

# A systematic review on regression test selection techniques

Emelie Engström, Per Runeson \*, Mats Skoglund

Department of Computer Science, Lund University, SE-221 00 Lund, Sweden

## ARTICLE INFO

### Article history:

Received 4 February 2009

Received in revised form 24 June 2009

Accepted 10 July 2009

Available online 18 July 2009

### Keywords:

Regression testing

Test selection

Systematic review

Empirical studies

## ABSTRACT

Regression testing is verifying that previously functioning software remains after a change. With the goal of finding a basis for further research in a joint industry-academia research project, we conducted a systematic review of empirical evaluations of regression test selection techniques. We identified 27 papers reporting 36 empirical studies, 21 experiments and 15 case studies. In total 28 techniques for regression test selection are evaluated. We present a qualitative analysis of the findings, an overview of techniques for regression test selection and related empirical evidence. No technique was found clearly superior since the results depend on many varying factors. We identified a need for empirical studies where concepts are evaluated rather than small variations in technical implementations.

© 2009 Elsevier B.V. All rights reserved.

## Contents

1. Introduction .....	15
2. Research method .....	15
2.1. Research questions .....	15
2.2. Sources of information .....	15
2.3. Search criteria .....	16
2.4. Study selection .....	16
2.5. Data extraction and synthesis .....	16
2.6. Qualitative assessment of empirical results .....	17
2.7. Threats to validity .....	17
3. Results .....	18
3.1. Primary studies .....	18
3.2. Analyses of the primary studies .....	18
3.3. Empirically evaluated techniques (RQ1) .....	21
3.3.1. Overview .....	21
3.3.2. Development history .....	21
3.3.3. Uniqueness of the techniques .....	23
3.4. Classification of techniques (RQ2) .....	24
3.5. Analysis of the empirical evidence (RQ3) .....	25
3.5.1. Types of empirical evidence .....	25
3.5.2. Evaluation criteria .....	26
3.6. Comparison of techniques (RQ4) .....	27
3.6.1. Cost reduction .....	27
3.6.2. Fault detection effectiveness .....	27
4. Discussion .....	28
4.1. The reviewed studies .....	28
4.2. Implications for future studies .....	28
5. Conclusions and future work .....	28
Acknowledgments .....	29
References .....	29

\* Corresponding author. Tel.: +46 46 222 93 25.

E-mail addresses: [emelie.engstrom@cs.lth.se](mailto:emelie.engstrom@cs.lth.se) (E. Engström), [per.runeson@cs.lth.se](mailto:per.runeson@cs.lth.se) (P. Runeson), [mats.skoglund@cs.lth.se](mailto:mats.skoglund@cs.lth.se) (M. Skoglund).

## 1. Introduction

Efficient regression testing is important, even crucial, for organizations with a large share of their cost in software development. It includes, among other tasks, determining which test cases need to be re-executed, i.e. regression test selection, in order to verify the behavior of modified software. Regression test selection involves a trade-off between the cost for re-executing test cases, and the risk for missing faults introduced through side effects of changes to the software. Iterative development strategies and reuse are common means of saving time and effort for the development. However both require frequent retesting of previously tested functions due to changes in related code. The need for efficient regression testing strategies is thus becoming more and more important.

A great deal of research effort has been spent on finding cost-efficient methods for different aspects of regression testing. Examples include test case selection based on code changes [1,6,13,17,20,22,43,49,62,64,67] and specification changes [38,40,54,68], evaluation of selection techniques [48], change impact analysis [44], regression tests for different applications, e.g. database applications [18], regression testing of GUIs and test automation [39], and test process enhancement [31]. To bring structure to the topics, researchers have typically divided the field of regression testing into (i) test selection, (ii) modification identification, (iii) test execution, and (iv) test suite maintenance. This review is focused on test selection techniques for regression testing.

Although techniques for regression test selection have been evaluated in previous work [3,15,36,65], no general solution has been put forward since no technique could possibly respond adequately to the complexity of the problem and the great diversity in requirements and preconditions in software systems and development organizations. Neither does any single study evaluate every aspect of the problem; e.g. Kim et al. [27] evaluate the effects of regression test application frequency, Elbaum et al. [11] investigate the impact that different modifications have on regression test selection techniques, several studies examine the ability to reduce regression testing effort [3,11,15,27,36,65,66] and to reveal faults [11,15,27,49].

In order to map the existing knowledge in the field, we launched a systematic review to collect and compare the existing empirical evidence on regression test selection. The use of systematic reviews in the software engineering domain has been subject to a growing interest in the last years. In 2004 Kitchenham proposed a guideline adapted to the specific characteristics of software engineering research. This guideline has been followed and evaluated [5,30,57] and updated accordingly in 2007 [29]. Kitchenham et al. recently published a review of 20 systematic reviews in software engineering during 2004–2007 [28].

Ideally, several empirical studies identified in a systematic review evaluate the same set of techniques under similar conditions on different subject programs. Then there would be a possibility to perform an aggregation of findings or even meta-analysis and thus enable drawing general conclusions. However, as the field of empirical software engineering is quite immature, systematic reviews have not given very clear pictures of the results. In this review we found that the existing studies were diverse, thus hindering proper quantitative aggregation. Instead we present a qualitative analysis of the findings, an overview of the existing techniques for regression test selection and of the amount and quality of empirical evidence.

There are surveys and reviews of software testing research published before, but none of these has the broad scope and the extensive approach of a systematic review. In 2004 Do et al. presented a survey of empirical studies in software testing in general [8]

including regression testing. Their study covered two journals and four conferences over 10 years (1994–2003). Other reviews of regression test selection are not exhaustive but compare a limited number of chosen regression test selection techniques. Rothermel and Harrold presented a framework for evaluating regression test techniques already in 1996 [48] and evaluated the, by that time, existing techniques. Juristo et al. aggregated results from unit testing experiments [25] of which some evaluated regression testing techniques, although with a more narrow scope. Binkley et al. reviewed research on the application of program slicing to the problem of regression testing [4]. Hartman et al. reported a survey and critical assessment of regression testing tools [21]. However, as far as we know, no systematic review on regression test selection research has been carried through since the one in 1996 [48]. An early report of this study was published in 2008 [12], which here is further advanced especially with respect to the detailed description of the techniques (Section 3.4), their development history and the analysis of the primary studies (Section 3.5).<sup>1</sup>

This paper is organized as follows. In Section 2 the research method used for our study is described. In Section 3 the empirical studies and our analyses are reported. In Section 4 the results are discussed, and in Section 5 the work is concluded.

## 2. Research method

### 2.1. Research questions

This systematic review aims at summarizing the current state of the art in regression test selection research by proposing answers to a set of questions below. The research questions stem from a joint industry-academia research project, which aims at finding efficient procedures for regression testing in practice. We searched for candidate regression test selection techniques that were empirically evaluated, and in case of lack of such techniques, to identify needs for future research. Further, as the focus is on industrial use, issues of scale-up to real-size projects and products are important in our review. The questions are:

- (RQ1) Which techniques for regression test selection in the literature have been evaluated empirically?
- (RQ2) Can these techniques be classified, and if so, how?
- (RQ3) Are there significant differences between these techniques that can be established using empirical evidence?
- (RQ4) Can technique *A* be shown to be superior to technique *B*, based on empirical evidence?

