiménez et al. (2011), Omar and Ghosh (2012), Just et al. (2012b), Nobre et al. (2012), Bluemke and Kulesza (2014a), Derezińska (2013), Deng et al. (2013), Reales et al. (2013), Zhang et al. (2013a), Gligoric et al. (2013), Just and Schweiggert (2015), Madeyski et al. (2014), Harman et al. (2014), Lacerda and Ferrari (2014), Ammann et al. (2014), Delamaro et al. (2014a, 2014c, 2014b), Belli and Beyazit (2015), Lima et al. (2016), Delgado-Pérez et al. (2017a), Delgado-Pérez et al. (2017c), Wang et al. (2017), Zhu et al. (2017), Gopinath et al. (2017), Parsai et al. (2016), Praphamonpong and Offutt (2017), Sun et al. (2017a), Derezińska and Halas (2015), Bluemke and Kulesza (2014b), Chen and Zhang (2018), Zhu et al. (2018), Bluemke and Kulesza (2013), Marshall et al. (1990), Wong et al. (1995), Siami-Namin and Andrews (2006), Ji et al. (2009), Zhang et al. (2010a), Cachia et al. (2013), Inozemtseva et al. (2013), Reuling et al. (2015), Al-Hajjaji et al. (2017), Just et al. (2017), Derezińska and Rudnik (2017), Derezińska (2016), Derezińska and Rudnik (2012), Oliveira et al. (2013), Henard et al. (2014), Marcozzi et al. (2018), Kaminski and Ammann (2009), Quyen et al. (2016), Kaminski et al. (2013), Iida and Takada (2017)
M-2	<i>Number of mutants to be executed</i>	Estimate the number of mutants that need to be executed.	Offutt et al. (1993), Wong and Mathur (1995), Untch et al. (1996), Barbosa et al. (2001), Vincenzi et al. (2001), Tuya et al. (2007), Siami-Namin et al. (2008), Kaminski and Ammann (2009), Polo et al. (2009), Untch (2009), Sridharan and Siami-Namin (2010), Papadakis and Malevris (2010a), Steimann and Thies (2010), Kaminski et al. (2011b), Domínguez-Jiménez et al. (2011), Omar and Ghosh (2012), Just et al. (2012b), Nobre et al. (2012), Bluemke and Kulesza (2014a), Derezińska (2013), Deng et al. (2013), Reales et al. (2013), Zhang et al. (2013a), Gligoric et al. (2013), Just and Schweiggert (2015), Madeyski et al. (2014), Harman et al. (2014), Lacerda and Ferrari (2014), Ammann et al. (2014), Delamaro et al. (2014a, 2014c, 2014b), Belli and Beyazit (2015), Lima et al. (2016), Delgado-Pérez et al. (2017a), Delgado-Pérez et al. (2017c), Wang et al. (2017), Zhu et al. (2017), Gopinath et al. (2017), Parsai et al. (2016), Praphamonpong and Offutt (2017), Sun et al. (2017a), Derezińska and Halas (2015), Bluemke and Kulesza (2014b), Chen and Zhang (2018), Zhu et al. (2018), Bluemke and Kulesza (2013), Marshall et al. (1990), Wong et al. (1995), Siami-Namin and Andrews (2006), Ji et al. (2009), Zhang et al. (2010a), Cachia et al. (2013), Inozemtseva et al. (2013), Reuling et al. (2015), Al-Hajjaji et al. (2017), Just et al. (2017), Derezińska and Rudnik (2017), Derezińska (2016), Derezińska and Rudnik (2012), Oliveira et al. (2013), Henard et al. (2014), Marcozzi et al. (2018), Kaminski and Ammann (2009), Quyen et al. (2016), Kaminski et al. (2013), Iida and Takada (2017)
M-3	<i>Mutant execution speedup</i>	Estimate how much time is required for mutant execution.	Offutt et al. (1992), Choi and Mathur (1993), Offutt and Lee (1994), Untch et al. (1997), Ma et al. (2005), Gligoric et al. (2010), Papadakis and Malevris (2011b), Just et al. (2012b), Gligoric et al. (2012), Durelli et al. (2012), Reales and Polo (2014), Just and Schweiggert (2015), Just et al. (2014a), Delgado-Pérez et al. (2017a), Wang et al. (2017), Zhu et al. (2017), Zhang et al. (2016), Devroey et al. (2016), Li et al. (2015), Derezińska and Halas (2015), Krauser et al. (1991), McMinn et al. (2016), Kim et al. (2012), Zhang et al. (2013a), Madeyski et al. (2014), Chen and Zhang (2018), Zhu et al. (2018), Untch et al. (1993), Reales and Polo (2013), Wright et al. (2013), Ma and Kim (2016), Gopinath et al. (2016b), Derezińska and Rudnik (2017), Weiss and Fleysgakker (1993), Fleysgakker and Weiss (1994), Bogacki and Walter (2006a), Offutt et al. (1992)
M-4	<i>Mutant execution efficiency</i>	Estimate the performance gains or losses when two or more execution configurations (e.g. processor and network infrastructure) are compared.	DeMillo and Offutt (1993), Liu et al. (2006)
M-5	<i>Test generation speedup</i>	Estimate how much time is required for test generation.	Offutt and Craft (1994), Offutt and Pan (1997), Anbalagan and Xie (2008), Ferrari et al. (2013), Kintis and Malevris (2013), Wright et al. (2014), Papadakis et al. (2014), Kintis and Malevris (2015), Fernandes et al. (2017), Kintis et al. (2017), Devroey et al. (2017), Offutt et al. (2006), Schuler et al. (2009), Papadakis and Le Traon (2013), Patel and Hierons (2016), Marcozzi et al. (2018), McMinn et al. (2018), Bashir and Nadeem (2017), Baudry et al. (2005)
M-6	<i>Number of equivalent mutants automatically detected</i>	Estimate the number of mutants that need be handled during mutant analysis.	Delgado-Pérez et al. (2017a), Fernandes et al. (2017), Just et al. (2011), Ma et al. (2005)
M-7	<i>Number of cycles (or generations) in evolutionary algorithms</i>	Estimate the effort required to run a specific evolutionary algorithm.	Ma et al. (2005), Domínguez-Jiménez et al. (2011), Kim et al. (2012), Just et al. (2011)
M-8	<i>Mutant compilation speedup</i>	Estimate how much time is required for mutant compilation.	Papadakis and Malevris (2009), Kaminski and Ammann (2009), Papadakis and Malevris (2012), Matnei Filho and Vergilio (2015), Zhang et al. (2010b)
M-9	<i>Mutant generation speedup</i>	Estimate how much time is required for mutant generation.	Devroey et al. (2017), Holling et al. (2016), Schuler et al. (2009), Patel and Hierons (2016)
M-10	<i>Completeness of generated test suite</i>	Estimate the effort required to evolve the current set suite to achieve the intended mutant coverage.	(continued on next page)
M-11	<i>Number of killable mutants automatically detected</i>	Estimate the number of mutants that need be handled during mutant analysis.	

**Table 6** (continued)

ID	Metric name	Intent	Studies that collected the metric
M-12	<i>Number of equivalent mutants generated</i>	Estimate the number of mutants that need be handled during mutant analysis.	Papadakis and Malevris (2010a), Kintis et al. (2010), Hu et al. (2011), Kaminski et al. (2011b), Wedyan and Ghosh (2012), Deng et al. (2013), Kintis and Malevris (2014), Delamaro et al. (2014b), Polo et al. (2009), Delamaro et al. (2014a), Delamaro et al. (2014c), Madeyski et al. (2014) Papadakis and Malevris (2010b)
M-13	<i>Number of cycles in test generation algorithms</i>	Estimate the effort required to run a specific test generation algorithm.	Steimann and Thies (2010), Wright et al. (2014), Madeyski et al. (2014), Devroey et al. (2017), Aichernig et al. (2013)
M-14	<i>Mutant analysis speedup</i>	Estimate how much time is required for mutant analysis.	Patrick et al. (2012)
M-15	<i>Probability to kill mutants</i>	To obtain a probability for mutant coverage based on the available test set.	Zhang et al. (2012, 2013b), Bashir and Nadeem (2017), Derezińska and Halas (2015), Kim et al. (2012), Sahinoğlu and Spafford (1990), Ma and Kim (2016), Reales and Polo (2014), Gopinath et al. (2018), Chen and Zhang (2018)
M-16	<i>Number of test case executions</i>	Estimate the number of test case executions for a given mutant set.	Kintis et al. (2017), Fernandes et al. (2017), McMinn et al. (2018)
M-17	<i>Number of duplicated mutants automatically detected</i>	Estimate the number of mutants that need be handled during mutant analysis.	Sun et al. (2017b)
M-18	<i>Variation in mutation score</i>	Estimate gains and losses with respect to mutation score.	

**Fig. 11.** Timeline for metrics used to measure cost reduction of mutation testing according to their uses in peer-reviewed studies.

of the metrics. Notice that the intents of some metrics overlap. Metrics M-6 (*Number of equivalent mutants automatically detected*), M-11 (*Number of killable mutants automatically detected*), M-12 (*Number of equivalent mutants generated*), and M-17 (*Number of duplicated mutants automatically detected*) have similar intents, despite the fact that they collect different values from the sets of mutants. Also notice that the definition for each metric can be derived from its names and intent. For instance, M-3 (*Mutant execution speedup*) can be defined as the time spent (in a given time unit) for mutant execution when compared to conventional execution (i.e. without applying a cost reduction technique).

Fig. 11 lists all identified metrics in a timeline. The first studies that used the metrics were published in 1990 (Marshall et al., 1990; Sahinoğlu and Spafford, 1990) and 1991 (Krauser et al., 1991). Marshall et al. (1990) estimated how many mutants would be executed (M-2) after using static data-flow analysis to avoid

creating certain mutants. Sahinoğlu and Spafford (1990) applied statistical methods to estimate the number of test cases required to reach a threshold mutation score (M-1), as well as to estimate the number of test case executions (M-16). Both studies used small programs. Krauser et al. (1991) ran simulation experiments to measure mutant execution speedup (M-3) on parallel machines. The timeline also shows that (1) few metrics (M-1 to M-6, and M-16) were used in the early days (1990 to 1994), and (2) many more metrics were used in the subsequent 14 years.

Fig. 12 shows the number of studies that used each metric. *Number of mutants to be executed*(M-2) was used the most (66 studies), followed by *Mutant execution speedup*(M-3) (36 studies) and *Number of test cases required to achieve mutant coverage*(M-1) (25 studies). If we consider only recent research, particularly studies published in the last five years (2014–2018), the most applied

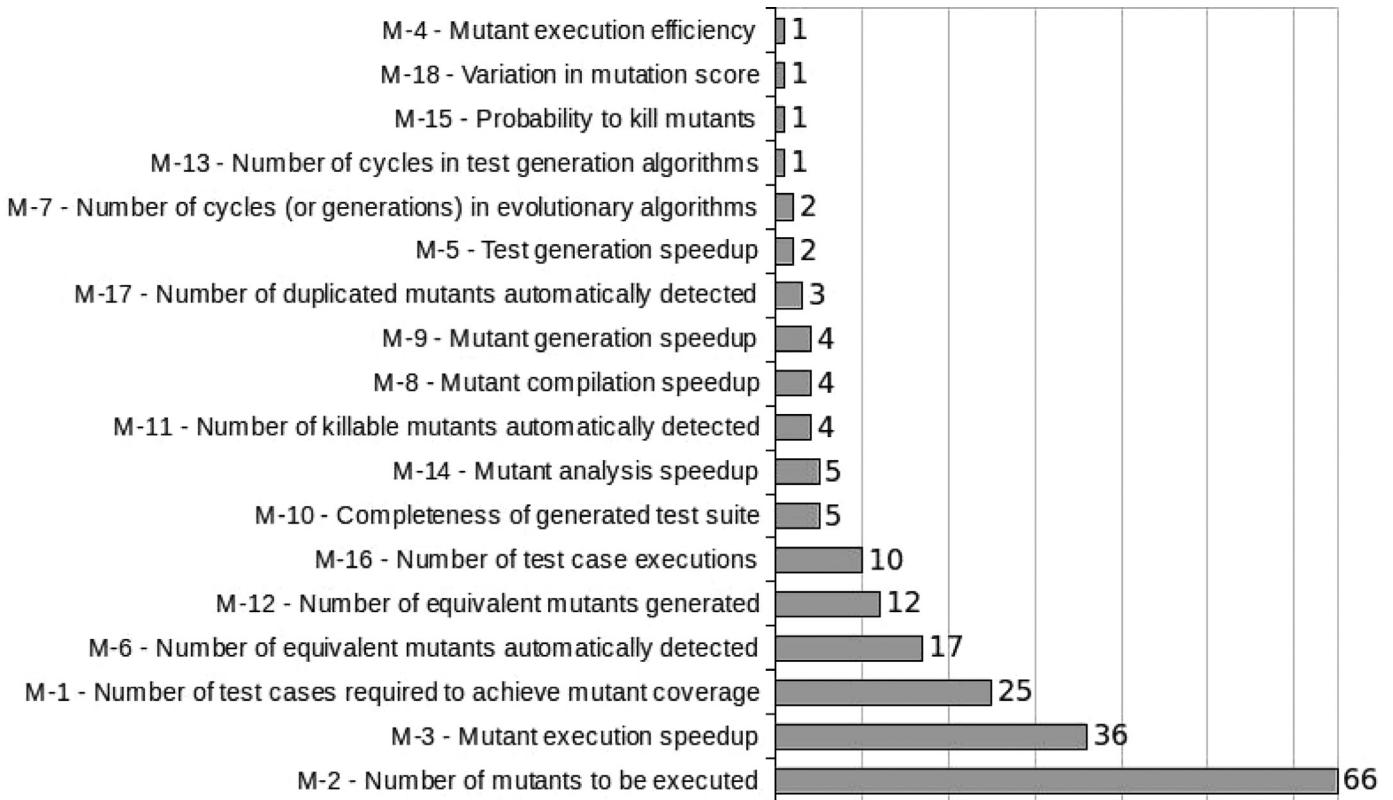


Fig. 12. Number of studies per metric.

metrics are the same as in Fig. 12. The distribution of studies per year is depicted in Fig. 13.

Regarding RQ2 (*Which metrics are used to measure the cost reduction of mutation testing?*), we noticed that a variety of metrics have been used, with emphasis in measuring the *Number of mutants to be executed*(M-2), *Mutant execution speedup*(M-3), and the *Number of test cases required to achieve mutant coverage*(M-1). These are important factors related to the cost of mutation testing. Nonetheless, the results of our SLR for RQ2 also indicates that a major challenge for the practical adoption of mutation testing, *dealing with equivalent mutants*, has not been addressed as much as the others and we recommend more research into the problem. Specifically, only 11% (17 of 151) of the studies used metric M-6 (*Number of equivalent mutants automatically detected*) and 8% (12 of 151) used M-12 (*Number of equivalent mutants generated*).

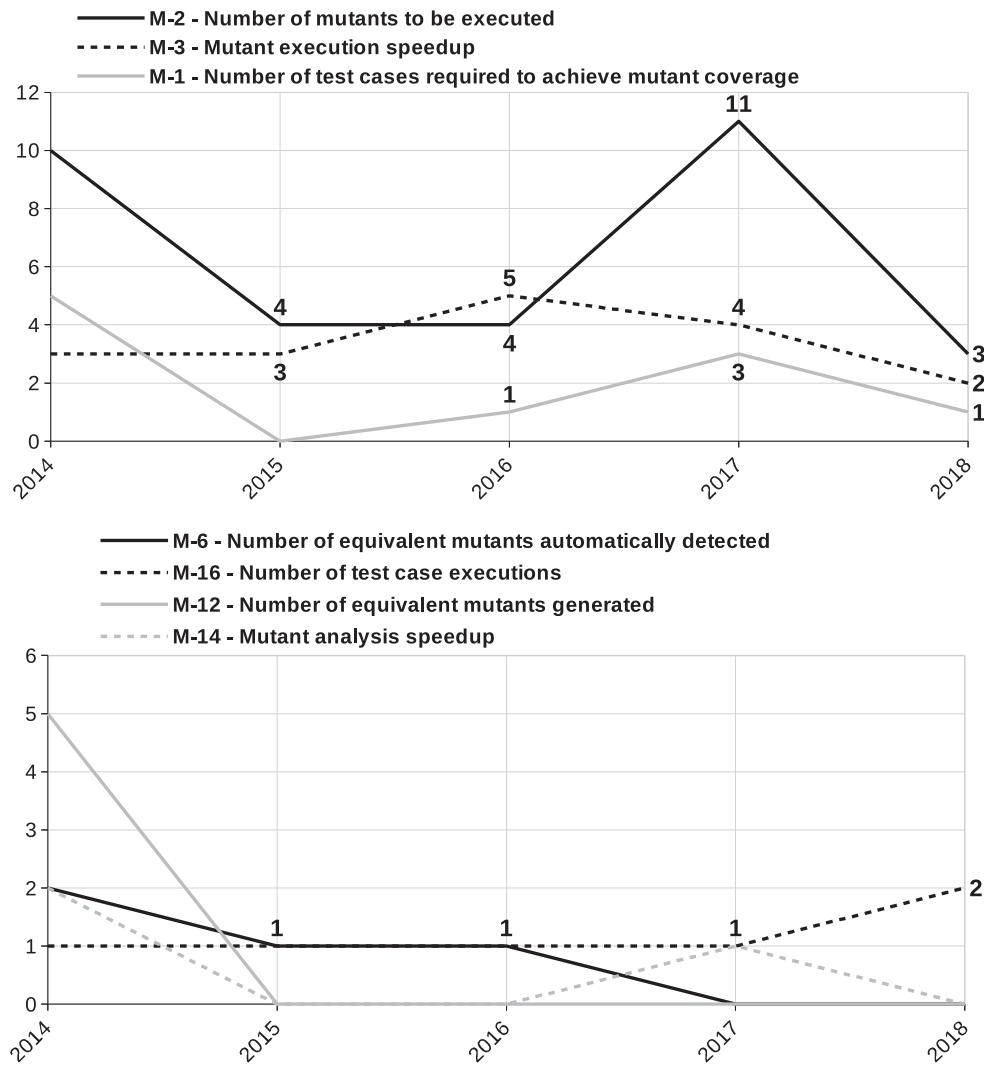
*Note about complex metrics not listed in this section:* The list of metrics identified in this SLR is not intended to be complete, because several studies applied customized, and somewhat complex, metrics in the experiments. Examples can be found in the studies by [Mresa and Bottaci \(1999\)](#), [Delamaro et al. \(2014a\)](#), and [Delgado-Pérez et al. \(2017c\)](#), and probably more. [Mresa and Bottaci \(1999\)](#) estimated the cost of test case generation to be a combination of factors such as the number of mutant executions and the number of redundant tests. They also defined metrics such as *relative test generation cost* and *relative equivalence detection*. [Delamaro et al. \(2014a\)](#) devised a weighted cost function that is based on the *relative cost* of a mutation operator instead of being based on its absolute cost. [Delgado-Pérez et al. \(2017c\)](#) measured the *degree of redundancy* of a mutation operator, which is defined as the ratio of the number of redundant mutants generated by an operator to the number of mutants generated by the remaining operators. In our analysis, we focused only on iden-

tifying metrics that are derived from quantifiable elements directly extracted from the mutant creation, execution and analysis steps. Typical examples are the time saved to run a set of mutants (reflected by the *mutant execution speedup* (M-3) metric), and the size of mutation-adequate test sets (reflected by the *number of test cases required to achieve mutant coverage* (M-1) metric). The measurements used most often are presented in the next section.