Answers to these research questions are searched in the published literature using the procedures of systematic literature reviews as proposed by Kitchenham [29].

### 2.2. Sources of information

In order to gain a broad perspective, as recommended in Kitchenham's guideline [29], we searched widely in electronic sources. The advantage of searching databases rather than a limited set of journals and conference proceedings is also empirically motivated by Dieste et al. [7]. The following seven databases were covered:

- Inspec (<[www.theiet.org/publishing/inspec](http://www.theiet.org/publishing/inspec)>).
- Compendex (<[www.engineeringvillage2.org](http://www.engineeringvillage2.org)>).

<sup>1</sup> In this extended analysis, some techniques that originally were considered different ones, were considered the same technique. Hence, the number of techniques differ from [10]. Further, the quality of two empirical studies was found insufficient in the advanced analysis, why two studies were removed.

- ACM Digital Library (<portal.acm.org>).
- IEEE eXplore (<ieeexplore.ieee.org>).
- ScienceDirect (<www.sciencedirect.com>).
- Springer LNCS (<www.springer.com/lncs>).
- Web of Science (<www.isiknowledge.com>).

These databases cover the most relevant journals and conference and workshop proceedings within software engineering, as confirmed by Dybå et al. [10]. Grey literature (technical reports, some workshop reports, and work in progress) was excluded from the analysis for two reasons: the quality of the grey literature is more difficult to assess and the volume of studies included in the first searches would have grown unreasonably. The searches in the sources selected resulted in overlap among the papers, where the duplicates were excluded primarily by manual filtering.

### 2.3. Search criteria

The initial search criteria were broad in order to include articles with different uses of terminology. The key words used were <regression> and (<test> or <testing>) and <software>, and the database fields of title and abstract were searched. The start year was set to 1969 to ensure that most relevant research within the field would be included, and the last date for inclusion is publications within 2006. The earliest primary study actually included was published in 1997. Kitchenham recommends that exclusion based on languages should be avoided [29]. However, only papers written in English are included. The initial search located 2 923 potentially relevant papers.

### 2.4. Study selection

In order to obtain independent assessments, four researchers were involved in a three-stage selection process, as depicted in Fig. 1.

In the first stage duplicates and irrelevant papers were excluded manually based on titles. In our case, the share of irrelevant papers was extremely large since papers on software for *statistical* regression testing or other regression testing could not be distinguished from papers on *software* regression testing in the database search. The term software did not distinguish between the two areas, since researchers on statistical regression testing often develop some software for their regression test procedures. After the first stage 450 papers remained.

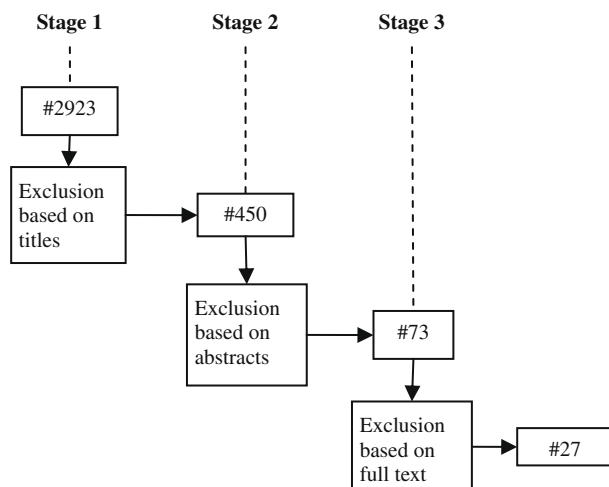


Fig. 1. Study selection procedure.

In the second stage, information in abstracts was analyzed and the papers were classified along two dimensions: research approaches and regression testing approaches. Research approaches were experiment, case study, survey, review, theory and simulation. The two latter types were excluded, as they do not present an empirical research approach, and the survey and review papers were not considered as being primary studies but rather related work to the systematic review. At this stage we did not judge the quality of the empirical data. Regression testing approaches were selection, reduction, prioritization, generation, execution and others. Only papers focusing on regression test *selection* were included.

In the third stage a full text analysis was performed on 73 papers and the empirical quality of the studies was further assessed. The following questions were asked in order to form quality criteria for which studies to exclude before the final data extraction:

- Is the study focused on a specific regression test selection method? For example, a paper could be excluded that presents a method that potentially could be used for regression testing, but is evaluated from another point of view.
- Are the metrics used and the results relevant for a comparison of methods? For example, a paper could be excluded which only reports on the ability to predict fault prone parts of the code, but not on the fault detection effectiveness or the cost of the regression test selection strategy.
- Is data collected and analyzed in a sufficiently rigorous manner? For example, a paper could be excluded if a subset of components was analyzed and conclusions were drawn based on those, without any motivation for the selection.

These questions are derived from a list of questions, used for a similar purpose, published by Dybå et al. [10]. However in our review context, quality requirements for inclusion had to be weaker than suggested by Dybå et al. in order to obtain a useful set of studies to compare. The selection strategy was in general more inclusive than exclusive. Only papers with very poorly reported or poorly conducted studies were excluded, as well as papers in which the comparisons made were considered irrelevant to the original goals of this study.

Abstract analysis and full text analysis were performed in a slightly iterative fashion. Firstly, the articles were independently assessed by two of the researchers. In case of disagreement, the third researcher acted as a checker. In many cases, disagreement was due to insufficient specification of the criteria. Hence, the criteria were refined and the analysis was continued.

In order to get a measure of agreement in the study selection procedure, the Kappa coefficient was calculated for the second stage, which comprised most judgments in the selection. In the second stage 450 abstracts were assessed by two researchers independently. In 41 cases conflicting assessments were made which correspond to the Kappa coefficient  $K = 0.78$ . According to Landis and Koch [33] this translates to a substantial strength of agreement.

### 2.5. Data extraction and synthesis

Using the procedure described in the previous section, 27 articles were finally selected that reported on 36 unique empirical studies, evaluating 28 different techniques. The definition of what constitutes a single empirical study, and what constitutes a unique technique is not always clear cut. The following definitions have been used in our study:

- Study: an empirical study applying a technique to one or more programs. Decisions on whether to split studies with multiple artifacts into different studies were based on the authors' own

classification of the primary studies. Mostly, papers including studies on both small and large programs are presented as two different studies.

- **Technique:** An empirically evaluated method for regression test selection. If the only difference between two methods is an adaption to a specific programming language (e.g. from C++ to Java) they are considered being the same technique.

Studies were classified according to type and size. Two types of studies are included in our review: experiments and case studies. We use the following definitions:

- **Experiment:** A study in which an intervention is deliberately introduced to observe its effects [55].
- **Case study:** An empirical inquiry that investigates a contemporary phenomenon within its real-life context, especially when the boundaries between the phenomenon and context are not clearly evident [69].