## 6.2. Results about cost reduction measurement

Figs. 14 – 17 show the results achieved by different studies with respect to different cost reduction measurements. These data reflects the first part of RQ3: “*What are the savings (benefit) ... for the techniques?*”. Results for the four most used metrics are displayed (M-2, M-3, M-1, and M-6), grouped by the three most common programming languages (Fortran, Java, and C). Figs. 18 and 19 show the mutation scores achieved in the same studies listed in Figs. 14 and 16 (i.e. for metrics M-2 and M-1). These data address the second part of RQ3: “*What are ... loss of effectiveness (as proxied by mutation score) for the techniques?*”. We do not present mutation score data for metrics M-3 and M-6, since they were not reported for techniques applied to speed up mutant execution and automatically detect equivalent mutants. Furthermore, some values are not available so are missing in Figs. 18 and 19.

In Figs. 14 through 19, the bars show the minimum, mean, and maximum percentages of either cost reduction or mutation score. Note that not all values could be extracted from the studies. From some studies, we extracted only minimum and maximum values (for example, [Wong and Mathur \(1995\)](#) in Fig. 14a). Sometimes we could extract only the mean values (for example, [Vincenzi et al. \(2001\)](#) in Fig. 14b). Furthermore, no values were



**Fig. 13.** Distribution of studies that applied the top seven metrics in the last five years.

extracted at all for some studies. This is because of the level of details of reported results, characteristics of the applied cost reduction techniques, and customized experimental settings. This also meant we could not create whisker charts, which would require more detailed data than were presented in several studies. Also note that some studies appear twice in the charts of Figs. 14–19. Authors of such studies applied more than one cost reduction technique, and performed individual measurements regarding cost reduction and mutation score. An example is by Wong and Mathur (1995) (Figs. 14a and 18a), who applied *random mutation* and *constrained mutation*. Additionally, Papadakis and Malevris (2010a) (Figs. 14b and 18b) explored *random mutation* and *higher order mutation*, and Kintis et al. (2010) (Fig. 17c) applied *higher order mutation* and *weak mutation*. Section 6.4 further discusses some characteristics of the measurements we retrieved.

Fig. 14 shows the results of studies that measured cost reduction using metric M-2 (*Number of mutants to be executed*) for Fortran (14a), C (14b), and Java (14c) programs. Fig. 15 shows the results regarding M-3 (*Mutant execution speedup*), while Figs. 16 and 17 show results for metrics M-1 (*Number of test cases required to achieve mutant coverage*) and M-6 (*Number of equivalent mutants*

*automatically detected*). Notice that M-1 was applied in a single study that targeted *Fortran* programs, therefore the chart is not shown in Fig. 16.

### 6.3. Analysis and discussion

Figs. 14 through 17 reveal wide variations in cost reductions from both the inter-study or the intra-study perspectives. Variations among mutation scores (Figs. 18 and 19) were less frequent, but still noticeable. Note that the inter-study perspective can provide an interesting baseline to compare results reported on groups of studies by different scientists who used the same cost reduction techniques. The intra-study perspective, on the other hand, provides a fine-grained notion of how varied results can be for different subjects of a particular study.

Regarding the *Number of mutants to be executed*, average reductions ranged from 26.70% (Just et al., 2012b) (Fig. 14c) to 99.77% (Steimann and Thies, 2010) (Fig. 14c). With respect to *Mutant execution speedup*, the average reductions ranged from 25.00% (Just and Schweiggert, 2015) (Fig. 15c) to 97.78% (Choi and Mathur, 1993) (Fig. 15a). For *Number of test cases required to achieve mutant coverage*, the average reductions ranged from 1.00% (Just et al., 2012b) (Fig. 16c) to 99.77% (Steimann and Thies, 2010) (Fig. 16c). Finally, regarding *Number of equivalent mutants*, the average reductions ranged from 1.00% (Just et al., 2012b) (Fig. 17c) to 99.77% (Steimann and Thies, 2010) (Fig. 17c).

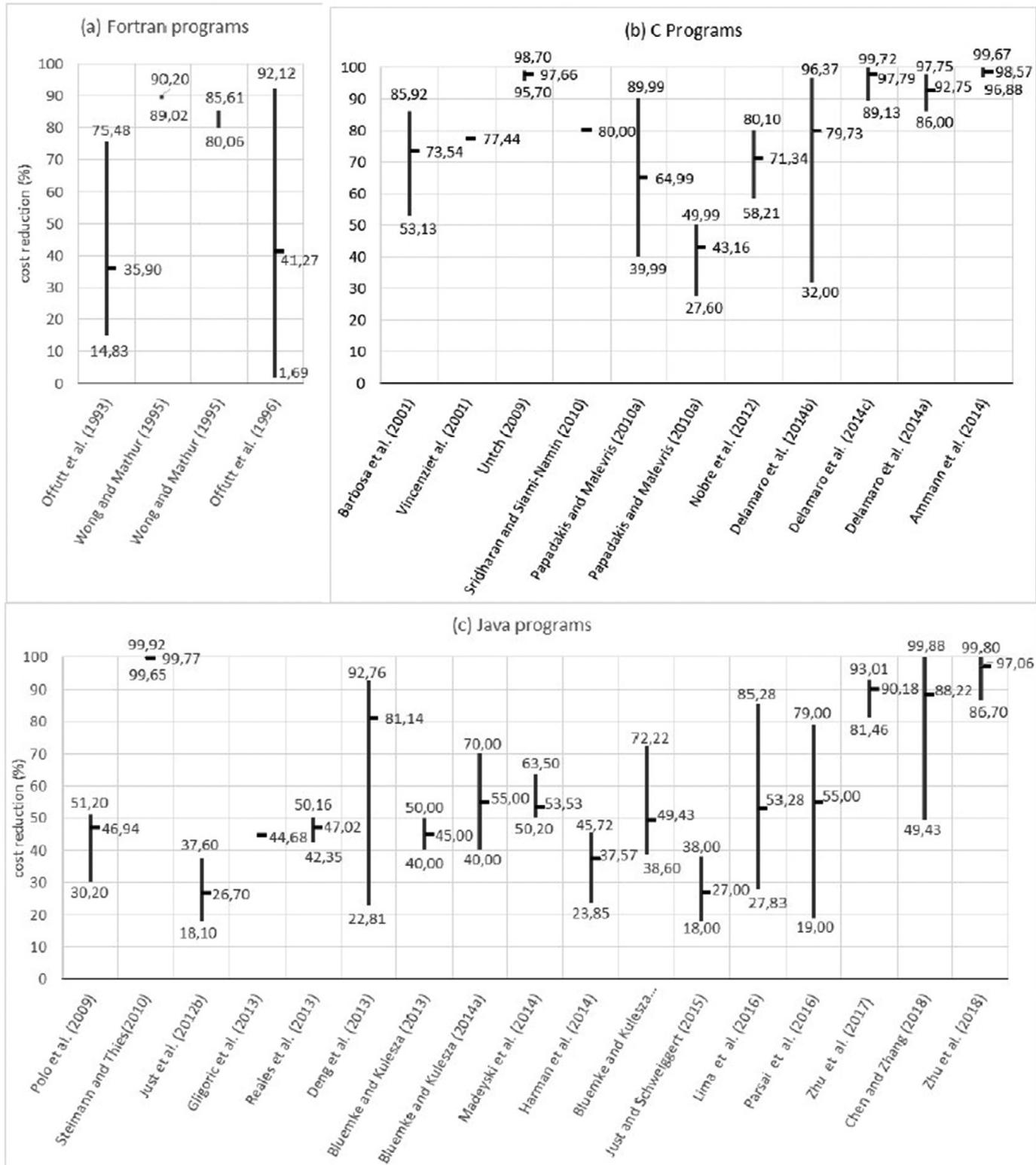
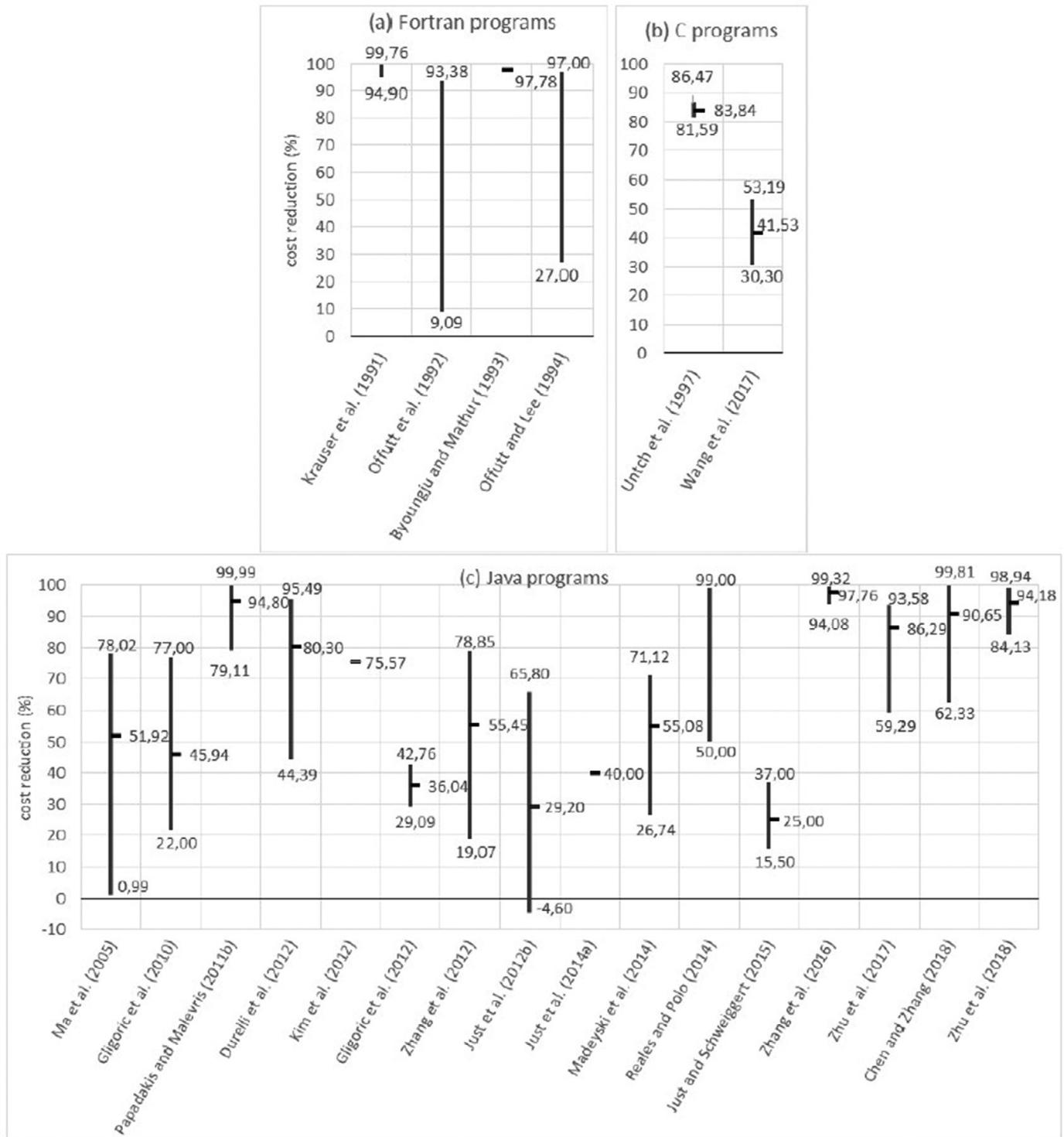


Fig. 14. Cost reduction regarding M-2 – Number of mutants to be executed.

mutant coverage and Number of equivalent mutants automatically detected, the averages ranged from 7.80% (Delamaro et al., 2014b) (Fig. 16a) to 92.60% (Siami-Namin et al., 2008) (Fig. 16a), and from 1.78% (Fernandes et al., 2017) (Fig. 17c) to 90.51% (Deng et al., 2013) (Fig. 17c).

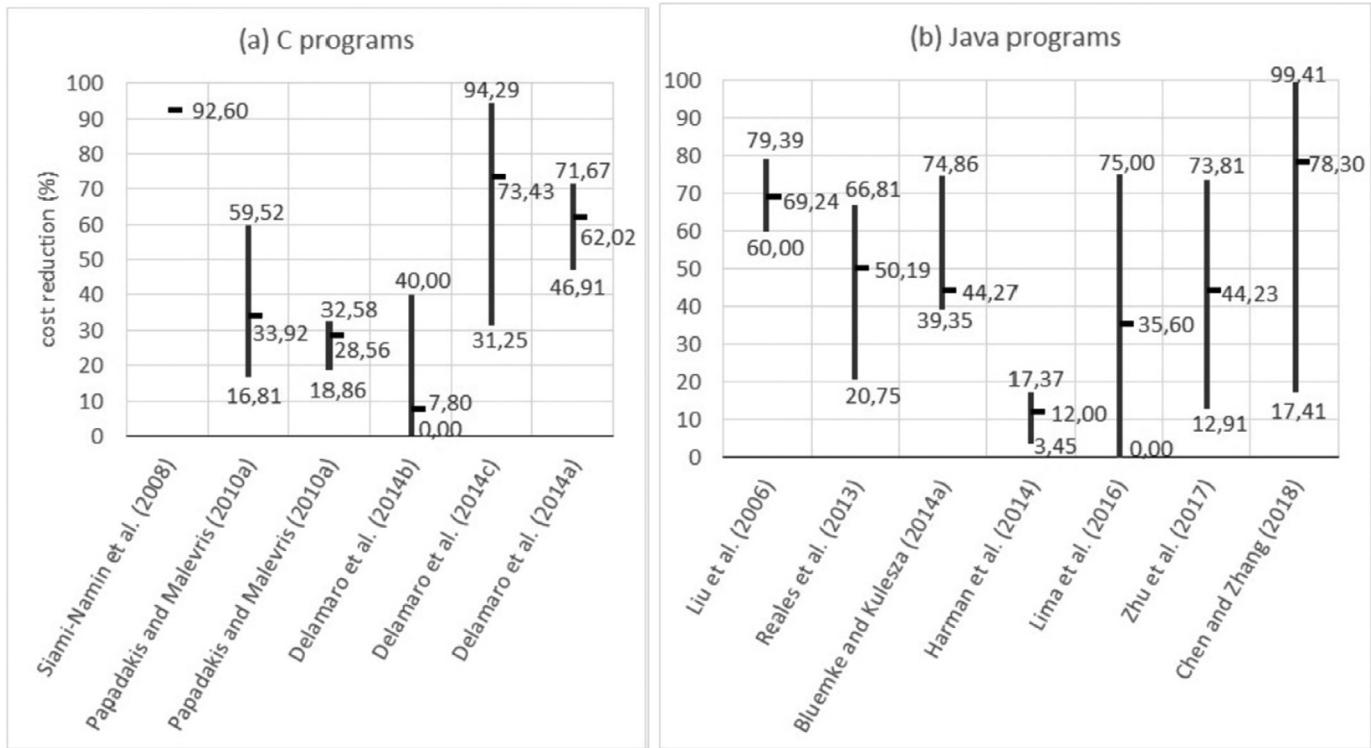
Regarding mutation scores, for both groups of studies (namely, the ones that measured Number of mutants to be executed, and Number of test cases required to achieve mutant coverage), the averages ranged from 0.800 (Reales et al., 2013) (Fig. 18c) to 1.000 (Delamaro et al., 2014b) (Fig. 18b). Before choosing a cost

**Fig. 15.** Cost reduction regarding M-3 – Mutant execution speedup.

reduction approach to adopt, it is important to consider the quality of the resulting test set regarding mutation-effectiveness, and hence its fault-revealing capability. Overall, the results presented in the charts show that while some studies achieved very high, even full, mutation score, some others did not satisfactorily maintain test effectiveness.

Taking an inter-study perspective reveals that the cost reduction variations presented in Figs. 14 through 17 are mainly due

to the range of techniques that have been used, or to the experimental settings of studies that used the same technique and tools to a common set of programs. For example, the cost reductions shown in Fig. 14a were obtained by applying *selective mutation* (Offutt et al., 1993; 1996), *random mutation*, and *constrained mutation* (Wong and Mathur, 1995). These three studies targeted the same language (*Fortran*) with the same tool (*Mothra*). However, Offutt et al. (1993, 1996)'s studies addressed different



**Fig. 16.** Cost reduction regarding M-1 – Number of test cases required to achieve mutant coverage.

*selective mutation* strategies (even though being applied to the same set of programs), while Wong and Mathur (1995) used another set of programs and a different test generation tool.

Regarding RQ3 (*What are the savings (benefits) and loss of effectiveness (regarding mutation score) for the techniques?*), we found wide variations in results from both an inter-study and an intra-study perspective. Such variations, together with the number of metrics applied (see Section 6.1), make it harder to establish a baseline of studies for comparison in new experiments. Nevertheless, researchers can benefit from our results by using them as a reference to obtain more detailed information from specific studies, and to better design experiments that are comparable and reproducible.