Surveys and literature reviews were also considered in the systematic review, e.g. [48,25], but rather as reference points for inclusion of primary studies than as primary studies as such.

Regarding size, the studies are classified as small, medium or large (S, M, L) depending on the study artifact sizes. A small study artifact has less than 2000 lines of code (LOC), a large study artifact has more than 100,000 LOC, and a medium-sized study artifact is in between. The class limits are somewhat arbitrarily defined. In most of the articles the lines of code metric is clearly reported and thus this is our main measurement of size. But in some articles sizes are reported in terms of number of methods or modules, reported as the authors' own statement about the size or are not reported at all.

The classification of the techniques is part of answering RQ2 and is further elaborated in Section 3.4.

## 2.6. Qualitative assessment of empirical results

The results from the different studies were qualitatively analyzed in categories of four key metrics: reduction of cost for test execution, cost for test case selection, total cost, and fault detection effectiveness, see Section 3.5.2. The “weight” of an empirical study was classified according to the scheme in Table 1. A study with more “weight” is considered contributing more to the overall conclusions. A unit of analysis in an experiment is mostly a version of a

piece of code, while in a case study; it is mostly a version of a system or subsystem.

The results from the different studies were then divided into six different categories according to the classification scheme in Table 2. The classification is based on the study “weight” and the size of the difference in a comparative empirical study. As the effect sizes were rarely reported in the studies, the sizes of the differences are also qualitatively assessed. The categorization of results was made by two researchers in parallel and uncertainties were resolved in discussions. Results are presented in Figs. 5–8 in Section 3.5.

## 2.7. Threats to validity

Threats to the validity of the systematic review are analyzed according to the following taxonomy: construct validity, reliability, internal validity and external validity.

Construct validity reflects to what extent the phenomenon under study really represents what the researchers have in mind and what is investigated according to the research questions. The main threat here is related to terminology. Since the systematic review is based on a hierarchical structure of terms – regression test/testing consists of the activities modification identification, test selection, test execution and test suite maintenance – we might miss other relevant studies on test selection that are not specifically aimed for regression testing. However, this is a consciously decided limitation, which has to be taken into account in the use of the results. Another aspect of the construct validity is assurance that we actually find all papers on the selected topic. We analyzed the list of publication fora and the list of authors of the primary studies to validate that no major forum or author was missed.

Reliability focuses on whether the data are collected and the analysis is conducted in a way that it can be repeated by other researchers with the same results. We defined a study protocol setting up the overall research questions, the overall structure of the study as well as initial definitions of criteria for inclusions/exclusion, classification and quality. The criteria were refined during the study based on the identification of ambiguity that could mislead the researchers.

In a systematic review, the decision process for inclusion and exclusion of primary studies is the major focus when it comes to reliability. Our countermeasures taken to reduce the reliability threat were to set up criteria and to use two researchers to classify papers in stages 2 and 3. In cases of disagreement, a third opinion is used. However, the Kappa analysis indicates strong agreements. One of the primary researchers was changed between stages 2 and 3. Still, the uncertainties in the classifications are prevalent and are a major threat to reliability, especially since the quality standards for empirical studies in software engineering are not high enough. Research databases are another threat to reliability [10]. The threat is reduced by using multiple databases; still the non-determinism of some database searches is a major threat to the reliability of any systematic review.

Internal validity is concerned with the analysis of the data. Since no statistical analysis was possible due to the inconsistencies between studies, the analysis is mostly qualitative. Hence we link

**Table 1**  
“Weight” of empirical study.

Type and size of study	Light empirical study “weight”	Medium empirical study “weight”
Experiment (small) Case study (small-medium)	Analysis units <10	Analysis units ≥10
Experiment (medium) Case study (large)	Analysis units <4	Analysis units ≥4

**Table 2**  
Classification scheme for qualitative assessment of the weight of empirical results.

	No difference	Difference of small size	Difference of large size
Medium empirical study “weight”	Strong indication of equivalence between the two compared techniques	Weak indication that one technique is superior to the other	Strong indication that one technique is superior to the other
Light empirical study “weight”	Weak indication of equivalence between the two compared techniques	No indication of differences or similarities	Weak indication that one technique is superior to the other



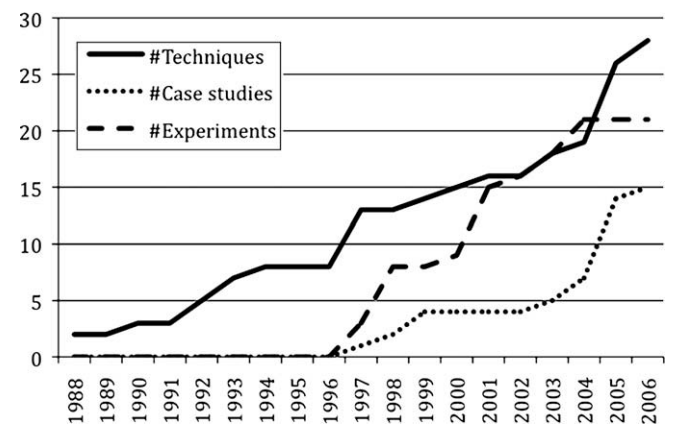


Fig. 2. Accumulated number of published techniques, case studies and experiments.

the conclusions as clearly as possible to the studies, which underpin our discussions.

External validity is about generalizations of the findings derived from the primary studies. Most studies are conducted on small programs and hence generalizing them to a full industry context is not possible. In the few cases where experiments are conducted on small programs as well as case studies on large programs, the external validity is reasonable, although there is room for substantial improvements.

### 3. Results

#### 3.1. Primary studies

The goal of this study was to find regression test selection techniques that are empirically evaluated. The papers were initially obtained in a broad search in seven databases covering the relevant journals, conference and workshop proceedings within software engineering. Then an extensive systematic selection process was carried out to identify papers describing empirical evaluations of regression test selection techniques. The results presented here thus give a good picture of the existing evidence base.

Out of 2923 titles initially screened, 27 papers (P1–P27) on empirical evaluations of techniques for regression test selection re-

mained until the final stage. These 27 papers report on 36 unique studies (S1–S36), see Table 3, and compare in total 28 different techniques for regression test selection for evaluation (T1–T28), see listing in Table 8, constituting the primary studies of this systematic review. Five reference techniques are also identified (REF1–REF5), e.g. *retest all* (all test cases are selected) and *random* (25) (25% of the test cases are randomly selected). In case the studies are reported partially or fully in different papers, we generally refer to the most recent one as this contains the most updated study. When referring to the techniques, we do on the contrary refer to the oldest, considering it being the original presentation of the technique.

In most of the studies, the analyses are based on descriptive statistics. Tabulated data or bar charts are used as a basis for the conclusions. In two studies (S23 and S24), published in the same paper (P16) [46], statistical analysis is conducted using ANOVA.

#### 3.2. Analyses of the primary studies

In order to explore the progress of the research field, and to validate that the selected primary studies reasonably cover our general expectations of which fora and which authors should be represented, we analyze, as an extension to RQ1, aspects of the primary studies as such: where they are published, who published them, and when. As defined in Section 2.5, a paper may report on multiple studies, and in some cases the same study is reported in more than one paper. Different researchers have different criteria for what constitutes a study. We have tried to apply a consistent definition of what constitutes a study. This distribution of studies over papers is shown in Table 4. Most papers (18 out of 27) report a single study, while few papers report more than one. Two papers report new analyses of earlier published studies. Note that many of the techniques are originally presented in papers without empirical evaluation, hence these papers are not included as primary studies in the systematic review, but are referenced in Section 3.3 as sources of information about the techniques as such (Table 8). The number of identified techniques in the primary studies is relatively high compared to the number of studies, 28 techniques were evaluated in 36 studies. Table 5 presents the distribution of the number of studies in which different techniques occur. One technique was present in 14 different studies, another technique in 8 studies, etc. 14 techniques only appear in one study, which is not satisfactory when trying to aggregate information from empirical evaluations of the techniques.

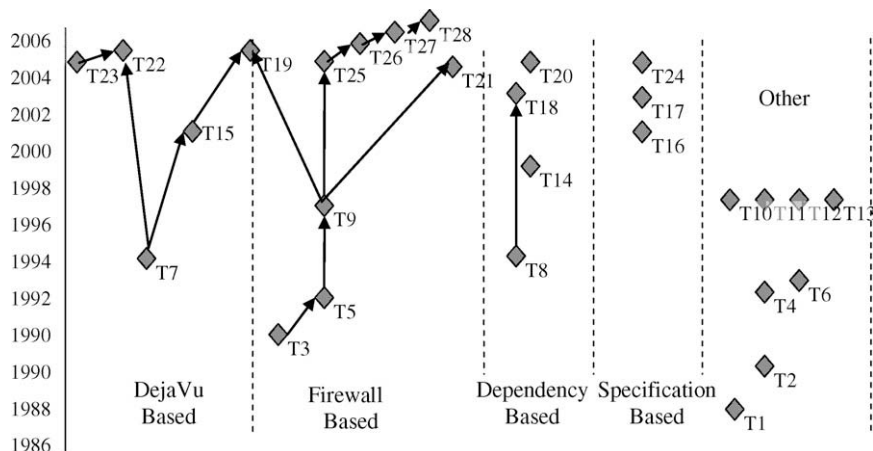


Fig. 3. Evolution of techniques.

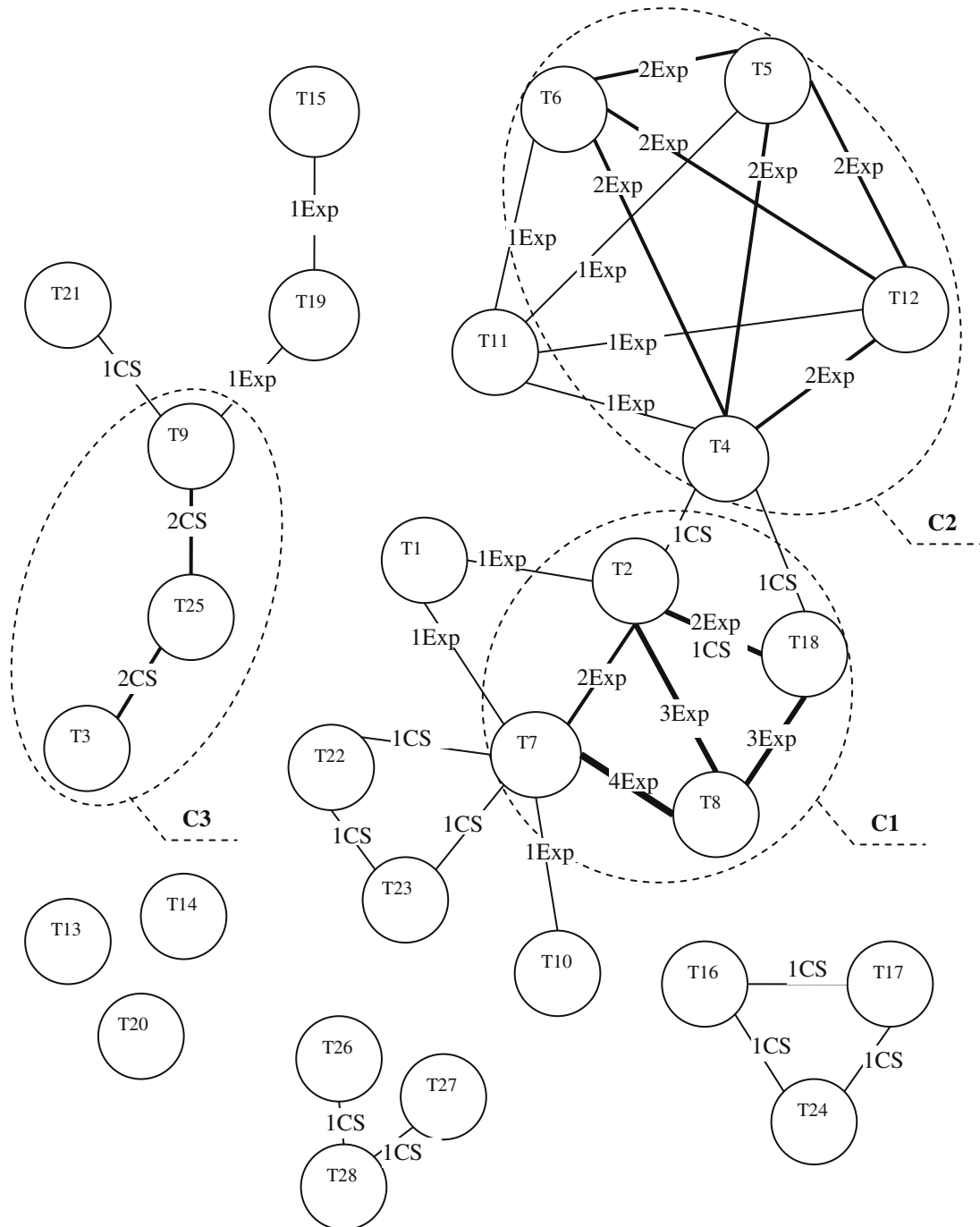


Fig. 4. Techniques related to each other through empirical comparisons.

Table 6 lists the different publication fora in which the articles have been published. It is worth noting regarding the publication fora that the empirical regression testing papers are published in a wide variety of journals and conference proceedings. Limiting the search to fewer journals and proceedings would have missed many papers.

The major software engineering journals and conferences are represented among the fora. It is not surprising that a conference on software maintenance is on the top, but we found, during the validity analysis, that the International Symposium on Software Testing and Analysis is not on the list at all. We checked the proceedings specifically and have also noticed that, for testing in general, empirical studies have been published there, as reported by

Do et al. [8], but apparently not on regression test selection during the studied time period.