#### 6.4. Further discussion about results for cost reduction measurement

This section presents more details about the measurements the studies used. These details affect how the results should be interpreted. We discuss six types of results: (1) results that represent close to (or even exactly) 100% cost reduction, (2) results that depend on either randomly or systematically defined thresholds, (3) results that are very context-specific, (4) imprecise results, and (5) results obtained from inconsistent subjects.

##### (i) Results that are close to (or even exactly) 100% cost reduction

Some studies reached 100% cost reduction or very close (Anbalagan and Xie, 2008; Papadakis and Malevris, 2009; Steimann and Thies, 2010; Kintis and Malevris, 2014). Although appealing, some were not generalizable and some did not measure all cost factors. For instance, Anbalagan and Xie (2008) automatically identified all equivalent mutants for a particular category of mutants: mutants of pointcut expressions of aspect-oriented programs written in the *AspectJ* language. Other costs that were not considered include the cost of running the code weaving compilation step (re-

quired by compilers of aspect-oriented languages), and the cost of running the mutant equivalence checking algorithm.

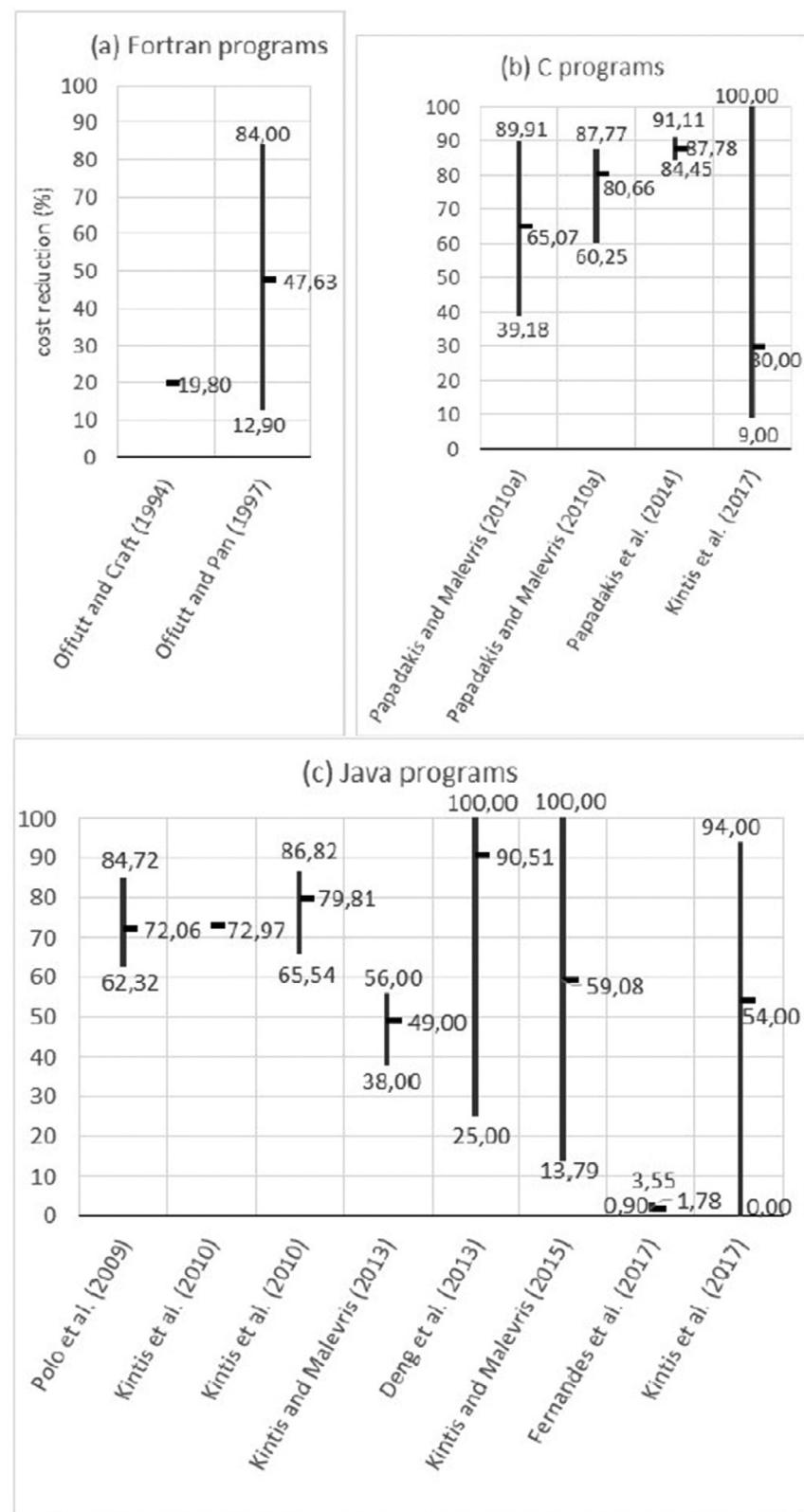
Another example is Papadakis and Malevris (2009)'s approach, which automatically generates tests to kill all non-equivalent mutants. That is, it produces mutation-adequate test suites. On the surface, this appears to eliminate the human cost of test data generation. However, the study did not include the manual mutant analysis cost, or the time to generate, compile, and execute mutants.

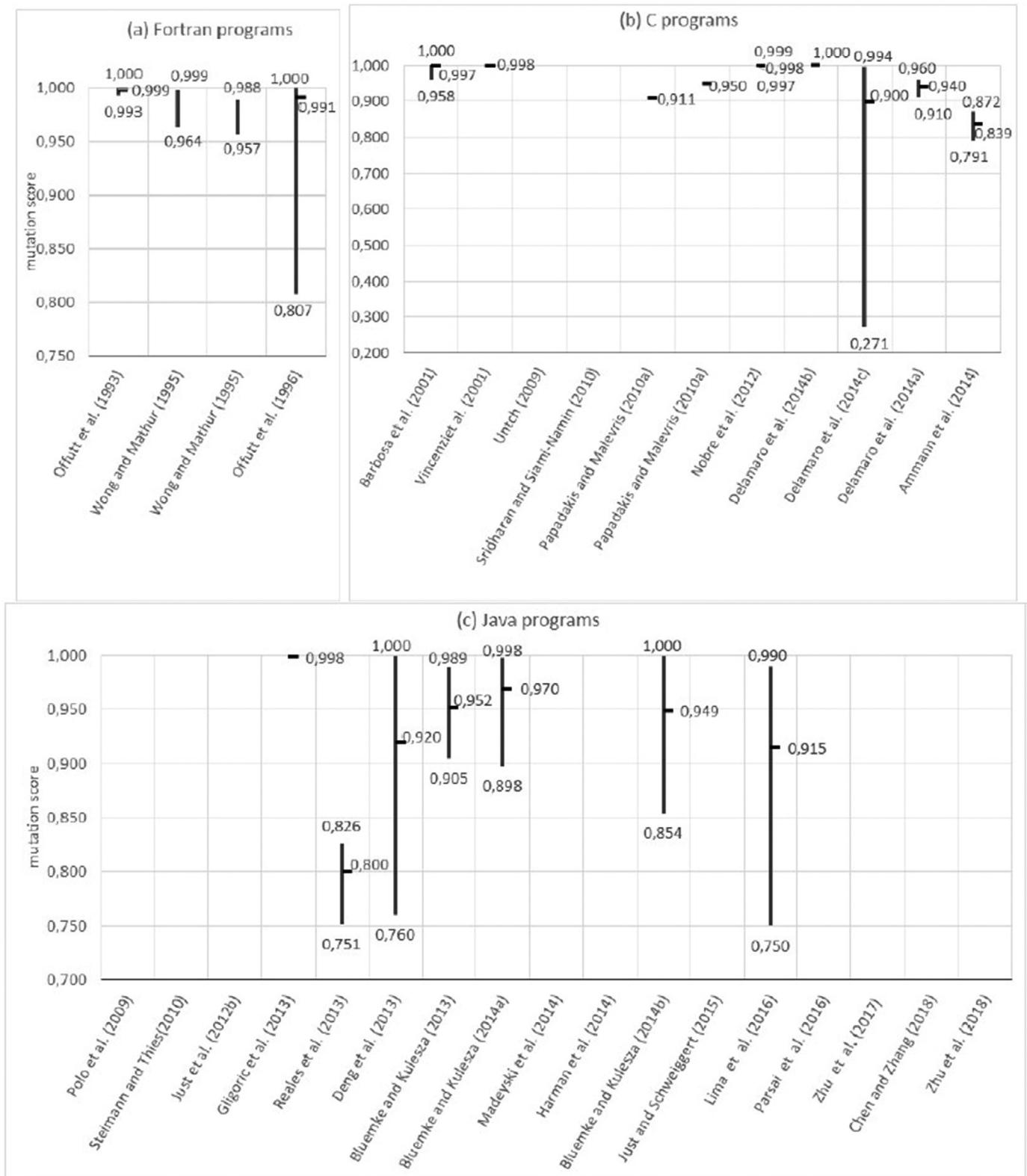
Steimann and Thies (2010) reached 99.77% cost reduction using the *Number of mutants to be executed*, but only for a single type of mutant (access modifier change). Kintis and Malevris (2014) reported cost reductions using *Number of equivalent mutants automatically detected* that widely ranged from 0 to 100% for Java programs. Patel and Hierons (2016) investigated mutation testing for non-deterministic programs. The authors performed a preliminary evaluation based on a single program and identified 100% of equivalent mutants. Both studies omitted other costs.

##### (ii) Results that depend on either randomly or systematically defined thresholds

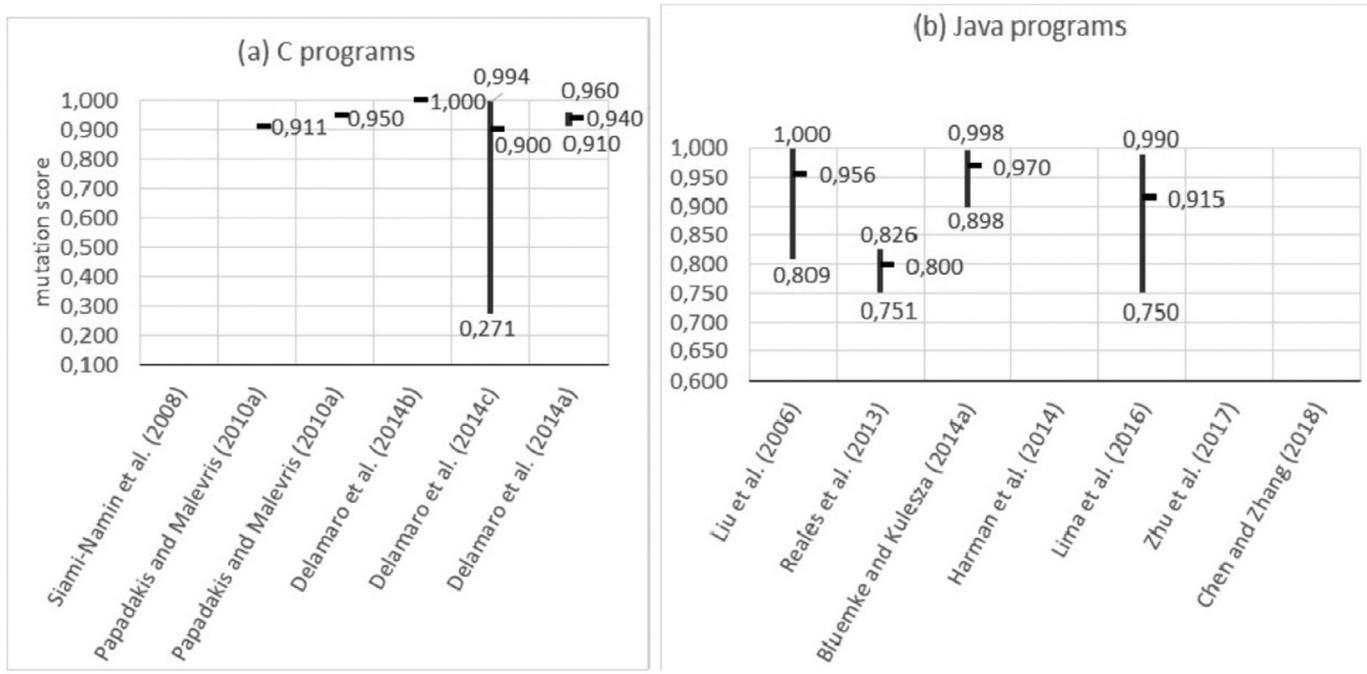
We did not present the cost reduction for some studies that used *evolutionary algorithms* (T-16), or for *optimization-related techniques* (T-15). This is because the results depended largely on the number of iterations of the algorithm, which is usually predefined arbitrarily based on factors such as timeout, achieved mutant coverage, and size of subjects. Some examples are the studies by Domínguez-Jiménez et al. (2009a), Fraser and Arcuri (2015), and Reuling et al. (2015).

Domínguez-Jiménez et al. (2009a) based their experiment on parameters such as population size and new individuals generated between different generations. Fraser and Arcuri (2015) focused on generating tests to achieve mutant coverage by applying a combination of techniques (*evolutionary algorithms*, *weak mutation*, and

**Fig. 17.** Cost reduction regarding M-6 – Number of equivalent mutants automatically detected.



**Fig. 18.** Mutation scores achieved in studies that applied metric M-2 – Number of mutants to be executed.



**Fig. 19.** Mutation scores achieved in studies that applied metric M-1 – Number of test cases required to achieve mutant coverage.

control flow analysis). Their studies compared mutation-based test case generation with branch coverage-based tests. One evaluation compared their approach with conventional mutation testing. They found an increase in mutation score, but the execution of the test generation algorithm was limited to four minutes.

Reuling et al. addressed mutation testing for software product lines (SPL). The authors set up three levels of SPL feature combinations, and were not able to use subjects that had more than 200 features (they ran out of time). Zhang et al. (2010b) investigated automatic test generation for Java programs and set a timeout after a predefined execution time or number of iterations, so it was impossible to extract cost reduction and mutation score values.

### (iii) Results that are very context-specific

Some studies not listed in Section 6.2 were done in very particular contexts. Examples include specific tools (Zhang et al., 2010b; Aichernig et al., 2013; Zhang et al., 2012), regression testing (Cachia et al., 2013; Chen and Zhang, 2018; Zhang et al., 2012) specific algorithms (Reales and Polo, 2013), and unusual sets of mutation operators (Kaminski and Ammann, 2011; Gopinath et al., 2018), or programming constructs (Wedyan and Ghosh, 2012; Ferrari et al., 2013). For instance, Zhang et al. (2010b) compared results from two test generation tools (*Pex* and *PexMutator*), but only *PexMutator* was designed for mutation testing, so there was no basis for comparison. Aichernig et al. (2013) explored test case generation based on state machine models and measured cost reduction based on results produced by a prior implementation by the same authors (Aichernig and Jöbstl, 2012).<sup>10</sup> Zhang et al. (2012) studied mutation testing for regression testing of Java programs and compared their results with results produced by the *JAVALANCHE* tool (Schuler et al., 2009).

Another example is by Reales and Polo (2013), who studied *Mutant execution speedup*. The authors discussed results for five alter-

native parallel execution implementations, but they were not compared with sequential mutant execution.

Kaminski and Ammann (2011) measured the *Number of test cases required to achieve mutant coverage* but only for logical expression mutants. Gopinath et al. (2018) created higher-order mutants but only used the statement deletion mutation operator. Wedyan and Ghosh (2012) and Ferrari et al. (2013) investigated mutation testing for *AspectJ* programs. They calculated cost reduction with respect to automatic detection of equivalent mutants of pointcut expressions, which are specific to aspect-oriented programming.

Other examples of varied contexts for mutation testing in which cost reduction has been measured are formal behavioural models (Devroey et al., 2017), and variability points in configurable systems (Al-Hajjaji et al., 2017).

### (iv) Imprecise results

Some studies reduced costs of mutation testing, but we could not extract precise cost reduction values from the reported results. An example was Kintis et al.'s study (Kintis et al., 2015), which classified 81% of mutants as killable with 71% precision. Even though this was 20% "superior" to previous approaches, we could not define a percentage of cost reduction from the reported data. Schuler et al.'s studies were similar (Schuler et al., 2009; Schuler and Zeller, 2013).

Kurtz et al. (2016) presented several cross-comparison experiments involving redundant, equivalent, and dominator mutants. However, they did not present precise values for cost reduction or mutation scores. Marcozzi et al. (2018) addressed the automatic identification of equivalent, redundant, and trivial mutants, which they called *polluting test objectives*. They did not compare their results with baseline or ground-truth numbers of equivalent or trivial mutants.

### (v) Results obtained from inconsistent subjects

We also found inconsistencies among the sets of subject programs used in some studies. For example, Papadakis and

<sup>10</sup> Aichernig and Jöbstl's study (Aichernig and Jöbstl, 2012) was subsumed by their 2013 study (Aichernig et al., 2013), as shown in Table A.2.

[Malevris \(2011b\)](#) measured *Mutant execution speedup* based on seven Java programs. In the same study, however, the authors compared results for automatic test generation based on 10 programs that included only five of the original seven programs. [Ma and Kim \(2016\)](#) studied cost reduction with *Number of test case executions* and *Mutant execution speedup*. While the cost reduction based on the first metric (*Number of test case executions*) was collected for three groups of subjects, the measurement of results for *Mutant execution speedup* was presented for a single group of programs. As another example, [Chen and Zhang \(2018\)](#) measured cost reduction based on four metrics: *Number of test cases required to achieve mutant coverage*; *Number of mutants to be executed*; *Mutant execution speedup*; and *Number of test case executions*. However, the paper presented results only for a subset of subjects for *Mutant execution speedup*.

## 7. Threats to validity

This section discusses threats to the validity of our SLR and how we mitigated them. The two main threats to validity are (i) the completeness of the results; and (ii) the classification of studies based on the lists of categories. Regarding the first, the results come from an analysis of 175 peer-reviewed studies published in journals, conference, and workshop proceedings. These studies subsumed other 22 studies. These studies do not include “gray literature,” including Masters and PhD theses, technical reports, technical specifications, and other unrefereed documents. This same focus on selecting only peer-reviewed studies was adopted in other SLRs ([Wohlin, 2014](#); [MacDonell et al., 2010](#)). Furthermore, other surveys on mutation testing (summarized in [Section 8](#)) show similar trends, particularly with respect to the growing number of published studies in the last 15 years ([Papadakis et al., 2019](#); [Jia and Harman, 2011](#); [Madeyski et al., 2014](#); [Silva et al., 2017](#)). The inclusion of additional, non-peer-reviewed, studies would probably not change these results, although it would increase the numbers in this paper.

The backward snowballing technique was used after the first search round with only one level of depth. [Wohlin \(2014\)](#) suggested snowballing should be an iterative process that ends only when no new studies are found. Therefore, our procedure to apply snowballing may have created another threat to completeness. To mitigate this threat, we used [Wohlin’s](#) alternative suggestion for contacting authors of primary studies to ask for additional primary studies. In particular, we contacted every author of selected primary studies to ask for additional papers to be analyzed (see [Section 3.3](#) for more details). As a result, we analyzed 135 more studies suggested in 42 replies, and added 25 items to our final set of studies. The same search procedure was adopted by [Jia and Harman \(2011\)](#), who called it a ‘transitive closure’ on the literature.

For threat (ii), the initial classification of studies based on the *do fewer*, *do smarter*, and *do faster* categories ([Offutt and Untch, 2000](#)) is no longer adequate. When these categories were first defined, fewer techniques were known. More recent techniques are hard to categorize into the *do fewer*, *do smarter*, and *do faster* groups. Category *T-15* for cost reduction techniques reflects this complexity. In addition, since 2000 it has become more common to combine multiple techniques, as can be seen in [Figs. 6](#) and [10](#). As far as possible, the classification we present in this paper was performed in a discerning and unbiased manner.

Still regarding threat (ii), as discussed in [Section 5.3](#), several studies were assigned category *T-15* (*optimization of generation, execution and analysis of mutants*) as an “others” category. It might be possible to further analyze those techniques to identify subcategories, increasing the list of cost reduction techniques. It might

also be reasonable to combine some categories, as discussed in [Section 5.4](#).

## 8. Related work

In some sense, every paper referenced here is “related.” For our related work section, therefore, we focus on other surveys on mutation testing, with a specific emphasis on surveys that focus on cost reduction.

[Polo and Reales \(2010\)](#) summarized a small set of reference papers that tried to reduce costs of the main steps in mutation testing: (1) mutant generation, (2) test case generation and execution, and (3) result analysis. The authors suggested techniques such as *T-1 (random mutation)* and *T-12 (selective mutation)* for step (1); *T-7 (minimization and prioritization of test sets)* and *T-3 (weak mutation)* for step (2); and *T-2 (higher order mutation)* for step (3). [Polo and Reales](#) did not use systematic search criteria to identify the techniques and they did not classify studies by technique. They only analyzed 10 cost reduction-related studies, all of which we include here. One interesting note regarding [Polo and Reales’s](#) work is that the authors grouped together *test case generation* and *number of test case executions* as a common goal for cost reduction, as we did with our primary goal *PG-4* ([Section 5.2](#)).

[Jia and Harman \(2011\)](#) published what is probably the most extensive survey of mutation testing in general to date. They identified several cost reduction techniques, including equivalent mutant detection. [Jia and Harman](#) identified the first peer-reviewed cost reduction paper as the 1982 paper that introduced *weak mutation* ([Howden, 1982](#)). We include this paper in [Fig. 6](#), but not as a selected paper, as it did not focus on cost reduction, did not implement the idea, and included no empirical validation. Our SLR search found the first study that used *weak mutation* to reduce cost to be [Marshall et al. \(1990\)](#) in 1990. [Jia and Harman](#) summarized 66 mutation testing publications they classified as cost reduction-related, published from 1982 through 2009, including 52 peer-reviewed papers and 14 non-indexed (mostly, gray literature) items such as PhD and Masters Theses, technical reports, and Doctoral Symposium papers. In the same timeframe, our paper includes 44 peer-reviewed studies specifically focused on cost reduction.

[Madeyski et al. \(2014\)](#) published results of an SLR on the equivalent mutant problem. The authors classified 24 studies in three categories, based on the goal: *detecting equivalent mutants*, *avoiding equivalent mutant generation*, and *suggesting equivalent mutants*. Moreover, they identified 17 methods that contribute to achieve such goals. This paper includes 15 of [Madeyski et al.’s](#) studies; the others either did not focus on cost reduction, were not peer-reviewed, or are included here as related work.

[Silva et al. \(2017\)](#) performed an SLR on the use of Search-based Software Testing (SBST) techniques in mutation testing. The authors identified 69 studies published between 1998 and 2014, focusing primarily on the meta-heuristics applied. As opposed to this study, [Silva et al.](#) provided an overall characterization of their topic without focusing on cost reduction-related studies. They also included non-peer-reviewed studies.

In the most recent related work, [Papadakis et al. \(2019\)](#) updated [Jia and Harman \(2011\)](#)’s survey to include studies published from 2008 to 2017 (inclusive). The authors stated they selected a total of 502 publications (although only 405 appear in the reference list), of which 260 addressed mutation testing problems (including supporting tools), while 242 either used mutation testing to assess test quality, or tackled problems unrelated to mutation testing. According to [Papadakis et al.](#), their study focused on conferences, symposia, workshops, and journals. When contrasted with our work, we searched the same databases that index the conference proceedings and journals they analyzed. Therefore,

papers in their dataset that are not in our dataset are probably due to differences in our selection criteria. Papadakis et al.'s paper did not report the selection criteria so we could not compare directly.

## 9. Conclusion, recommendations, and future work

Mutation testing is a highly investigated and very effective way to generate tests and to assess test quality. However, it is also very expensive. The four major cost factors are (i) a large number of mutants are generated and executed, even for small programs; (ii) test data generation; (iii) test suites are large; and (iv) equivalent mutants. In response, researchers have developed many approaches to reduce the cost of mutation. This paper summarized results of a systematic literature review that characterized the history and the state-of-the-art of cost reduction for mutation testing. We updated and extended our prior paper (Ferrari et al., 2018a), analyzed the evolution of research on the topic, and summarized its underlying goals, techniques, metrics used, and results achieved.

We analyzed a total of 175 peer-reviewed studies, of which 153 present either original or updated approaches and results for mutation cost reduction. The 21 techniques identified use six main goals to reduce mutation cost. Apart from the *Optimization-related* category (which groups a set of techniques that could not be otherwise classified), the techniques investigated the most are *control-flow analysis*, *evolutionary algorithms*, *selective mutation*, *higher order mutation*, and *data-flow analysis*.

Experimental results were measured with 18 metrics, plus a range of complex, study-specific metrics that were not addressed in this review. The mostly commonly used metrics are the *Number of mutants to be executed*, *Mutant execution speedup*, and the *Number of test cases required to achieve mutant coverage*.

**Summary of recommendations:** Based on the research questions investigated in this work, we summarize several recommendations.

- *Targeting the primary goals:* Our results reveal that *reducing the number of mutants* (PG-1) was the most common primary goal, followed by *executing faster* (PG-3), then *automatically detecting equivalent mutants* (PG-2). If we look at the last 10 years considered in this research (2009 to 2018), we found 113 studies, 15 of which focused on PG-2 (Papadakis and Le Traon, 2013; Schuler et al., 2009; Schuler and Zeller, 2013; Kintis and Malevris, 2015; Papadakis et al., 2014; Kintis and Malevris, 2013; Ferrari et al., 2013; Kintis et al., 2015; Patel and Hierons, 2016; Holling et al., 2016; Kintis et al., 2017; Devroey et al., 2017; Marcozzi et al., 2018; Delgado-Pérez et al., 2017; McMinn et al., 2018). Thus, only 13% of the studies from the past decade directly addressed a problem that is widely considered the major obstacle to practical adoption of mutation testing (Jia and Harman, 2011; Madeyski et al., 2014). The total number of studies related to equivalent mutants corroborates this observation. In a recent survey, Papadakis et al. (2019) also found that equivalent mutant detection remains an open problem. Therefore, we recommend more research into the problem.

We found 12 studies (11 in the past decade) that addressed *avoiding the creation of certain mutants* (PG-5). Rules to prevent creation of equivalent and redundant mutant can be embedded into mutation operators. Given that this goal impacts the largest number of other goals, including reducing the number of equivalent mutants, we recommend more research on this problem.

- *Applying combined techniques:* We found that mutation cost reduction research has become more interdisciplinary over time.

It is true that traditional mutation-specific techniques such as changing the set of operators (e.g. *selective mutation* and *one-op mutation*), *weak mutation*, *higher-order mutation*, and *minimal mutation* are still studied. However, we also found frequent and increasing interest in techniques borrowed from other Computer Science subareas such as software analysis and artificial intelligence.

We also found several studies that combined techniques to achieve substantial savings. For example, *weak mutation* and *optimization* were applied together to reduce the number of mutants to be executed and mutant execution time, with 97% and 94% savings (Zhu et al., 2018). As another example, *control-flow analysis* has been combined with *constrained mutation* and *random mutation* in industrial setting (Petrović and Ivanković, 2018). This suggests that stronger collaboration with researchers from different subareas can be very productive.

- *Standardizing and sharing experimental methods, materials and reports:* While collecting cost reduction results from the selected studies, we found lots of differences in experimental design that affected the results. The level of detail in the papers also varied greatly. As a consequence, we could not compare studies with similar goals (e.g. to run meta-analysis of results). Moreover, very few experimental materials are available online. Therefore, we recommend the community to create benchmarks that include all software artifacts, including software under test, test cases that reveal actual faults, test cases that kill mutants, mutant programs, mutation execution results in the form of killing matrices, as well as tools to produce and handle such artifacts and directions on their use. This would greatly help researchers design experiments that can be compared and reproduced.

**Future work:** As with any SLR, it is unlikely that we found all primary studies, as acknowledged as a validity threat in Section 7. Even though several search strategies were used, well-known limitations of scientific writing and repositories mean no search can be perfect. Thus we hope to see future updates to this systematic literature review, performed by either ourselves or other researchers.

## Acknowledgements

We would like to thank all 42 authors that suggested studies to be included into the SLR dataset. We also thank George Mason University for hosting Fabiano Ferrari during part of his participation in this research. Finally, we are also thankful for the financial support, detailed as follows: Fabiano Ferrari was supported by the Fundação de Amparo à Pesquisa do Estado de São Paulo - Brasil, grant #2016/21251-0; and Conselho Nacional de Desenvolvimento Científico e Tecnológico - Brasil, grant #306310/2016-3; Alessandro Pizzoleto was supported by the Coordenação de Aperfeiçoamento de Pessoal de Nível Superior - Finance Code 001. Jeff Offutt was partly supported by KKS (The Knowledge Foundation), by project 20130085, Testing of Critical System Characteristics (TOCSYC); and partly by a Google educational grant, Self-Paced Learning Increases Retention and Capacity (SPARC).

## Appendix A. Selected studies and subsumed studies

This appendix includes tables that list all primary studies that were selected in this SLR (Table A.1), and a table that lists all subsumed studies (Table A.2). Tables of selected studies are sorted by author name and year of study publication, while tables of subsumed studies are sorted by year of publication of the updated or extended studies.

**Table A.1**  
Selected studies.

Author/Year	Title	Publisher	Journal/Event
Abuljadayel and Wedyan (2018)	An Approach for the Generation of Higher Order Mutants Using Genetic Algorithms	Modern Education and Computer Science Press	International Journal of Intelligent Systems and Applications
Aichernig et al. (2013)	Incremental Refinement Checking for Test Case Generation	Springer	International Conference on Quality Software (QSIC)
Aichernig et al. (2014)	Model-Based Mutation Testing of an Industrial Measurement Device	Springer	International Conference Tests and Proofs (TAP)
Al-Hajjaji et al. (2017)	Efficient Mutation Testing in Configurable Systems	IEEE	International Workshop on Variability and Complexity in Software Design (VACE)
Alexander et al. (2002)	Mutation of Java Objects	IEEE	International Symposium on Software Reliability Engineering (ISSRE)
Ammann et al. (2014)	Establishing Theoretical Minimal Sets of Mutants	IEEE	International Conference on Software Testing, Verification and Validation (ICST)
Anbalagan and Xie (2008)	Automated Generation of Pointcut Mutants for Testing Pointcuts in AspectJ Programs	IEEE	International Symposium on Software Reliability Engineering (ISSRE)
Adamopoulos et al. (2004)	How to Overcome the Equivalent Mutant Problem and Achieve Tailored Selective Mutation Using Co-evolution	Springer	Genetic and Evolutionary Computation Conference (GECCO)
Ayari et al. (2007)	Automatic Mutation Test Input Data Generation via Ant Colony	ACM	Genetic and Evolutionary Computation Conference (GECCO)
Barbosa et al. (2001)	Toward the Determination of Sufficient Mutant Operators for C	Wiley	Software Testing, Verification and Reliability
Bashir and Nadeem (2017)	Improved Genetic Algorithm to Reduce Mutation Testing Cost	IEEE	IEEE Access
Baudry et al. (2005)	From Genetic to Bacteriological Algorithms for Mutation-Based Testing	Wiley	Software Testing, Verification and Reliability
Belli and Beyazit (2015)	Exploiting Model Morphology for Event-Based Testing	IEEE	IEEE Transactions on Software Engineering
Bluemke and Kulesza (2013)	Reduction of Computational Cost in Mutation Testing by Sampling Mutants	Springer	International Conference on Dependability and Complex Systems (DepCos-RELCOMEX)
Bluemke and Kulesza (2014b)	Reductions of Operators in Java Mutation Testing	Springer	International Conference on Dependability and Complex Systems (DepCos-RELCOMEX)
Bluemke and Kulesza (2014a)	Reduction in Mutation Testing of Java Classes	IEEE	International Conference on Software Engineering and Applications (ICSOFT-EA)
Bogacki and Walter (2006a)	Aspect-oriented Response Injection: an Alternative to Classical Mutation Testing	Springer	IFIP Working Conference on Software Engineering Techniques: Design for Quality (SET)
Choi and Mathur (1993)	High-Performance Mutation Testing	Elsevier	Journal of Systems and Software
Cachia et al. (2013)	Towards Incremental Mutation Testing	Elsevier	Validation Strategies for Software Evolution Workshop (VSSE)
Chen and Zhang (2018)	Speeding up Mutation Testing via Regression Test Selection: An Extensive Study	IEEE	International Conference on Software Testing, Verification and Validation (ICST)
Delamaro et al. (2014a)	Experimental Evaluation of SDL and One-Op Mutation for C	IEEE	International Conference on Software Testing, Verification and Validation (ICST)
Delamaro et al. (2014c)	Designing Deletion Mutation Operators	IEEE	International Conference on Software Testing, Verification and Validation (ICST)
Delamaro et al. (2014b)	Growing a Reduced Set of Mutation Operators	IEEE	Brazilian Symposium on Software Engineering (SBES)
Delgado-Pérez et al. (2017b)	GiGAN: Evolutionary Mutation Testing for C++ Object-oriented Systems	ACM	ACM Symposium on Applied Computing (SAC)
Delgado-Pérez et al. (2017c)	Assessment of C++ Object-Oriented Mutation Operators: A Selective Mutation Approach	Wiley	Software Testing, Verification and Reliability
Delgado-Pérez et al. (2017a)	Assessment of Class Mutation Operators for C++ with the MuCPP Mutation System	Elsevier	Information and Software Technology
Delgado-Pérez et al. (2017)	Using Evolutionary Computation to Improve Mutation Testing	Springer	International Work-Conference on Artificial Neural Networks (IWANN)
DeMillo et al. (1991)	Compiler-integrated Program Mutation	IEEE	IEEE Annual Computer Software and Applications Conference (COMPSAC)
DeMillo and Offutt (1991)	Constraint-based Automatic Test Data Generation	IEEE	IEEE Transactions on Software Engineering
DeMillo and Offutt (1993)	Experimental Results from an Automatic Test Case Generator	ACM	ACM Transactions on Software Engineering and Methodology
Deng et al. (2013)	Empirical Evaluation of the Statement Deletion Mutation Operator	IEEE	International Conference on Software Testing, Verification and Validation (ICST)
Denisov and Pankevich (2018)	Mull It Over: Mutation Testing Based on LLVM	IEEE	International Workshop on Mutation Analysis (Mutation)
Derezińska and Rudnik (2012)	Quality Evaluation of Object-Oriented and Standard Mutation Operators Applied to C# Programs	Springer	International Conference on Modelling Techniques and Tools for Computer Performance Evaluation (TOOLS)
Derezińska (2013)	A Quality Estimation of Mutation Clustering in C# Programs	Springer	International Conference on Dependability and Complex Systems (DepCos-RELCOMEX)
Derezińska and Hałas (2015)	Improving Mutation Testing Process of Python Programs	Springer	Computer Science On-line Conference Software Engineering in Intelligent Systems (CSOC)

(continued on next page)