Table 7 lists authors with more than one publication. In addition to these 17 authors, five researchers have authored or co-authored one paper each. On the top of the author's list, we find the names of the most prolific researchers in the field of regression test selection (Rothermel and Harrold). It is interesting to notice from the point of view of conducting empirical software engineering research that there are two authors on the top list with industry affiliation (Robinson and Smiley).

The regression test selection techniques have been published from 1988 to 2006, as shown in Fig. 2 and Table 8. The first empirical evaluations were published in 1997 (one case study and three

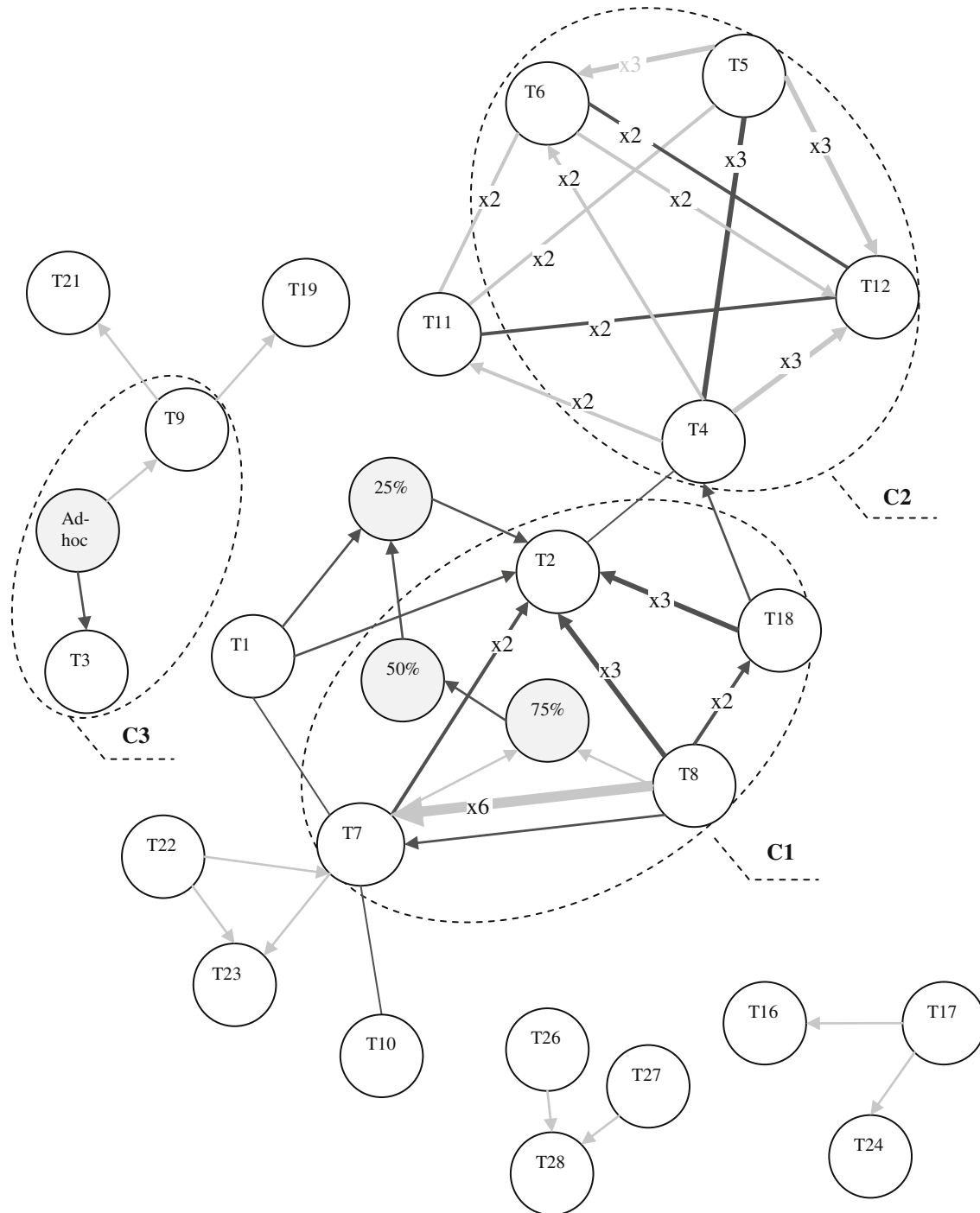


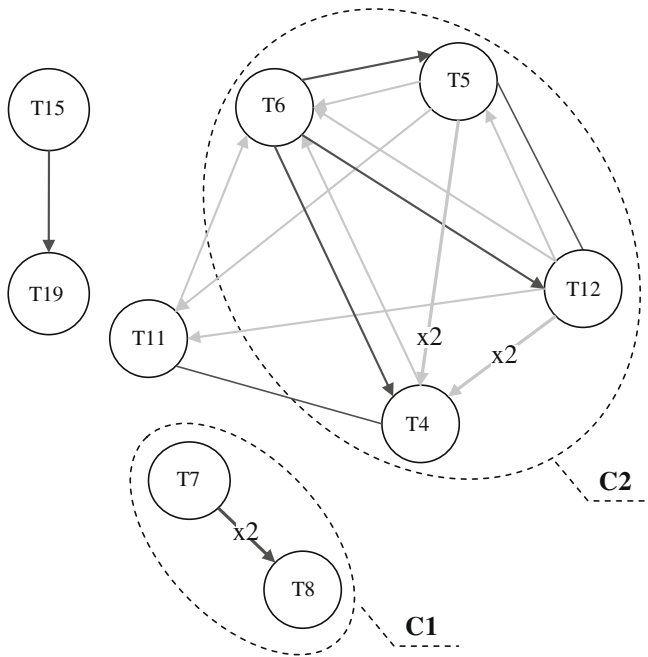
Fig. 5. Empirical results for cost reduction, including test execution time, test suite reduction and precision.

experiments), hence the empirical evaluations have entered the scene relatively late. Twelve out of the 28 techniques have been originally presented and evaluated in the same paper: T12–S11 and T13–S32 (1997); T14–S33–S34 (1999); T18–S5 (2003); T19–S13 (2004); T20–S14; T21–S25; T23–S31; T25–S29–S30 and T26–S35 (2005); T27–S33 and T28–S35 (2006).

We conclude from this analysis that there are only a few studies comparing many techniques in the same study, making it hard to find empirical data for a comprehensive comparison. However, some small- and medium-sized artifacts have appeared as a de-

facto benchmark in the field [8], enabling comparison of some techniques to some extent.

Most of the expected publication fora are represented, and one that is not represented, but was expected, was specifically double checked. Similarly, well known researchers in the field were among the authors, hence we consider the selected primary studies as being a valid set. It is clear from the publication analysis that the techniques published during the later years are published with empirical evaluations to a higher degree than those published during the earlier years, which is a positive trend in



**Fig. 6.** Empirical results for test selection time.

searching for empirically evaluated techniques as defined in RQ1.