**Table A.1 (continued)**

Author/Year	Title	Publisher	Journal/Event
Derezińska (2016)	Evaluation of Deletion Mutation Operators in Mutation Testing of C# Programs	Springer	International Conference on Dependability and Complex Systems (DepCoS-RELCOMEX)
Derezińska and Rudnik (2017)	Evaluation of Mutant Sampling Criteria in Object-Oriented Mutation Testing	Polskie Towarzystwo Informatyczne	Federated Conference on Computer Science and Information Systems (ACSIS)
Devroey et al. (2016)	Featured Model-based Mutation Analysis	ACM	International Conference on Software Engineering (ICSE)
Devroey et al. (2017)	Automata Language Equivalence vs. Simulations for Model-Based Mutant Equivalence: An Empirical Evaluation	IEEE	International Conference on Software Testing, Verification and Validation (ICST)
Domínguez-Jiménez et al. (2009a)	GAMERA: An Automatic Mutant Generation System for WS-BPEL Compositions	IEEE	IEEE European Conference on Web Services (ECOWS)
Domínguez-Jiménez et al. (2011)	Evolutionary Mutation Testing	Elsevier	Information and Software Technology
Durelli et al. (2012)	Toward Harnessing High-Level Language Virtual Machines for Further Speeding Up Weak Mutation Testing	IEEE	International Conference on Software Testing, Verification and Validation (ICST)
Fernandes et al. (2017)	Avoiding Useless Mutants	ACM	International Conference on Generative Programming: Concepts and Experiences (GPCE)
Ferrari et al. (2013)	Towards the Practical Mutation Testing of AspectJ Programs	Elsevier	Science of Computer Programming
Matnei Filho and Vergilio (2015)	A Mutation and Multi-objective Test Data Generation Approach for Feature Testing of Software Product Lines	IEEE	Brazilian Symposium on Software Engineering (SBES)
Fleyshgakker and Weiss (1994)	Efficient Mutation Analysis: A New Approach	ACM	International Symposium on Software Testing and Analysis (ISSTA)
Fraser and Zeller (2012)	Mutation-Driven Generation of Oracles and Unit Tests	IEEE	IEEE Transactions on Software Engineering
Fraser and Arcuri (2015)	Achieving Scalable Mutation-based Generation of Whole Test Suites	Springer	Empirical Software Engineering
Gligoric et al. (2010)	MuTMuT: Efficient Exploration for Mutation Testing of Multithreaded Code	IEEE	International Conference on Software Testing, Verification and Validation (ICST)
Gligoric et al. (2012)	Efficient Mutation Testing of Multithreaded Code	Wiley	Software Testing, Verification and Reliability
Gligoric et al. (2013)	Selective Mutation Testing for Concurrent Code	ACM	International Symposium on Software Testing and Analysis (ISSTA)
Gopinath et al. (2016b)	Topsy-Turvy: A Smarter and Faster Parallelization of Mutation Analysis	ACM	International Conference on Software Engineering (ICSE)
Gopinath et al. (2017)	Mutation Reduction Strategies Considered Harmful	IEEE	IEEE Transactions on Software Reliability
Gopinath et al. (2018)	If You Can't Kill a Supermutant, You Have a Problem	IEEE	International Workshop on Mutation Analysis (Mutation)
Harman et al. (2000)	The Relationship Between Program Dependence and Mutation Analysis	Kluwer	Mutation 2000 Symposium
Harman et al. (2014)	Angels and Monsters: An Empirical Investigation of Potential Test Effectiveness and Efficiency Improvement from Strongly Subsuming Higher Order Mutation	ACM	International Conference on Automated Software Engineering (ASE)
Hierons et al. (1999)	Using Program Slicing to Assist in the Detection of Equivalent Mutants	Wiley	Software Testing, Verification and Reliability
Henard et al. (2014)	Mutation-Based Generation of Software Product Line Test Configurations	Springer	International Symposium Search-Based Software Engineering (SSBSE)
Holling et al. (2016)	Nequivack: Assessing Mutation Score Confidence	IEEE	International Workshop on Mutation Analysis (Mutation)
Hu et al. (2011)	An Analysis of OO Mutation Operators	IEEE	International Workshop on Mutation Analysis (Mutation)
Iida and Takada (2017)	Reducing Mutants with Mutant Killable Precondition	IEEE	International Workshop on Mutation Analysis (Mutation)
Inozemtseva et al. (2013)	Using Fault History to Improve Mutation Reduction	ACM	International Symposium on the Foundations of Software Engineering (FSE)
Jackson and Woodward (2000)	Parallel Firm Mutation of Java Programs	Kluwer	Mutation 2000 Symposium
Ji et al. (2009)	A Novel Method of Mutation Clustering Based on Domain Analysis	Knowledge Systems Institute Graduate School	International Conference on Software Engineering and Knowledge Engineering (SEKE)
Just et al. (2011)	Using Conditional Mutation to Increase the Efficiency of Mutation Analysis	ACM	International Workshop on Automation of Software Test (AST)
Just et al. (2012b)	Using Non-Redundant Mutation Operators and Test Suite Prioritization to Achieve Efficient and Scalable Mutation Analysis	IEEE	International Symposium on Software Reliability Engineering (ISSRE)
Just et al. (2014a)	Efficient Mutation Analysis by Propagating and Partitioning Infected Execution States	ACM	International Symposium on Software Testing and Analysis (ISSTA)
Just and Schweiggert (2015)	Higher Accuracy and Lower Run Time: Efficient Mutation Analysis Using Non-redundant Mutation Operators	Wiley	Software Testing, Verification and Reliability

(continued on next page)

**Table A.1** (continued)

Author/Year	Title	Publisher	Journal/Event
Just et al. (2017)	Inferring Mutant Utility from Program Context	ACM	International Symposium on Software Testing and Analysis (ISSTA)
Kaminski and Ammann (2009)	Using a Fault Hierarchy to Improve the Efficiency of DNF Logic Mutation Testing	IEEE	International Conference on Software Testing, Verification and Validation (ICST)
Kaminski and Ammann (2011)	Reducing Logic Test Set Size while Preserving Fault Detection	Wiley	Software Testing, Verification and Reliability
Kaminski et al. (2011b)	A Logic Mutation Approach to Selective Mutation for Programs and Queries	Elsevier	Information and Software Technology
Kaminski et al. (2013)	Improving logic-based testing	Elsevier	Journal of Systems and Software
Kim et al. (2012)	Combining Weak and Strong Mutation for a Noninterpretive Java Mutation System	Wiley	Software Testing, Verification and Reliability
Kintis et al. (2010)	Evaluating Mutation Testing Alternatives: A Collateral Experiment	IEEE	Asia-Pacific Software Engineering Conference (APSEC)
Kintis and Malevris (2013)	Identifying More Equivalent Mutants via Code Similarity	IEEE	Asia-Pacific Software Engineering Conference (APSEC)
Kintis and Malevris (2014)	Using Data Flow Patterns for Equivalent Mutant Detection	IEEE	International Workshop on Mutation Analysis (Mutation)
Kintis et al. (2015)	Employing Second-order Mutation for Isolating First-order Equivalent Mutants	Wiley	Software Testing, Verification and Reliability
Kintis and Malevris (2015)	MEDIC: A Static Analysis Framework for Equivalent Mutant Identification	Elsevier	Information and Software Technology
Kintis et al. (2017)	Detecting Trivial Mutant Equivalences via Compiler Optimisations	IEEE	IEEE Transactions on Software Engineering
Krauser et al. (1991)	High Performance Software Testing on SIMD Machines	IEEE	IEEE Transactions on Software Engineering
Kurtz et al. (2016)	Analyzing the Validity of Selective Mutation with Dominator Mutants	ACM	International Symposium on the Foundations of Software Engineering (FSE)
Lacerda and Ferrari (2014)	Towards the Establishment of a Sufficient Set of Mutation Operators for AspectJ Programs	SBC	Brazilian Symposium on Systematic and Automated Software Testing (SAST)
Li et al. (2015)	Mutation Testing in Practice Using Ruby	IEEE	International Workshop on Mutation Analysis (Mutation)
Lima et al. (2016)	Evaluating Different Strategies for Reduction of Mutation Testing Costs	ACM	Brazilian Symposium on Systematic and Automated Software Testing (SAST)
Liu et al. (2006)	An Approach to Test Data Generation for Killing Multiple Mutants	IEEE	International Conference on Software Maintenance (ICSM)
Ma et al. (2005)	MuJava: An Automated Class Mutation System	Wiley	Software Testing, Verification and Reliability
Ma and Kim (2016)	Mutation Testing Cost Reduction by Clustering Overlapped Mutants	Elsevier	Journal of Systems and Software
Madeyski et al. (2014)	Overcoming the Equivalent Mutant Problem: A Systematic Literature Review and a Comparative Experiment of Second Order Mutation	IEEE	IEEE Transactions on Software Engineering
Marcozzi et al. (2018)	Time to Clean Your Test Objectives	IEEE	International Conference on Software Engineering (ICSE)
Marshall et al. (1990)	Static Dataflow-aided Weak Mutation Analysis (SDAWM)	Elsevier	Information and Software Technology
Reales and Polo (2013)	Parallel Mutation Testing	wiley	Software Testing, Verification and Reliability
Reales et al. (2013)	Validating Second-Order Mutation at System Level	IEEE	IEEE Transactions on Software Engineering
Reales and Polo (2014)	Reducing Mutation Costs Through Uncovered Mutants	Wiley	Software Testing, Verification and Reliability
McMinn et al. (2016)	Virtual Mutation Analysis of Relational Database Schemas	IEEE	International Workshop on Automation of Software Test (AST)
McMinn et al. (2018)	Automatic Detection and Removal of Ineffective Mutants for the Mutation Analysis of Relational Database Schemas	IEEE	IEEE Transactions on Software Engineering
Mresa and Bottaci (1999)	Efficiency of Mutation Operators and Selective Mutation Strategies: an Empirical Study	Wiley	Software Testing, Verification and Reliability
Siami-Namin and Andrews (2006)	Finding Sufficient Mutation Operators via Variable Reduction	IEEE	International Workshop on Mutation Analysis (Mutation)
Siami-Namin et al. (2008)	Sufficient Mutation Operators for Measuring Test Effectiveness	ACM	International Conference on Software Engineering (ICSE)
Nobre et al. (2012)	Reducing Interface Mutation Costs with Multiobjective Optimization Algorithms	IGI Global	International Journal of Natural Computing Research
Offutt et al. (1992)	Mutation Testing of Software Using a MIMD Computer	Wiley	International Conference on Parallel Processing (ICPP)
Offutt et al. (1993)	An Experimental Evaluation of Selective Mutation	IEEE	International Conference on Software Engineering (ICSE)
Offutt and Craft (1994)	Using Compiler Optimization Techniques to Detect Equivalent Mutants	Wiley	Software Testing, Verification and Reliability
Offutt and Lee (1994)	An Empirical Evaluation of Weak Mutation	IEEE	IEEE Transactions on Software Engineering
Offutt et al. (1996)	An Experimental Determination of Sufficient Mutant Operators	ACM	ACM Transactions on Software Engineering and Methodology

(continued on next page)

**Table A.1** (continued)

Author/Year	Title	Publisher	Journal/Event
Offutt and Pan (1997)	Automatically Detecting Equivalent Mutants and Infeasible Paths	Wiley	Software Testing, Verification and Reliability
Offutt et al. (1999)	The Dynamic Domain Reduction Procedure for Test Data Generation	Wiley	Software Practice and Experience
Offutt et al. (2006)	The Class-Level Mutants of MuJava	ACM	International Workshop on Automation of Software Test (AST)
Oliveira et al. (2013)	A Coevolutionary Algorithm to Automatic Test Case Selection and Mutant in Mutation Testing	IEEE	IEEE Congress on Evolutionary Computation (CEC)
Omar and Ghosh (2012)	An Exploratory Study of Higher Order Mutation Testing in Aspect-Oriented Programming	IEEE	International Symposium on Software Reliability Engineering (ISSRE)
Papadakis and Malevris (2009)	An Effective Path Selection Strategy for Mutation Testing	IEEE	Asia-Pacific Software Engineering Conference (APSEC)
Papadakis and Malevris (2010a)	An Empirical Evaluation of the First and Second Order Mutation Testing Strategies	IEEE	International Workshop on Mutation Analysis (Mutation)
Papadakis and Malevris (2010b)	Automatic Mutation Test Case Generation via Dynamic Symbolic Execution	IEEE	International Symposium on Software Reliability Engineering (ISSRE)
Papadakis and Malevris (2011a)	Automatic Mutation based Test Data Generation	ACM	Genetic and Evolutionary Computation Conference (GECCO)
Papadakis and Malevris (2011b)	Automatically Performing Weak Mutation with the Aid of Symbolic Execution, Concolic Testing and Search-based Testing	Springer	Software Quality Journal
Papadakis and Malevris (2012)	Mutation Based Test Case Generation Via a Path Selection Strategy	Elsevier	Information and Software Technology
Papadakis and Le Traon (2013)	Mutation Testing Strategies using Mutant Classification	ACM	ACM Symposium on Applied Computing (SAC)
Papadakis et al. (2014)	Mitigating the Effects of Equivalent Mutants with Mutant Classification Strategies	Elsevier	Science of Computer Programming
Parsai et al. (2016)	Evaluating Random Mutant Selection at Class-level in Projects with Non-adequate Test Suites	ACM	International Conference on Evaluation and Assessment in Software Engineering (EASE)
Patel and Hierons (2016)	Resolving the Equivalent Mutant Problem in the Presence of Non-determinism and Coincidental Correctness	Springer	International Conference on Testing Software and Systems (ICTSS)
Patrick et al. (2012)	MESSI: Mutant Evaluation by Static Semantic Interpretation	IEEE	International Conference on Software Testing, Verification and Validation (ICST)
Petrović and Ivanković (2018)	State of Mutation Testing at Google	IEEE	International Conference on Software Engineering - Software Engineering in Practice Track (ICSE-SEIP)
Praphamontripong and Offutt (2017)	Finding Redundancy in Web Mutation Operators	IEEE	International Workshop on Mutation Analysis (Mutation)
Quyen et al. (2016)	Improving Mutant Generation for Simulink Models Using Genetic Algorithm	IEEE	International Conference on Electronics, Information, and Communications (ICEIC)
Reuling et al. (2015)	Fault-based Product-line Testing: Effective Sample Generation Based on Feature-Diagram Mutation	ACM	International Conference on Software Product Line (SPLC)
Sahinoğlu and Spafford (1990)	A Bayes Sequential Statistical Procedure for Approving Software Products	Elsevier	IFIP Conference on Approving Software Products (ASP)
Schuler et al. (2009)	Efficient Mutation Testing by Checking Invariant Violations	ACM	International Symposium on Software Testing and Analysis (ISSTA)
Schuler and Zeller (2013)	Covering and Uncovering Equivalent Mutants	Wiley	Software Testing, Verification and Reliability
Sridharan and Siami-Namin (2010)	Prioritizing Mutation Operators based on Importance Sampling	IEEE	International Symposium on Software Reliability Engineering (ISSRE)
Steimann and Thies (2010)	From Behaviour Preservation to Behaviour Modification: Constraint-Based Mutant Generation	ACM	International Conference on Software Engineering (ICSE)
Sun et al. (2017a)	An Empirical Study on Mutation Testing of WS-BPEL Programs	Oxford University Press	Computer Journal
Sun et al. (2017b)	A Path-aware Approach to Mutant Reduction in Mutation Testing	Elsevier	Information and Software Technology
Tuya et al. (2007)	Mutating database queries	Elsevier	Information and Software Technology
Untch et al. (1993)	Mutation Analysis Using Mutant Schemata	ACM	International Symposium on Software Testing and Analysis (ISSTA)
Untch et al. (1997)	TUMS: Testing Using Mutant Schemata	ACM	Annual Southeast Regional Conference (ACM-SE)
Untch (2009)	On Reduced Neighborhood Mutation Analysis Using a Single Mutagenic Operator	ACM	Annual Southeast Regional Conference (ACM-SE)
Polo et al. (2009)	Decreasing the Cost of Mutation Testing with Second-Order Mutants	Wiley	Software Testing, Verification and Reliability
Vincenzi et al. (2001)	Unit and Integration Testing Strategies for C Programs Using Mutation	Wiley	Software Testing, Verification and Reliability
Vincenzi et al. (2002)	Bayesian-Learning Based Guidelines to Determine Equivalent Mutants	World Scientific Publishing	International Journal of Software Engineering and Knowledge Engineering

(continued on next page)

**Table A.1** (continued)

Author/Year	Title	Publisher	Journal/Event
Wang et al. (2017)	Faster Mutation Analysis via Equivalence Modulo States	ACM	International Symposium on Software Testing and Analysis (ISSTA)
Wedyan and Ghosh (2012)	On Generating Mutants for AspectJ Programs	Elsevier	Information and Software Technology
Weiss and Fleysgakker (1993)	Improved Serial Algorithms for Mutation Analysis	ACM	International Symposium on Software Testing and Analysis (ISSTA)
Wong et al. (1994)	Constrained Mutation in C Programs	Brazilian Computer Society	Brazilian Symposium on Software Engineering (SBES)
Wong and Mathur (1995)	Reducing the Cost of Mutation Testing: An Empirical Study	Elsevier	Journal of Systems and Software
Wright et al. (2013)	Efficient Mutation Analysis of Relational Database Structure Using Mutant Schemata and Parallelisation	IEEE	International Workshop on Mutation Analysis (Mutation)
Zhang et al. (2010a)	Is Operator-Based Mutant Selection Superior to Random Mutant Selection?	ACM	International Conference on Software Engineering (ICSE)
Zhang et al. (2010b)	Test Generation via Dynamic Symbolic Execution for Mutation Testing	IEEE	International Conference on Software Maintenance (ICSM)
Zhang et al. (2012)	Regression Mutation Testing	ACM	International Symposium on Software Testing and Analysis (ISSTA)
Zhang et al. (2013b)	Faster Mutation Testing Inspired by Test Prioritization and Reduction	ACM	International Symposium on Software Testing and Analysis (ISSTA)
Zhang et al. (2013a)	Operator-Based and Random Mutant Selection: Better Together	IEEE	International Conference on Automated Software Engineering (ASE)
Zhang et al. (2016)	Predictive Mutation Testing	ACM	International Symposium on Software Testing and Analysis (ISSTA)
Zhu et al. (2017)	Speeding-Up Mutation Testing via Data Compression and State Infection	IEEE	International Workshop on Mutation Analysis (Mutation)
Zhu et al. (2018)	An Investigation of Compression Techniques to Speed up Mutation Testing	IEEE	International Conference on Software Testing, Verification and Validation (ICST)

**Table A.2**

Subsumed studies.