### 3.3. Empirically evaluated techniques (RQ1)

### 3.3.1. Overview

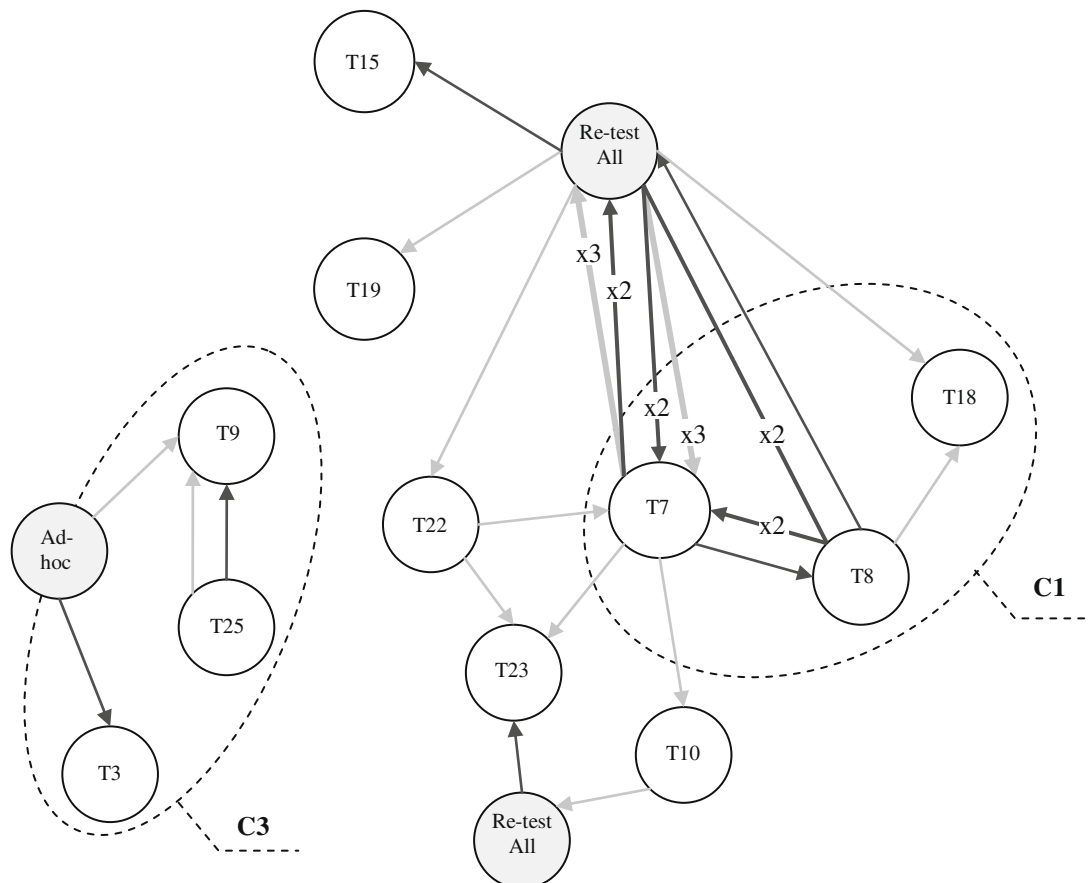
**Table 8** lists the 28 different regression test selection techniques (T1–T28), in chronological order of their first publication. In case the studies are reported partially or fully in different papers, we generally refer to the original one. In case a later publication has added details that are needed for exact specification of the technique, both references are used.

This list is specifically the answer to the first research question: which techniques for regression test selection existing in the literature have been evaluated empirically (RQ1). In this review, the techniques, their origin and description, are identified in accordance with what is stated in each of the selected papers, although adapted according to our definition of what constitutes a unique technique in Section 2.5.

### 3.3.2. Development history

The historical development chain gives some advice on which techniques are related and how they are developed, see Fig. 3. There are three major paths, beginning with T3, T7 and T8, respectively.

One group of techniques is the *firewall* techniques where dependencies to modified software parts are isolated inside a firewall. Test cases covering the parts within the firewall are selected for re-execution. The first firewall technique (T3) for procedural languages was presented by Leung and White in 1990 [35]. An empirical evaluation used a changed version (T5). The technique was adapted to object-oriented languages T9 in two similar ways



**Fig. 7.** Empirical results for total time.





**Table 3**

Primary studies, S1–S36, published in papers P1–P27, evaluation techniques T1–T28.