Updated / Extended Studies		Subsumed studies	
Author/Year	Title	Author/Year	Title
Krauser et al. (1991)	High Performance Software Testing on SIMD Machines	Krauser et al. (1988)	High Performance Testing on SIMD Machines
Offutt and Lee (1994)	An Empirical Evaluation of Weak Mutation	Offutt and Lee (1991)	How Strong is Weak Mutation?
Untch et al. (1993)	Mutation Analysis Using Mutant Schemata	Untch (1992)	Mutation-based Software Testing Using Program Schemata
Wong and Mathur (1995)	Reducing the Cost of Mutation Testing: An Empirical Study	Mathur and Wong (1993)	Evaluation of the Cost of Alternative Mutation Testing Strategies
Wong and Mathur (1995)	Reducing the Cost of Mutation Testing: An Empirical Study	Wong et al. (1995)	Mutation Versus All-uses: An Empirical Evaluation of Cost, Strength and Effectiveness
Offutt and Pan (1997)	Automatically Detecting Equivalent Mutants and Infeasible Paths	Offutt and Pan (1996)	Detecting Equivalent Mutants and the Feasible Path Problem
Alexander et al. (2002)	Mutation of Java Objects	Bieman et al. (2001)	A Technique for Mutation of Java Objects
Bogacki and Walter (2006a)	Aspect-oriented Response Injection: an Alternative to Classical Mutation Testing	Bogacki and Walter (2006b)	Evaluation of Test Code Quality with Aspect-Oriented Mutations
Dominguez-Jiménez et al. (2009a)	GAmera: An Automatic Mutant Generation System for WS-BPEL Compositions	Dominguez-Jiménez et al. (2009b)	A Framework for Mutant Genetic Generation for WS-BPEL
Papadakis and Malevris (2011b)	Automatically Performing Weak Mutation with the Aid of Symbolic Execution, Concolic Testing and Search-based Testing	Papadakis et al. (2010)	Towards Automating the Generation of Mutation Tests
Fraser and Zeller (2012)	Mutation-Driven Generation of Oracles and Unit Tests	Fraser and Zeller (2010)	Mutation-Driven Generation of Oracles and Unit Tests
Schuler and Zeller (2013)	Covering and Uncovering Equivalent Mutants	Schuler and Zeller (2010)	(Un-)Covering Equivalent Mutants
Kaminski et al. (2013)	Improving logic-based testing	Kaminski et al. (2011a)	Better Predicate Testing
Reales and Polo (2014)	Reducing Mutation Costs Through Uncovered Mutants	Reales and Polo (2012)	Do Redundant Mutants Affect the Effectiveness and Efficiency of Mutation Analysis?
Just and Schweiggert (2015)	Higher Accuracy and Lower Run Time: Efficient Mutation Analysis Using Non-redundant Mutation Operators	Just et al. (2012a)	Mutant Execution Cost Reduction: Through MUSIC (Mutant Schema Improved with Extra Code)
Aichernig et al. (2013)	Incremental Refinement Checking for Test Case Generation	Aichernig and Jöbstl (2012)	Efficient Refinement Checking for Model-Based Mutation Testing
Kintis et al. (2015)	Employing Second-order Mutation for Isolating First-order Equivalent Mutants	Kintis et al. (2012)	Isolating First Order Equivalent Mutants via Second Order Mutation
Just et al. (2014a)	Efficient Mutation Analysis by Propagating and Partitioning Infected Execution States	Just (2014)	The Major Mutation Framework: Efficient and Scalable Mutation Analysis for Java

(continued on next page)

**Table A.2** (continued)

Updated / Extended Studies		Subsumed studies	
Author/Year	Title	Author/Year	Title
Kintis et al. (2017)	Detecting Trivial Mutant Equivalences via Compiler Optimisations	Papadakis et al. (2015)	Trivial Compiler Equivalence: A Large Scale Empirical Study of a Simple, Fast and Effective Equivalent Mutant Detection Technique
Gopinath et al. (2017)	Mutation Reduction Strategies Considered Harmful	Gopinath et al. (2016a)	On the Limits of Mutation Reduction Strategies
Petrović and Ivanković (2018)	State of Mutation Testing at Google	Petrović et al. (2018)	An Industrial Application of Mutation Testing: Lessons, Challenges, and Research Directions
McMinn et al. (2018)	Automatic Detection and Removal of Ineffective Mutants for the Mutation Analysis of Relational Database Schemas	Wright et al. (2014)	The Impact of Equivalent, Redundant and Quasi Mutants on Database Schema Mutation Analysis

## References

- Abuljadayel, A., Wedyan, F., 2018. An approach for the generation of higher order mutants using genetic algorithms. *Int. J. Intell. Syst. Appl.* 10 (1), 34–45.
- Acree, A.T., Budd, T.A., DeMillo, R.A., Lipton, R.J., Sayward, F.G., 1979. Mutation Analysis. Technical Report GIT-ICS-79/08. School of Information and Computer Science, Georgia Institute of Technology, Atlanta, GA, USA.
- Adamopoulos, K., Harman, M., Hierons, R.M., 2004. How to overcome the equivalent mutant problem and achieve tailored selective mutation using co-evolution. In: Proceedings of the 3rd Genetic and Evolutionary Computation Conference (GECCO). Springer, Seattle, WA, USA. II-1338–1349 (LNCS v.3103).
- Aichernig, B.K., Auer, J., Jöbstl, E., Korošec, R., Krenn, W., Schlick, R., Schmidt, B.V., 2014. Model-based mutation testing of an industrial measurement device. In: Proceedings of the 8th International Conference Tests and Proofs (TAP). Springer, York, UK, pp. 1–19.
- Aichernig, B.K., Jöbstl, E., 2012. Efficient refinement checking for model-based mutation testing. In: Proceedings of the 12th International Conference on Quality Software (QSIC). IEEE Computer Society, Xi'an, China, pp. 21–30.
- Aichernig, B.K., Jöbstl, E., Kegele, M., 2013. Incremental refinement checking for test case generation. In: Proceedings of the 7th International Conference Tests and Proofs (TAP). Springer, Budapest, Hungary, pp. 1–19.
- Al-Hajjaji, M., Krüger, J., Benduhn, F., Leich, T., Saake, G., 2017. Efficient mutation testing in configurable systems. In: Proceedings of the IEEE/ACM 2nd International Workshop on Variability and Complexity in Software Design (VACE). IEEE, Buenos Aires, Argentina, pp. 2–8.
- Alexander, R.T., Bieman, J.M., Ghosh, S., Ji, B., 2002. Mutation of Java objects. In: Proceedings of the 13th International Symposium on Software Reliability Engineering (ISSRE). IEEE, Annapolis, MD, USA, pp. 341–351.
- Ammann, P., Delamaro, M.E., Offutt, A.J., 2014. Establishing theoretical minimal sets of mutants. In: Proceedings of the 7th Conference on Software Testing, Verification and Validation (ICST). IEEE, Cleveland, OH, USA, pp. 21–30.
- Anbalagan, P., Xie, T., 2008. Automated generation of pointcut mutants for testing pointcuts in AspectJ programs. In: Proceeding of the 19th International Symposium on Software Reliability Engineering (ISSRE). IEEE, pp. 239–248.
- Andrews, J.H., Briand, L.C., Labiche, Y., 2005. Is mutation an appropriate tool for testing experiments? In: Proceedings of the 27th International Conference on Software Engineering (ICSE). ACM, St. Louis, MO, USA, pp. 402–411.
- Ayari, K., Bouktif, S., Antoniou, G., 2007. Automatic mutation test input data generation via ant colony. In: Proceedings of the 9th Genetic and Evolutionary Computation Conference (GECCO). ACM, London, UK, pp. 1074–1081.
- Barbosa, E.F., Maldonado, J.C., Vincenzi, A.M.R., 2001. Toward the determination of sufficient mutant operators for C. *Softw. Test. Verif. Reliabil.* 11 (2), 113–136.
- Bashir, M.B., Nadeem, A., 2017. Improved genetic algorithm to reduce mutation testing cost. *IEEE Access* 5, 3657–3674.
- Baudry, B., Fleurey, F., Jézéquel, J.-M., Le Traon, Y., 2005. From genetic to bacteriological algorithms for mutation-based testing. *Softw. Test. Verif. Reliabil.* 15 (2), 73–96.
- Belli, F., Beyazit, M., 2015. Exploiting model morphology for event-based testing. *IEEE Trans. Softw. Eng.* 41 (2), 113–134.
- Bieman, J.M., Ghosh, S., Alexander, R.T., 2001. A technique for mutation of Java objects. In: Proceedings of the 16th International Conference on Automated Software Engineering (ASE) - Short Paper. IEEE, San Diego, CA, USA, pp. 337–340.
- Biolchini, J., Mian, P.G., Natali, A.C.C., Travassos, G.H., 2005. Systematic Review in Software Engineering. Tech. Report RT-ES 679/05. Systems Engineering and Computer Science Dept., COPPE/UFRJ, Rio de Janeiro, RJ, Brazil.
- Bluemke, I., Kulesza, K., 2013. Reduction of computational cost in mutation testing by sampling mutants. In: Proceedings of the 8th International Conference on Dependability and Complex Systems (DepCoS-RELCOMEX). Springer, Brunów, Poland, pp. 41–51.
- Bluemke, I., Kulesza, K., 2014a. Reduction in mutation testing of Java classes. In: Proceedings of the 9th International Conference on Software Engineering and Applications (ICSOFT-EA). IEEE, Vienna, Austria, pp. 297–304.
- Bluemke, I., Kulesza, K., 2014b. Reductions of operators in Java mutation testing. In: Proceedings of the 9th International Conference on Dependability and Complex Systems (DepCoS-RELCOMEX). Springer, Brunów, Poland, pp. 93–102.
- Bogacki, B., Walter, B., 2006. Aspect-oriented response injection: an alternative to classical mutation testing. In: Proceedings of the IFIP Working Conference on Software Engineering Techniques: Design for Quality (SET). Springer, Warsaw, Poland, pp. 273–282.
- Bogacki, B., Walter, B., 2006. Evaluation of test code quality with aspect-oriented mutations. In: Proceedings of the 7th International Conference Extreme Programming and Agile Processes in Software Engineering (XP) - Posters and Demonstrations. Springer, Oulu, Finland, pp. 202–204.
- Budd, T., Angluin, D., 1982. Two notions of correctness and their relation to testing. *Acta Inform.* 8 (1), 31–45.
- Cachia, M.A., Micallef, M., Colombo, C., 2013. Towards incremental mutation testing. In: Proceedings of the 2013 Validation Strategies for Software Evolution Workshop (VSSE). Elsevier, Rome, Italy, pp. 2–11.
- Chen, L., Zhang, L., 2018. Speeding up mutation testing via regression test selection: an extensive study. In: Proceedings of the 11th International Conference on Software Testing, Verification and Validation (ICST). IEEE, Västerås, Sweden, pp. 58–69.
- Choi, B., Mathur, A.P., 1993. High-performance mutation testing. *J. Syst. Softw.* 20 (2), 135–152.
- Delamaro, M.E., Deng, L., Durelli, V.H.S., Li, N., Offutt, A.J., 2014. Experimental evaluation of SDL and One-Op mutation for C. In: Proceedings of the 7th International Conference on Software Testing, Verification and Validation (ICST). IEEE, Cleveland, OH, USA, pp. 203–212.
- Delamaro, M.E., Deng, L., Li, N., Durelli, V.H.S., Offutt, A.J., 2014. Growing a reduced set of mutation operators. In: Proceedings of the 28th Brazilian Symposium on Software Engineering (SBES). IEEE, Maceió, AL, Brazil, pp. 81–90.
- Delamaro, M.E., Offutt, A.J., Ammann, P., 2014. Designing deletion mutation operators. In: Proceedings of the 7th International Conference on Software Testing, Verification and Validation (ICST). IEEE, Cleveland, OH, USA, pp. 11–20.
- Delgado-Pérez, P., Medina-Bulo, I., Merayo, M.G., 2017. Using evolutionary computation to improve mutation testing. In: Proceedings of the 14th International Work-Conference on Artificial Neural Networks (IWANN). Springer, Cadiz, Spain, pp. 381–391.
- Delgado-Pérez, P., Medina-Bulo, I., Palomo-Lozano, F., García-Domínguez, A., Domínguez-Jiménez, J.J., 2017. Assessment of class mutation operators for C++ with the MuCPP mutation system. *Inf. Softw. Technol.* 169–184.
- Delgado-Pérez, P., Medina-Bulo, I., Segura, S., García-Domínguez, A., Juan, J., 2017. GiGAn: evolutionary mutation testing for C++ object-oriented systems. In: Proceedings of the 32nd ACM Symposium on Applied Computing (SAC). ACM, Marrakech, Morocco, pp. 1387–1392.
- Delgado-Pérez, P., Segura, S., Medina-Bulo, I., 2017. Assessment of C++ object-oriented mutation operators: a selective mutation approach. *Softw. Test. Verif. Reliabil.* 27 (4–5), 1630–1649.
- DeMillo, R.A., Krauser, E.W., Mathur, A.P., 1991. Compiler-integrated program mutation. In: Proceedings of the 15th IEEE Annual Computer Software and Applications Conference (COMPSAC). IEEE, Tokyo, Japan, pp. 351–356.
- DeMillo, R.A., Lipton, R.J., Sayward, F.G., 1978. Hints on test data selection: help for the practicing programmer. *IEEE Comput.* 11 (4), 34–43.
- DeMillo, R.A., Offutt, A.J., 1991. Constraint-based automatic test data generation. *IEEE Trans. Softw. Eng.* 17 (9), 900–910.
- DeMillo, R.A., Offutt, A.J., 1993. Experimental results from an automatic test case generator. *ACM Trans. Softw. Eng. Methodol.* 2 (2), 109–127.
- Deng, L., Offutt, A.J., Li, N., 2013. Empirical evaluation of the statement deletion mutation operator. In: Proceedings of the 6th International Conference on Software Testing, Verification and Validation (ICST). IEEE, Luxembourg City, Luxembourg, pp. 84–93.
- Denisov, A., Pankevich, S., 2018. Mull it over: mutation testing based on LLVM. In: Proceedings of the 13th International Workshop on Mutation Analysis (Mutation). IEEE, Västerås, Sweden, pp. 25–31.
- Derezińska, A., 2013. A quality estimation of mutation clustering in C# programs. In: Proceedings of the 8th International Conference on Dependability and Complex Systems (DepCoS-RELCOMEX). Springer, Brunów, Poland, pp. 119–129(AISC v.224).
- Derezińska, A., 2016. Evaluation of deletion mutation operators in mutation testing of C# programs. In: Proceedings of the 11th International Conference on