Study ID	Publication ID	Reference	Techniques	Artifacts	Type of study	Size of study
S1	P1	Baradhi and Mansour [2]	T4, T5, T6, T11, T12 REF1	Own unspecified	Exp	S
S2	P2	Bible et al. [3]	T7, T8 REF1	7x Siemens, Small constructed programs, constructed, realistic non-coverage based test suites	Exp	S
S3	P2	Bible et al. [3]	T7, T8 REF1	Space, Real application, real faults, constructed test cases	Exp	S
S4	P2	Bible et al. [3]	T7, T8 REF1	Player, One module of a large software system constructed realistic test suites	Exp	M
S5	P3	Elbaum et al. [11]	T2, T4, T18 REF1	Bash, Grep, Flex and Gzip, Real, non-trivial C program, constructed test suites	CS (Mult)	M
S6	P4	Frankl et al. [14]	T7, T10 REF1	7xSiemens, Small constructed programs, constructed, realistic, non-coverage based test suites	Exp	S
S7	P5	Graves et al. [15]	T1, T2, T7 REF1, REF2, REF3, REF4	7xSiemens, Small constructed programs, constructed, realistic non-coverage based test suites; space, Real application, real faults, constructed test cases; player, One module of a large software system constructed realistic test suites	Exp	S M
S8	P6	Harrold et al. [19]	T15 REF1	Siena, Jedit, JMeter, RegExp, Real programs, constructed faults	Exp	S
S9	P7	Kim et al. [27]	T2, T7, T8 REF1, REF2, REF3, REF4	7xSiemens, Small constructed programs, constructed, realistic non-coverage based test suites; Space, Real application, real faults, constructed test cases	Exp	S
S10	P8	Koju et al. [32]	T15 REF1	Classes in .net framework, Open-source, real test cases	Exp	S
S11	P9	Mansour et al. [36]	T4, T5, T6, T12	20 small sized Modules	Exp	S
S12	P10	Mao and Lu [38]	T16, T17, T24 REF1	Triangle, eBookShop, ShipDemo, Small Constructed programs	CS	S
S13	P11	Orso et al. [41]	T9, T15, T19 REF1	Jaba, Daikon, JBoss, Real-life programs, original test suites	Exp	M L
S14	P12	Pasala and Bhowmick [42]	T20 REF1	Internet Explorer (client), IIS (web server), application (app. Server), An existing browser based system, real test cases	CS	NR
S15	P13	Rothermel and Harrold [49]	T7 REF1	7xSiemens, Small constructed programs, constructed, realistic non-coverage based test suites	Exp	S
S16	P13	Rothermel and Harrold [49]	T7 REF1	Player, One module of a large software system constructed realistic test suites	Exp	M
S17	P14	Rothermel and Harrold [50]	T7 REF1	7xSiemens, Small constructed programs, constructed, realistic non-coverage based test suites	Exp	S
S18	P14	Rothermel and Harrold [50]	T7 REF1	7xSiemens, Small constructed programs, constructed, realistic non-coverage based test suites	Exp	S
S19	P14	Rothermel and Harrold [50]	T7 REF1	7xSiemens, Small constructed programs, constructed, realistic non-coverage based test suites;	Exp	S
S20	P14	Rothermel and Harrold [50]	T7 REF1	Player, One module of a large software system constructed realistic test suites	Exp	M
S21	P14	Rothermel and Harrold [50]	T7 REF1	Commerercial, Real application, real test suite	Exp	S
S22	P15	Rothermel et al. [45]	T8, T18 REF1	Emp-server, Open-source, server component, constructed test cases; Bash Open-source, real and constructed test cases	Exp	M
S23	P16	Rothermel et al. [46]	T2, T8, T18 REF1	Bash, Open-source, real and constructed test cases	Exp	M
S24	P16	Rothermel et al. [46]	T2, T8, T18 REF1	Emp-server, Open-source, server component, constructed test cases	Exp	M
S25	P17	Skoglund and Runeson [56]	T9, T21 REF1	Swedbank, Real, large scale, distributed, component-based, J2EE system, constructed, scenario-based test cases	CS	L
S26	P18	Vokolos and Frankl [65]	T10 REF1	ORACOLO2, Real industrial subsystems, real modifications, constructed test cases	CS	M
S27	P19	White and Robinson [61]	T3 REF5	14 real ABB projects, Industrial, Real-time system	CS	L
S28	P19	White and Robinson [61]	T9 REF5	2 real ABB projects, Industrial, Real-time system	CS	L
S29	P20	White et al. [60]	T3, T9, T25	OO-telecommunication software system	CS	S
S30	P20	White et al. [60]	T3, T9, T25	OO – real-time software system	CS	L
S31	P21	Willmor and Embury [63]	T7, T22, T23 REF1	Compiere, James, Mp3cd browser, Open-source systems, real modifications	CS	NR
S32	P22	Wong et al. [66]	T13 REF1	Space, Real application, real faults, constructed test cases	CS	S
S33	P23	Wu et al. [67]	T14 REF1	ATM-simulator, small constructed program	CS	S
S34	P23	Wu et al. [67]	T14 REF1	Subsystem of a fully networked supervisory control and data analysis system	CS	M
S35	P24, P25, P26	Zheng et al. [71,72] and Zheng [70]	T26, T28 REF1	ABB-internal, Real C/C++ application	CS	M
S36	P27, P25	Zheng et al. [73,72]	T27, T28 REF1	ABB-internal, Real C/C++ application	CS	M

### 3.3.3. Uniqueness of the techniques

There is a great variance regarding the uniqueness of the techniques identified in the studied papers. Some techniques may be regarded as novel at the time of their first presentation, while others may be regarded as only variants of already existing techniques. For example in [3] a regression test selection technique is evaluated, T8, and the technique used is based on modified entities in the subject programs. In another evaluation reported in [11], it is stated that the same technique is used as in [3] but is adapted to use a different scope of what parts of the subjects programs are in-

cluded in the analysis, T18. In [3] the complete subject programs are included in the analysis; while in [11] core functions of the subject programs are ignored. This difference of scope probably has an effect on the test cases selected using the two different approaches. The approach in which core functions are ignored is likely to select fewer test cases compared to the approach in which all parts of the programs are included. It is not obvious whether the two approaches should be regarded as two different techniques or if they should be regarded as two very similar variants of the same technique. We chose the former option.

**Table 4**

Distribution of number of papers after the number of studies each paper reports.

# reported studies in each paper	# papers	# studies
0 (re-analysis of another study)	2	0
1	18	18
2	5	10
3	1	3
5	1	5
Total	27	36

**Table 5**

Distribution of techniques after occurrences in number of studies.

Represented in number of studies	Number of techniques
14	1
8	1
5	2
4	1
3	2
2	7
1	14
Total	28

Some techniques evaluated in the reviewed papers are specified to be used for a specific type of software, e.g. Java, T15 and T19 [19,41], component-based software, T17, T20, T24 and T28 [38,42,72,73], or database-driven applications, T22, [63]. It is not clear whether they should be considered as one technique applied to two types of software, or two distinctly different techniques. For example, a technique specified for Java, T15, is presented and evaluated in [19]. In [58] the same technique is used on MSIL (Microsoft Intermediate Language) code, however adapted to handle programming language constructs not present in Java. Thus, it can be argued that the results of the two studies cannot be synthesized in order to draw conclusions regarding the performance of neither the technique presented in [19], nor the adapted version, used in [32]. However, we chose to classify them as the same technique.

There are also techniques specified in a somewhat abstract manner, e.g. techniques that handle object-oriented programs in general, e.g. T14 [67]. However, when evaluating a technique, the abstract specification of a technique must be concretized to handle the specific type of subjects selected for the evaluation. The concretization may look different depending on the programming language used for the subject programs. T14 is based on dependencies

**Table 7**

Researchers and number of publications.

Name	#	Name	#
Rothermel G.	9	Baradhi G.	2
Harrold M. J.	5	Frankl P. G.	2
Robinson B.	5	Kim J. M.	2
Zheng J.	4	Mansour N.	2
Elbaum S. G.	3	Orso A.	2
Kallakuri P.	3	Porter A.	2
Malishevsky A.	3	White L.	2
Smiley K.	3	Vokolos F.	2
Williams L.	3		

between functions in object-oriented programs in general. The technique is evaluated by first tailoring the abstract specification of the technique to C++ programs and then performing the evaluation on subject programs in C++. However, it is not clear how the tailoring of the specification should be performed to evaluate the technique using other object-oriented programming languages, e.g. C# or Java. Thus, due to differences between programming languages, a tailoring made for one specific programming language may have different general performance than a tailoring made for another programming language.

### 3.4. Classification of techniques (RQ2)

In response to our second research question (RQ2), we are looking for some kind of classification of the regression test selection techniques. As indicated in Fig. 3, there exist many variants of techniques, gradually evolved over time. Some suggested that classifications of regression test techniques exist. Rothermel and Harrold present a framework for analyzing regression test selection techniques [48], including evaluation criteria for the techniques: inclusiveness, precision, efficiency and generality. Graves et al. [15] present a classification scheme where techniques are classified as Minimization, Safe, Dataflow-Coverage-based, Ad hoc/Random or Retest-All techniques. Orso et al. [41] separate between techniques that operate at a higher granularity, e.g. method or class (called high-level), and techniques that operate at a finer granularity, e.g. statements (called low-level). In this review we searched for classifications in the papers themselves with the goal of finding common properties in order to be able to reason about groups of regression testing techniques.