- Dependability and Complex Systems (DepCoS-RELCOMEX). Springer, Brunnów, Poland, pp. 97–108.
- Derezińska, A., Hałas, K., 2015. Improving mutation testing process of python programs. In: Proceedings of the 4th Computer Science On-line Conference Software Engineering in Intelligent Systems (CSOC). Springer, pp. 233–242.
- Derezińska, A., Rudnik, M., 2012. Quality evaluation of object-oriented and standard mutation operators applied to C# programs. In: Proceedings of the 50th International Conference on Modelling Techniques and Tools for Computer Performance Evaluation (TOOLS). Springer, Prague, Czech Republic, pp. 42–57.
- Derezińska, A., Rudnik, M., 2017. Evaluation of mutant sampling criteria in object-oriented mutation testing. In: Proceedings of the Federated Conference on Computer Science and Information Systems (ACSIS). Polskie Towarzystwo Informatyczne, Prague, Czech Republic, pp. 1315–1324.
- Devroey, X., Perrouin, G., Papadakis, M., Legay, A., Schobbens, P.-Y., Heymans, P., 2016. Featured model-based mutation analysis. In: Proceedings of the 38th International Conference on Software Engineering (ICSE). ACM, Austin, TX, USA, pp. 655–666.
- Devroey, X., Perrouin, G., Papadakis, M., Legay, A., Schobbens, P.Y., Heymans, P., 2017. Automata language equivalence vs. simulations for model-based mutant equivalence: an empirical evaluation. In: Proceedings of the 10th IEEE International Conference on Software Testing, Verification and Validation (ICST). IEEE, Tokyo, Japan, pp. 424–429.
- Do, H., Rothermel, G., 2006. On the use of mutation faults in empirical assessments of test case prioritization techniques. *IEEE Trans. Softw. Eng.* 32 (9), 733–752.
- Dominguez-Jiménez, J.J., Estero-Botaro, A., García-Domínguez, A., Medina-Bulo, I., 2009. GAmara: an automatic mutant generation system for WS-BPEL compositions. In: Proceedings of the 7th IEEE European Conference on Web Services (ECOWS). IEEE, Eindhoven, The Netherlands, pp. 97–106.
- Dominguez-Jiménez, J.J., Estero-Botaro, A., García-Domínguez, A., Medina-Bulo, I., 2011. Evolutionary mutation testing. *Inf. Softw. Technol.* 53, 1108–1123.
- Domínguez-Jiménez, J.J., Estero-Botaro, A., Medina-Bulo, I., 2009. A framework for mutant genetic generation for WS-BPEL. In: Proceedings of the 35th Conference on Current Trends in Theory and Practice of Computer Science (SOFSEM). IEEE, Spindler Mlyn, Czech Republic, pp. 229–240.
- Durrelli, V.H.S., Offutt, A.J., Delamaro, M.E., 2012. Toward harnessing high-level language virtual machines for further speeding up weak mutation testing. In: Proceedings of the 5th International Conference on Software Testing, Verification and Validation (ICST). IEEE, Montreal, QC, Canada, pp. 681–690.
- Fernandes, L., Ribeiro, M., Carvalho, L., Gheyi, R., Mongiovì, M., Santos, A., Cavalcanti, A., Ferrari, F.C., Maldonado, J.C., 2017. Avoiding useless mutants. In: Proceedings of the 16th ACM International Conference on Generative Programming: Concepts and Experiences (GPCE). ACM, Vancouver, BC, Canada, pp. 187–198.
- Ferrari, F.C., Pizzoleto, A.V., Offutt, A.J., 2018. A systematic review of cost reduction techniques for mutation testing: preliminary results. In: Proceedings of the 13th International Workshop on Mutation Analysis (Mutation). IEEE, Västerås, Sweden, pp. 1–10.
- Ferrari, F.C., Pizzoleto, A.V., Offutt, A.J., Fernandes, L., Ribeiro, M., 2018. SLR on Cost Reduction for Mutation Testing – Companion Website. <http://goo.gl/edyFin>.
- Ferrari, F.C., Rashid, A., Maldonado, J.C., 2013. Towards the practical mutation testing of AspectJ programs. *Sci. Comput. Programm.* 78 (9), 1639–1662.
- Fleyshgakker, V.N., Weiss, S.N., 1994. Efficient mutation analysis: a new approach. In: Proceedings of the International Symposium on Software Testing and Analysis (ISSTA). ACM, Seattle, WA, USA, pp. 185–195.
- Frankl, P.G., Weiss, S.N., Hu, C., 1997. All-uses vs mutation testing: an experimental comparison of effectiveness. *J. Syst. Softw.* 38 (3), 235–253.
- Fraser, G., Arcuri, A., 2015. Achieving scalable mutation-based generation of whole test suites. *Empir. Softw. Eng.* 20 (3), 783–812.
- Fraser, G., Zeller, A., 2010. Mutation-driven generation of oracles and unit tests. In: Proceedings of the 19th International Symposium on Software Testing and Analysis (ISSTA). ACM, Trento, Italy, pp. 147–158.
- Fraser, G., Zeller, A., 2012. Mutation-driven generation of oracles and unit tests. *IEEE Trans. Softw. Eng.* 38 (2), 278–292.
- Gligoric, M., Jagannath, V., Luo, Q., Marinov, D., 2012. Efficient mutation testing of multithreaded code. *Softw. Test. Verif. Reliabil.* 23 (5), 375–403.
- Gligoric, M., Jagannath, V., Marinov, D., 2010. MuTMuT: efficient exploration for mutation testing of multithreaded code. In: Proceedings of the 3th International Conference on Software Testing, Verification and Validation (ICST). IEEE, Paris, France, pp. 55–64.
- Gligoric, M., Zhang, L., Pereira, C., Pokam, G., 2013. Selective mutation testing for concurrent code. In: Proceedings of the 22nd International Symposium on Software Testing and Analysis (ISSTA). ACM, Lugano, Switzerland, pp. 224–234.
- Gopinath, R., Ahmed, I., Alipour, M.A., Jensen, C., Groce, A., 2017. Mutation reduction strategies considered harmful. *IEEE Trans. Softw. Reliabil.* 66 (3), 854–874.
- Gopinath, R., Alipour, M.A., Ahmed, I., Jensen, C., Groce, A., 2016. On the limits of mutation reduction strategies. In: Proceedings of the 38th International Conference on Software Engineering (ICSE). ACM, Austin, TX, USA, pp. 511–522.
- Gopinath, R., Jensen, C., Groce, A., 2016. Topsy-turvy: a smarter and faster parallelization of mutation analysis. In: Proceedings of the 38th International Conference on Software Engineering Companion (ICSE). ACM, Austin, TX, USA, pp. 740–743.
- Gopinath, R., Mathis, B., Zeller, A., 2018. If you can't kill a supermutant, you have a problem. In: Proceedings of the 13th International Workshop on Mutation Analysis (Mutation). IEEE, Västerås, Sweden, pp. 18–24.
- Hamlet, R.G., 1977. Testing programs with the aid of a compiler. *IEEE Trans. Softw. Eng.* 3 (4), 279–290.
- Harman, M., Hierons, R.M., Danicic, S., 2000. The relationship between program dependence and mutation analysis. In: Proceedings of the Mutation 2000 Symposium. Kluwer Academic Publishers, San Jose, CA, USA, pp. 5–13.
- Harman, M., Jia, Y., Reales, P., Polo, M., 2014. Angels and monsters: an empirical investigation of potential test effectiveness and efficiency improvement from strongly subsuming higher order mutation. In: Proceedings of the 29th International Conference on Automated Software Engineering (ASE). ACM, Västerås, Sweden, pp. 397–408.
- Henard, C., Papadakis, M., Le Traon, Y., 2014. Mutation-based generation of software product line test configurations. In: Proceedings of the 6th International Symposium Search-Based Software Engineering (SSBSE). Springer, Fortaleza, CE, Brazil, pp. 92–106.
- Hernandes, E.C.M., Zamboni, A.B., Fabbri, S.C.P.F., Di Thommazo, A., 2012. Using GQM and TAM to evaluate StArt - a tool that supports systematic review. *CLEI Electron. J.* 15 (1), 13–25.
- Hierons, R.M., Harman, M., Danicic, S., 1999. Using program slicing to assist in the detection of equivalent mutants. *Softw. Test. Verif. Reliabil.* 9 (4), 233–262.
- Holling, D., Banescu, S., Probst, M., Petrovska, A., Pretschner, A., 2016. Nequivack: assessing mutation score confidence. In: Proceedings of the 11th International Workshop on Mutation Analysis (Mutation). IEEE, Chicago, IL, USA, pp. 152–161.
- Howden, W.E., 1982. Weak mutation testing and completeness of test sets. *IEEE Trans. Softw. Eng.* 8 (4), 371–379.
- Hu, J., Li, N., Offutt, A.J., 2011. An analysis of OO mutation operators. In: Proceedings of the 6th International Workshop on Mutation Analysis (Mutation). IEEE, Berlin, Germany, pp. 334–341.
- Iida, C., Takada, S., 2017. Reducing mutants with mutant killable precondition. In: Proceedings of the 12th International Workshop on Mutation Analysis (Mutation). IEEE, Tokyo, Japan, pp. 128–133.
- Inozemtseva, L., Hemmati, H., Holmes, R., 2013. Using fault history to improve mutation reduction. In: Proceedings of the 9th Joint Meeting on Foundations of Software Engineering: New Ideas Track (ESEC/FSE). ACM, Saint Petersburg, Russia, pp. 639–642.
- Jackson, D., Woodward, M.R., 2000. Parallel firm mutation of Java programs. In: Proceedings of the Mutation 2000 Symposium. Kluwer Academic Publishers, San Jose, CA, USA, pp. 55–61.
- Ji, C., Chen, Z., Xu, B., Zhao, Z., 2009. A novel method of mutation clustering based on domain analysis. In: Proceedings of the 21th International Conference on Software Engineering and Knowledge Engineering (SEKE). Knowledge Systems Institute Graduate School, Boston, MA, USA, pp. 1–6.
- Jia, Y., Harman, M., 2011. An analysis and survey of the development of mutation testing. *IEEE Trans. Softw. Eng.* 37 (5), 649–678.
- Just, R., 2014. The major mutation framework: efficient and scalable mutation analysis for Java. In: Proceedings of the 23rd International Symposium on Software Testing and Analysis (ISSTA). ACM, San Jose, CA, USA, pp. 433–436.
- Just, R., Ernst, M.D., Fraser, G., 2014. Efficient mutation analysis by propagating and partitioning infected execution states. In: Proceedings of the 23rd International Symposium on Software Testing and Analysis (ISSTA). ACM, San Jose, CA, USA, pp. 315–326.
- Just, R., Jalali, D., Inozemtseva, L., Ernst, M.D., Holmes, R., Fraser, G., 2014. Are mutants a valid substitute for real faults in software testing? In: Proceedings of the 22nd International Symposium on Foundations of Software Engineering (FSE). ACM, Hong Kong, China, pp. 654–665.
- Just, R., Kapfhammer, G.M., Schweiggert, F., 2011. Using conditional mutation to increase the efficiency of mutation analysis. In: Proceedings of the 6th International Workshop on Automation of Software Test (AST). ACM, Waikiki, Honolulu, HI, USA, pp. 50–56.
- Just, R., Kapfhammer, G.M., Schweiggert, F., 2012. Do redundant mutants affect the effectiveness and efficiency of mutation analysis? In: Proceedings of the 5th IEEE International Conference on Software Testing, Verification and Validation (ICST). IEEE, Montreal, QC, Canada, pp. 720–725.
- Just, R., Kapfhammer, G.M., Schweiggert, F., 2012. Using non-redundant mutation operators and test suite prioritization to achieve efficient and scalable mutation analysis. In: Proceedings of the IEEE 23th International Symposium on Software Reliability Engineering (ISSRE). IEEE, Dallas, TX, USA, pp. 11–20.
- Just, R., Kurtz, B., Ammann, P., 2017. Inferring mutant utility from program context. In: Proceedings of the 26th International Symposium on Software Testing and Analysis (ISSTA). ACM, Santa Barbara, CA, USA, pp. 284–294.
- Just, R., Schweiggert, F., 2015. Higher accuracy and lower run time: efficient mutation analysis using non-redundant mutation operators. *Softw. Test. Verif. Reliabil.* 25 (5–7), 490–507.
- Kaminski, G., Ammann, P., 2009. Using a fault hierarchy to improve the efficiency of DNF logic mutation testing. In: Proceedings of the 2nd International Conference on Software Testing, Verification and Validation (ICST). IEEE, Denver, CO, USA, pp. 386–395.
- Kaminski, G., Ammann, P., 2011. Reducing logic test set size while preserving fault detection. *Softw. Test. Verif. Reliabil.* 21 (3), 155–193.
- Kaminski, G., Ammann, P., Offutt, A.J., 2011. Better predicate testing. In: Proceedings of the 6th International Workshop on Automation of Software Test (AST). ACM, Honolulu, HI, USA, pp. 57–63.
- Kaminski, G., Ammann, P., Offutt, A.J., 2013. Improving logic-based testing. *J. Syst. Softw.* 86 (8), 2002–2012.
- Kaminski, G., Phraphamontripang, U., Ammann, P., Offutt, A.J., 2011. A logic mutation approach to selective mutation for programs and queries. *Inf. Softw. Technol.* 53 (10), 1137–1152.

- Kim, S.-W., Ma, Y.-S., Kwon, Y.-R., 2012. Combining weak and strong mutation for a noninterpretive Java mutation system. *Softw. Test. Verif. Reliabil.* 23 (8), 647–668.
- Kintis, M., Malevris, N., 2013. Identifying more equivalent mutants via code similarity. In: Proceedings of the 20th Asia-Pacific Software Engineering Conference (APSEC). IEEE, Bangkok, Thailand, pp. 180–188.
- Kintis, M., Malevris, N., 2014. Using data flow patterns for equivalent mutant detection. In: Proceedings of the 9th International Workshop on Mutation Analysis (Mutation). IEEE, Cleveland, OH, USA, pp. 196–205.
- Kintis, M., Malevris, N., 2015. MEDIC: a static analysis framework for equivalent mutant identification. *Inf. Softw. Technol.* 68, 1–17.
- Kintis, M., Papadakis, M., Jia, Y., Malevris, N., Le Traon, Y., Harman, M., 2017. Detecting trivial mutant equivalences via compiler optimisations. *IEEE Trans. Softw. Eng.* 44 (4), 1–25.
- Kintis, M., Papadakis, M., Malevris, N., 2010. Evaluating mutation testing alternatives: a collateral experiment. In: Proceedings of the 17th Asia-Pacific Software Engineering Conference (APSEC). IEEE, Sydney, Australia, pp. 300–309.
- Kintis, M., Papadakis, M., Malevris, N., 2012. Isolating first order equivalent mutants via second order mutation. In: Proceedings of the 5th IEEE International Conference on Software Testing, Verification and Validation (ICST). IEEE, Montreal, QC, Canada, pp. 701–710.
- Kintis, M., Papadakis, M., Malevris, N., 2015. Employing second-order mutation for isolating first-order equivalent mutants. *Softw. Test. Verif. Reliabil.* 25 (5–7), 508–535.
- Kitchenham, B., 2004. Procedures for Performing Systematic Reviews. Joint Technical Report TR/SE-0401 (Keele) - 040001T.1 (NICTA). Software Engineering Group - Department of Computer Science - Keele University; and Empirical Software Engineering - National ICT Australia Ltd, Staffordshire, UK; and Eversleigh, Australia.
- Kitchenham, B.A., Dybå, T., Jørgensen, M., 2004. Evidence-based software engineering. In: Proceedings of the 26th International Conference on Software Engineering (ICSE). IEEE, Edinburgh, Scotland, pp. 273–281.
- Krauser, E.W., Mathur, A.P., Rego, V.J., 1988. High performance testing on SIMD machines. In: Proceedings of the 2nd Workshop on Software Testing, Verification, and Analysis. IEEE, Banff, AB, Canada, pp. 171–177.
- Krauser, E.W., Mathur, A.P., Rego, V.J., 1991. High performance software testing on SIMD machines. *IEEE Trans. Softw. Eng.* 17 (5), 403–423.
- Kurtz, B., Ammann, P., Offutt, A.J., Delamaro, M.E., Kurtz, M., Gökçe, N., 2016. Analyzing the validity of selective mutation with dominator mutants. In: Proceedings of the 24th International Symposium on Foundations of Software Engineering (FSE). ACM, Seattle, WA, USA, pp. 571–582.
- Lacerda, J.T.S., Ferrari, F.C., 2014. Towards the establishment of a sufficient set of mutation operators for AspectJ programs. In: Proceedings of the 8th Brazilian Workshop on Systematic and Automated Software Testing (SAST). Brazilian Computer Society, Maceio, AL, Brazil, pp. 21–30.
- Li, N., Praphamontripong, U., Offutt, A.J., 2009. An experimental comparison of four unit test criteria: mutation, edge-pair, all-uses and prime path coverage. In: Proceedings of the 4th International Workshop on Mutation Analysis (Mutation). IEEE, Denver, CO, USA, pp. 220–229.
- Li, N., West, M., Escalona, A., Durelli, V.H.S., 2015. Mutation testing in practice using ruby. In: Proceedings of the 10th International Workshop on Mutation Analysis (Mutation). IEEE, Graz, Austria, pp. 1–6.
- Lima, J.A.P., Guizzo, G., Vergilio, S.R., Silva, A.P.C., Filho, H.L.J., Ehrenfried, H.V., 2016. Evaluating different strategies for reduction of mutation testing costs. In: Proceedings of the 1st Brazilian Symposium on Systematic and Automated Software Testing (SAST). ACM, Maringá, PR, Brazil, pp. 4:1–4:10.
- Liu, M.-H., Gao, Y.-F., Shan, J.-H., Liu, J.-H., Zhang, L., Sun, J.-S., 2006. An approach to test data generation for killing multiple mutants. In: Proceedings of the 22th International Conference on Software Maintenance (ICSM). IEEE, Philadelphia, PA, USA, pp. 113–122.
- Ma, Y.-S., Kim, S.-W., 2016. Mutation testing cost reduction by clustering overlapped mutants. *J. Syst. Softw.* 115 (C), 18–30.
- Ma, Y.-S., Offutt, A.J., Kwon, Y.R., 2005. MuJava: an automated class mutation system. *Softw. Test. Verif. Reliabil.* 15 (2), 97–133.
- MacDonell, S., Shepperd, M., Kitchenham, B., Mendes, E., 2010. How reliable are systematic reviews in empirical software engineering? *IEEE Trans. Softw. Eng.* 36 (5), 676–687.
- Madeyski, L., Orzeszyna, W., Torkar, R., Józala, M., 2014. Overcoming the equivalent mutant problem: a systematic literature review and a comparative experiment of second order mutation. *IEEE Trans. Softw. Eng.* 40 (1), 23–42.
- Marcozzi, M., Bardin, S., Kosmatov, N., Papadakis, M., Prevost, V., Correnson, L., 2018. Time to clean your test objectives. In: Proceedings of the 40th International Conference on Software Engineering (ICSE). IEEE, Gothenburg, Sweden, pp. 456–467.
- Marshall, A.C., Hedley, D., Riddell, I.J., Hennell, M.A., 1990. Static dataflow-aided weak mutation analysis (SDAWM). *Inf. Softw. Technol.* 32 (1), 99–104.
- Mathur, A.P., 1991. Performance, effectiveness, and reliability issues in software testing. In: Proceedings of the 15th IEEE Annual Computer Software and Applications Conference (COMPSAC). IEEE, Tokyo, Japan, pp. 604–605.
- Mathur, A.P., 2007. Foundations of Software Testing, 1st Addison-Wesley Professional, Toronto, ON, Canada.
- Mathur, A.P., Wong, W.E., 1993. Evaluation of the cost of alternative mutation testing strategies. In: Proceedings of the 7th Brazilian Symposium on Software Engineering (SBES). Brazilian Computer Society, João Pessoa, PB, Brazil, pp. 320–335.
- Mathur, A.P., Wong, W.E., 1994. An empirical comparison of data flow and mutation based test adequacy criteria. *Softw. Test. Verif. Reliabil.* 4 (1), 9–31.
- Matnei Filho, R.A., Vergilio, S.R., 2015. A mutation and multi-objective test data generation approach for feature testing of software product lines. In: Proceedings of the 29th Brazilian Symposium on Software Engineering (SBES). IEEE, Belo Horizonte, MG, Brazil, pp. 21–30.
- McMinn, P., Kapfhammer, G.M., Wright, C.J., 2016. Virtual mutation analysis of relational database schemas. In: Proceedings of the 11th International Workshop on Automation of Software Test (AST). ACM, Austin, TX, USA, pp. 36–42.
- McMinn, P., Wright, C.J., McCurdy, C.J., Kapfhammer, G., 2018. Automatic detection and removal of ineffective mutants for the mutation analysis of relational database schemas (in press). *IEEE Trans. Softw. Eng.* 1–1.
- Mresa, E.S., Bottaci, L., 1999. Efficiency of mutation operators and selective mutation strategies: an empirical study. *Softw. Test. Verif. Reliabil.* 9 (4), 205–232.
- Nobre, T., Vergilio, S.R., Pozo, A., 2012. Reducing interface mutation costs with multiobjective optimization algorithms. *Int. J. Nat. Comput. Res.* 21–40.
- Offutt, A.J., Craft, W.M., 1994. Using compiler optimization techniques to detect equivalent mutants. *Softw. Test. Verif. Reliabil.* 4 (3), 131–154.
- Offutt, A.J., Jin, Z., Pan, J., 1999. The dynamic domain reduction approach to test data generation. *Software Pract. Exp.* 29 (2), 167–193.
- Offutt, A.J., Lee, A., Rothermel, G., Untch, R.H., Zapf, C., 1996. An experimental determination of sufficient mutant operators. *ACM Trans. Softw. Eng. Methodol.* 5 (2), 99–118.
- Offutt, A.J., Lee, S.D., 1991. How strong is weak mutation? In: Proceedings of the International Symposium on Software Testing and Analysis (ISSTA). ACM, Victoria, BC, Canada, pp. 200–213.
- Offutt, A.J., Lee, S.D., 1994. An empirical evaluation of weak mutation. *IEEE Trans. Softw. Eng.* 20 (5), 337–344.
- Offutt, A.J., Ma, Y.-S., Kwon, Y.-R., 2006. The class-level mutants of MuJava. In: Proceedings of the International Workshop on Automation of Software Test (AST). ACM, Shanghai, China, pp. 78–84.
- Offutt, A.J., Pan, J., 1996. Detecting equivalent mutants and the feasible path problem. In: Proceedings of the 11th Annual Conference on Computer Assurance (COMPASS). ACM, Gaithersburg, MD, USA, pp. 224–236.
- Offutt, A.J., Pan, J., 1997. Automatically detecting equivalent mutants and infeasible paths. *Softw. Test. Verif. Reliabil.* 7 (3), 165–192.
- Offutt, A.J., Pargas, R.P., Fichter, S.V., Khambekar, P.K., 1992. Mutation testing of software using a MIMD computer. In: Proceedings of the International Conference on Parallel Processing (ICPP). John Wiley & Sons, Chicago, Illinois, USA.
- Offutt, A.J., Rothermel, G., Zapf, C., 1993. An experimental evaluation of selective mutation. In: Proceedings of the 15th International Conference on Software Engineering (ICSE). IEEE, Baltimore, MD, USA, pp. 100–107.
- Offutt, A.J., Untch, R.H., 2000. Mutation 2000: uniting the orthogonal. In: Proceedings of the Mutation 2000 Symposium. Kluwer Academic Publishers, San Jose, CA, USA, pp. 34–44.
- Oliveira, A.A.L., Camilo-Junior, C.G., Vincenzi, A.M.R., 2013. A coevolutionary algorithm to automatic test case selection and mutant in mutation testing. In: Proceedings of the 2013 IEEE Congress on Evolutionary Computation (CEC). IEEE, Cancun, Mexico, pp. 829–836.
- Omar, E., Ghosh, S., 2012. An exploratory study of higher order mutation testing in aspect-oriented programming. In: Proceedings of the 23th International Symposium on Software Reliability Engineering (ISSRE). IEEE, Dallas, TX, USA, pp. 1–10.
- Papadakis, M., Delamaro, M., Le Traon, Y., 2014. Mitigating the effects of equivalent mutants with mutant classification strategies. *Sci. Comput. Programm.* 298–319.
- Papadakis, M., Jia, Y., Harman, M., Le Traon, Y., 2015. Trivial compiler equivalence: a large scale empirical study of a simple, fast and effective equivalent mutant detection technique. In: Proceedings of the 37th International Conference on Software Engineering (ICSE). ACM, Florence, Italy, pp. 936–946.
- Papadakis, M., Delamaro, M., Zhang, J., Jia, Y., Le Traon, Y., Harman, M., 2019. Mutation testing advances: an analysis and survey. In: Memon, A.M. (Ed.), Advances in Computers, 112. Elsevier, Amsterdam, The Netherlands, pp. 275–378.
- Papadakis, M., Le Traon, Y., 2013. Mutation testing strategies using mutant classification. In: Proceedings of the 28th ACM Symposium on Applied Computing (SAC). ACM, Coimbra, Portugal, pp. 1223–1229.
- Papadakis, M., Malevris, N., 2009. An effective path selection strategy for mutation testing. In: Proceedings of the 16th Asia-Pacific Software Engineering Conference (APSEC). IEEE, Batu Ferringhi, Penang, Malaysia, pp. 422–429.
- Papadakis, M., Malevris, N., 2010. An empirical evaluation of the first and second order mutation testing strategies. In: Proceedings of the 5th International Workshop on Mutation Analysis (Mutation). IEEE, Paris, France, pp. 90–99.
- Papadakis, M., Malevris, N., 2010. Automatic mutation test case generation via dynamic symbolic execution. In: Proceedings of the IEEE 21th International Symposium on Software Reliability Engineering (ISSRE). IEEE, San Jose, CA, USA, pp. 121–130.
- Papadakis, M., Malevris, N., 2011. Automatic mutation based test data generation. In: Proceedings of the 13th Genetic and Evolutionary Computation Conference (GECCO): Search-Based Software Engineering Track (Poster Session). ACM, Dublin, Ireland, pp. 247–248.
- Papadakis, M., Malevris, N., 2011. Automatically performing weak mutation with the aid of symbolic execution, concolic testing and search-based testing. *Softw. Qual. J.* 19 (4), 691–723.
- Papadakis, M., Malevris, N., 2012. Mutation based test case generation via a path selection strategy. *Inf. Softw. Technol.* 54 (9), 915–932.
- Papadakis, M., Malevris, N., Kallia, M., 2010. Towards automating the generation of mutation tests. In: Proceedings of the 5th Workshop on Automation of Software Test (AST). ACM, Cape Town, South Africa, pp. 111–118.
- Parsai, A., Murgia, A., Demeyer, S., 2016. Evaluating random mutant selection at class-level in projects with non-adequate test suites. In: Proceedings of the 20th

- International Conference on Evaluation and Assessment in Software Engineering (EASE). ACM, Limerick, Ireland, pp. 1–10.
- Patel, K., Hierons, R.M., 2016. Resolving the equivalent mutant problem in the presence of non-determinism and coincidental correctness. In: Proceedings of the 28th International Conference on Testing Software and Systems (ICTSS). Springer, Graz, Austria, pp. 123–138(LNCS v.9976).
- Patrick, M., Oriol, M., Clark, J.A., 2012. MESSI: mutant evaluation by static semantic interpretation. In: Proceedings of the 5th International Conference on Software Testing, Verification and Validation (ICST). IEEE, Montreal, QC, Canada, pp. 711–719.
- Petrović, G., Ivanković, M., 2018. State of mutation testing at Google. In: Proceedings of the 40th International Conference on Software Engineering - Software Engineering in Practice Track (ICSE-SEIP) (to appear). IEEE, Gothenburg, Sweden, pp. 163–171.
- Petrović, G., Ivanković, M., Kurtz, B., Ammann, P., Just, R., 2018. An industrial application of mutation testing: lessons, challenges, and research directions. In: Proceedings of the 13th International Workshop on Mutation Analysis (Mutation). IEEE, Västerås, Sweden, pp. 47–53.
- Polo, M., Piattini, M., García-Rodríguez, I., 2009. Decreasing the cost of mutation testing with second-order mutants. *Softw. Test. Verif. Reliabil.* 19 (2), 111–131.
- Polo, M., Reales, P., 2010. Mutation testing cost reduction techniques: a survey. *IEEE Softw.* 27 (3), 80–86.
- Praphamontripong, U., Offutt, A.J., 2017. Finding redundancy in web mutation operators. In: Proceedings of the 12th International Workshop on Mutation Analysis (Mutation). IEEE, Tokyo, Japan, pp. 134–142.
- Quyen, N.T.H., Tung, K.T., Hanh, L.T.M., Binh, N.T., 2016. Improving mutant generation for simulink models using genetic algorithm. In: Proceedings of the International Conference on Electronics, Information, and Communications (ICEIC). IEEE, Da Nang, Vietnam, pp. 1–4.
- Reales, P., Polo, M., 2012. Mutant execution cost reduction: through MUSIC (Mutant schema improved with extra code). In: Proceedings of the 5th International Conference on Software Testing, Verification and Validation (ICST). IEEE, Montreal, QC, Canada, pp. 664–672.
- Reales, P., Polo, M., 2013. Parallel mutation testing. *Softw. Test. Verif. Reliabil.* 23 (4), 315–350.
- Reales, P., Polo, M., 2014. Reducing mutation costs through uncovered mutants. *Softw. Test. Verif. Reliabil.* 25 (5–7), 464–489.
- Reales, P., Polo, M., Alemán, J.L.F., 2013. Validating second-order mutation at system level. *IEEE Trans. Softw. Eng.* 39 (4), 570–587.
- Reuling, D., Bürdek, J., Rotärmel, S., Lochau, M., Kelter, U., 2015. Fault-based product-line testing: effective sample generation based on feature-diagram mutation. In: Proceedings of the 19th International Conference on Software Product Line (SPLC). ACM, Nashville, TN, USA, pp. 131–140.
- Sahinoğlu, M., Spafford, E.H., 1990. A Bayes sequential statistical procedure for approving software products. In: Proceedings of the IFIP Conference on Approving Software Products (ASP). Elsevier, Banff, AB, Canada, pp. 43–56.
- Schuler, D., Dallmeier, V., Zeller, A., 2009. Efficient mutation testing by checking invariant violations. In: Proceedings of the 18th International Symposium on Software Testing and Analysis (ISSTA). ACM, Chicago, IL, USA, pp. 69–80.
- Schuler, D., Zeller, A., 2010. (Un-)Covering equivalent mutants. In: Proceedings of the 3th International Conference on Software Testing, Verification and Validation (ICST). IEEE, Paris, France, pp. 45–54.
- Schuler, D., Zeller, A., 2013. Covering and uncovering equivalent mutants. *Softw. Test. Verif. Reliabil.* 23 (5), 353–374.
- Siami-Namin, A., Andrews, J.H., 2006. Finding sufficient mutation operators via variable reduction. In: Proceedings of the 2nd Workshop on Mutation Analysis (Mutation). IEEE, Raleigh, NC, USA, pp. 1–10.
- Siami-Namin, A., Andrews, J.H., Murdoch, D.J., 2008. Sufficient mutation operators for measuring test effectiveness. In: Proceedings of the 30th International Conference on Software Engineering (ICSE). ACM, Leipzig, Germany, pp. 351–360.
- Silva, R.A., Souza, S.R.S., Souza, P.S.L., 2017. A systematic review on search based mutation testing. *Inf. Softw. Technol.* 81, 19–35.
- Sridharan, M., Siami-Namin, A., 2010. Prioritizing mutation operators based on importance sampling. In: Proceedings of the IEEE 21th International Symposium on Software Reliability Engineering (ISSRE). IEEE, San Jose, CA, USA, pp. 378–387.
- Steimann, F., Thies, A., 2010. From behaviour preservation to behaviour modification: constraint-based mutant generation. In: Proceedings of the 32th International Conference on Software Engineering (ICSE). ACM, Cape Town, South Africa, pp. 425–434.
- Sun, C., Pan, L., Wang, Q., Liu, H., Zhang, X., 2017. An empirical study on mutation testing of WS-BPEL programs. *Comput. J.* 60 (1), 143–158.
- Sun, C., Xue, F., Liu, H., Zhang, X., 2017. A path-aware approach to mutant reduction in mutation testing. *Inf. Softw. Technol.* 81, 65–81.
- Tuya, J., Suárez-Cabal, M.J., de la Riva, C., 2007. Mutating database queries. *Inf. Softw. Technol.* 49 (4), 398–417.
- Untch, R.H., 1992. Mutation-based software testing using program schemata. In: Proceedings of the 30th Annual Southeast Regional Conference (ACM-SE). ACM, Raleigh, NC, USA, pp. 285–291.
- Untch, R.H., 2009. On reduced neighborhood mutation analysis using a single mutagenic operator. In: Proceedings of the 47th Annual Southeast Regional Conference (ACM-SE). ACM, Clemson, SC, USA, pp. 71–75.
- Untch, R.H., Harrold, M.J., Offutt, A.J., 1997. TUMS: testing using mutant schemata. In: Proceedings of the 35th Annual Southeast Regional Conference (ACM-SE). ACM, Murfreesboro, TN, USA, pp. 174–181.
- Untch, R.H., Offutt, A.J., Harrold, M.J., 1993. Mutation analysis using mutant schemata. In: Proceedings of the International Symposium on Software Testing and Analysis (ISSTA). ACM, Cambridge, MA, USA, pp. 139–148.
- Vincenzi, A.M.R., Maldonado, J.C., Barbosa, E.F., Delamaro, M.E., 2001. Unit and integration testing strategies for C programs using mutation. *Softw. Test. Verif. Reliabil.* 11 (4), 249–268.
- Vincenzi, A.M.R., Nakagawa, E.Y., Maldonado, J.C., Delamaro, M.E., Romero, R.A.F., 2002. Bayesian-learning based guidelines to determine equivalent mutants. *Int. J. Softw. Eng. Knowl. Eng.* 12 (6), 675–689.
- Wang, B., Xiong, Y., Shi, Y., Zhang, L., Hao, D., 2017. Faster mutation analysis via equivalence modulo states. In: Proceedings of the 26th International Symposium on Software Testing and Analysis (ISSTA). ACM, Santa Barbara, CA, USA, pp. 295–306.
- Wedyan, F., Ghosh, S., 2012. On generating mutants for AspectJ programs. *Inf. Softw. Technol.* 900–914.
- Weiss, S.N., Fleysgakker, V.N., 1993. Improved serial algorithms for mutation analysis. In: Proceedings of the International Symposium on Software Testing and Analysis (ISSTA). ACM, Cambridge, MA, USA, pp. 149–158.
- Wohlin, C., 2014. Guidelines for snowballing in systematic literature studies and a replication in software engineering. In: Proceedings of the 18th International Conference on Evaluation and Assessment in Software Engineering (EASE). ACM, London, UK, pp. 1–10.
- Wong, W.E., Maldonado, J.C., Delamaro, M.E., Mathur, A.P., 1994. Constrained mutation in C programs. In: Proceedings of the 8th Brazilian Symposium on Software Engineering (SBES). Brazilian Computer Society, pp. 439–452.
- Wong, W.E., Mathur, A.P., 1995. Reducing the cost of mutation testing: an empirical study. *J. Syst. Softw.* 31 (3), 185–196.
- Wong, W.E., Mathur, A.P., Maldonado, J.C., 1995. Mutation versus all-uses: an empirical evaluation of cost, strength and effectiveness. *Softw. Qual. Product.* 258–265.
- Woodward, M.R., Halewood, K., 1988. From weak to strong, dead or alive? An analysis of some mutation testing issues. In: Proceedings of the 2nd Workshop on Software Testing, Verification, and Analysis. IEEE, Banff, AB, Canada, pp. 152–158.
- Wright, C.J., Kapfhammer, G.M., McMinn, P., 2013. Efficient mutation analysis of relational database structure using mutant schemata and parallelisation. In: Proceedings of the 8th International Workshop on Mutation Analysis (Mutation). IEEE, Luxembourg City, Luxembourg, pp. 63–72.
- Wright, C.J., Kapfhammer, G.M., McMinn, P., 2014. The impact of equivalent, redundant and quasi mutants on database schema mutation analysis. In: Proceedings of the 14th International Conference on Quality Software (QSIC). IEEE, Dallas, TX, USA, pp. 57–66.
- Zhang, J., Wangi, Z., Zhang, L., Hao, D., Zang, L., Cheng, S., Zhang, L., 2016. Predictive mutation testing. In: Proceedings of the 25th International Symposium on Software Testing and Analysis (ISSTA). ACM, Saarbrücken, Germany, pp. 342–353.
- Zhang, L., Gligoric, M., Marinov, D., Khurshid, S., 2013. Operator-based and random mutant selection: better together. In: Proceedings of the 28th International Conference on Automated Software Engineering (ASE). IEEE, Palo Alto, CA, USA, pp. 92–102.
- Zhang, L., Hou, S., Hu, J., Xie, T., Mei, H., 2010. Is operator-based mutant selection superior to random mutant selection? In: Proceedings of the 32th International Conference on Software Engineering (ICSE). ACM, Cape Town, South Africa, pp. 435–444.
- Zhang, L., Marinov, D., Khurshid, S., 2013. Faster mutation testing inspired by test prioritization and reduction. In: Proceedings of the 22nd International Symposium on Software Testing and Analysis (ISSTA). ACM, Lugano, Switzerland, pp. 235–245.
- Zhang, L., Marinov, D., Zhang, L., Khurshid, S., 2012. Regression mutation testing. In: Proceedings of the 21st International Symposium on Software Testing and Analysis (ISSTA). ACM, Minneapolis, MN, USA, pp. 331–341.
- Zhang, L., Xie, T., Zhang, L., Tillmann, N., de Halleux, J., Mei, H., 2010. Test generation via dynamic symbolic execution for mutation testing. In: Proceedings of the 26th International Conference on Software Maintenance (ICSM). IEEE, Timisoara, Romania, pp. 1–10.
- Zhu, Q., Panichella, A., Zaidman, A., 2017. Speeding-up mutation testing via data compression and state infection. In: Proceedings of the 12th International Workshop on Mutation Analysis (Mutation). IEEE, Tokyo, Japan, pp. 103–109.
- Zhu, Q., Panichella, A., Zaidman, A., 2018. An investigation of compression techniques to speed up mutation testing. In: Proceedings of the 11th International Conference on Software Testing, Verification and Validation (ICST). IEEE, Västerås, Sweden, pp. 274–284.

**MSc. Alessandro Viola Pizzoleto** is a PhD candidate at the Computing Department of the Universidade Federal de São Carlos (UFSCar), Brazil. He received a Master's degree in computer science from the Universidade Estadual Paulista (UNESP), Brazil. His main research interests are related to software testing and experimental software engineering.

**Dr. Fabiano Cutigli Ferrari** is an associate professor at the Computing Department of the Universidade Federal de São Carlos (UFSCar), Brazil. He received a PhD degree in computer science from the Universidade de São Paulo (ICMC/USP). His PhD has been supported by FAPESP, and he was a visiting student at Lancaster University/UK supported by CAPES, AOSD-Europe Project (2007–2008), and DiVA Project (2010). He has also been a visiting researcher at George Mason University (2017–2018). His main research interests are related to software testing, experimental

software engineering, adaptive systems, and software metrics. He is a member of the ACM. A list of his main publications can be found at <http://lattes.cnpq.br/3154345471250570>.

**Dr. Jeff Offutt** is a Professor of Software Engineering at George Mason University. He has published over 185 refereed research papers (h-index of 64), and invented numerous widely used test techniques. Offutt is editor-in-chief of Wiley's journal of Software Testing, Verification and Reliability, co-founded the IEEE International Conference on Software Testing (ICST), and co-authored Introduction to Software Testing. He was awarded the Outstanding Faculty Award from the State Council of Higher Education for Virginia in 2019, GMU's Teaching Excellence Award, Teaching With Technology, in 2013, and his 2014 software engineering education paper was chosen by ACM as a notable paper. Current projects include the CS4all educational project, the SPARC educational project, mujava, Testing of Critical System Characteristics (TOCSYC) and PILOT projects at University of Skovde, automatic repair of SQL queries, testing mobile and web applications, test automation, and usable se-

curity. Offutt received the PhD in Computer Science from the Georgia Institute of Technology in 1988. He is on the web at <https://cs.gmu.edu/~offutt/>.

**Leo Fernandes** is a Ph.D. candidate at the Informatics Center of the Universidade Federal de Pernambuco (CIn/UFPE), Brazil. He received a Master's degree in computer science from the Universidade Federal de Pernambuco (UFPE). His main research interests are related to software testing, functional programming, parallel programming and software architecture.

**Dr. Márcio Ribeiro** is an assistant professor at the Computing Institute of the Universidade Federal de Alagoas (UFAL), Brazil. He received a Ph.D. degree in computer science from the Universidade Federal de Pernambuco (UFPE). His main research interests are related to software product lines and families, aspect-oriented and object-oriented programming, experimental software engineering, refactoring, software modularity and software static analysis. A list of his main publications can be found at <http://lattes.cnpq.br/9300936571715992>.