One property found is regarding the type of input required by the techniques. The most common type of required input is source code text, e.g. T2–T8, T10–T14 and T18. Other types of codes ana-

**Table 6**

Number of papers in different publication fora.

Publication fora	Type	#	%
International Conference on Software Maintenance	Conference	5	18.5
ACM Transactions of Software Engineering and Methodology	Journal	3	11.1
International Symposium on Software Reliability Engineering	Conference	3	11.1
International Conference on Software Engineering	Conference	3	11.1
Asia-Pacific Software Engineering Conference	Conference	2	7.4
International Symposium on Empirical Software Engineering	Conference	2	7.4
IEEE Transactions of Software Engineering	Journal	1	3.7
Journal of Systems and Software	Journal	1	3.7
Software Testing Verification and Reliability	Journal	1	3.7
Journal of Software Maintenance and Evolution	Journal	1	3.7
ACM SIGSOFT Symposium on Foundations of SE	Conference	1	3.7
Automated Software Engineering	Conference	1	3.7
Australian SE Conference	Conference	1	3.7
International Conf on COTS-based Software Systems	Conference	1	3.7
International Conference on Object-Oriented Programming, Systems, Languages, and Applications	Conference	1	3.7
Total		27	100

**Table 8**

Techniques for regression test selection.

Technique	Origin	Description	Evaluated in study
T1	Harrold and Soffa [20]	Dataflow-coverage-based	S7
T2	Fischer et al. [13] and Hartman and Robson [22]	Modification-focused, minimization, branch and bound algorithm	S5, S7, S9, S23, S24
T3	Leung and White [35]	Procedural-design firewall	S27, S29, S30
T4	Gupta et al. [16]	Coverage-focused, slicing	S1, S5, S11
T5	White and Leung [62]	Firewall	S1, S11
T6	Agraval et al. [1]	Incremental	S1, S11
T7	Rothermel and Harrold [47]	Viewing statements, DejaVu	S2–S4, S6, S7, S9, S15–S21, S31
T8	Chen and Rosenblum [6]	Modified entity – TestTube	S2–S4, S9, S22–S24
T9	Pei et al. [43] and White and Abdullah [59]	High level – identifies changes at the class and interface level	S13, S25, S28–S30
T10	Vokolos and Frankl [64]	Textual differing – Pythia	S6, S26
T11	Mansour and Fakhri [37]	Genetic algorithm	S1
T12	Mansour and Fakhri [37]	Simulated annealing	S1, S11
T13	Wong et al. [66]	Hybrid: modification, minimization and prioritization-based selection	S32
T14	Wu et al. [67]	Analysis of program structure and function-calling sequences	S33, S34
T15	Rothermel et al. [51], Harrold et al. [19] and Koju et al. [32]	Edge level – identifies changes at the edge level	S8, S10, S13
T16	Orso et al. [40]	Use of metadata to represent links between changes and test cases	S12
T17	Sajeev et al. [54]	Use of UML (OCL) to describe information changes	S12
T18	Elbaum et al. [11]	Modified-non-core, same as T8 but ignoring core functions	S5, S22
T19	Orso et al. [41]	Partitioning and selecting two phases	S13
T20	Pasala and Bhowmick [42]	Runtime dependencies captured and modeled into a graph (CIG)	S14
T21	Skoglund and Runeson [56]	Change-based selection	S25
T22	Willmor and Embury [63]	Test selection for DB-driven applications (extension of T7) combined safety	S31
T23	Willmor and Embury [63]	Database safety	S31
T24	Mao and Lu [38]	Enhanced representation of change information	S12
T25	White et al. [60]	Extended firewall additional data-paths	S29, S30
T26	Zheng [71]	I-BACCI v.1	S35
T27	Zheng et al. [73]	I-BACCI v.2 (firewall + BACCI)	S36
T28	Zheng et al. [73]	I-BACCI v.3	S35, S36
REF1	Leung and White [34]	Retest-all	S1–S10, S12–S24, S26, S31–S36
REF2		Random (25)	S7, S9
REF3		Random (50)	S7, S9
REF4		Random (75)	S7, S9
REF5		Intuitive, experience-based selection	S27, S28

lyzed by the techniques are intermediate codes for virtual machines, e.g. T9, T15 and T21, or machine codes, e.g. T24 and T26. Some techniques require input of a certain format, e.g. T16 (meta data) and T17 (OCL). Techniques may also be classified according to the type of code used in the analysis (Java, C++...). A third type of classification that could be extracted from the papers is regarding the programming language paradigm. Some techniques are specified for use with procedural code, e.g. T2, T7, T8, and T18, while other techniques are specified for the object-oriented paradigm, e.g. T9, T14–T17, and T21–T23. Some techniques are independent of programming language, e.g. T1, T20, and T26–T28.

The most found property assigned to regression test selection techniques is whether they are *safe* or *unsafe*. With a safe technique the defects found with the full test suite are also found with the test cases picked by the regression test selection technique. This property may be used to classify all regression test selection techniques into either *safe* or *unsafe* techniques. Retest all is an example of a safe technique since it selects all test cases, hence, it is guaranteed that all test cases that reveal defects are selected. Random selection of test cases is an example of an unsafe technique since there is a risk of test cases revealing defects being missed. In our study seven techniques were stated by the authors to be safe: T7, T8, T10, T15, and T22–T24. However, the safety characteristic is hard to achieve in practice, as it e.g. assumes determinism in program and test execution.

A major problem in addition to finding a classification scheme is applying the scheme to the techniques. The information regarding the different properties is usually not available in the publications. Hence, we may only give examples of techniques having the properties mentioned above based on what the authors state in their

publications. The properties reported for each technique are presented in Table 9.

### 3.5. Analysis of the empirical evidence (RQ3)

Once we have defined which empirical studies exist and a list of the techniques they evaluate, we continue with the third research question on whether there are significant differences between the techniques (RQ3). We give an overview of the primary studies as such in Section 3.5.1. Then we focus on the metrics and evaluation criteria used in different studies (Section 3.5.2).

#### 3.5.1. Types of empirical evidence

Table 10 overviews the primary studies by research method, and the size of the system used as subject. We identified 21 unique controlled experiments and 15 unique case studies. Half of the experiments are conducted on the same set of small programs [23], often referred to as the Siemens programs, which are made available through the software infrastructure repository<sup>2</sup> presented by Do et al. [8]. The number of large scale real-life evaluations is sparse. In this systematic review we found four evaluations (S25, S27, S28, and S30). Both types of studies have benefits and encounter problems, and it would be of interest to study the link between them, i.e. does a technique which is shown to have great advantages in a small controlled experiment show the same advantages in a large scale case study. Unfortunately no complete link was found in this review. However, the move from small toy programs to med-

<sup>2</sup> <http://sir.unl.edu>.

**Table 9**

Overview of properties for each technique.

Technique	Applicability		Method			Properties	
	Type of language	Type of software	Input	Approach	Granularity	Detection ability	Cost reduction
T1	Ind		IM	CF	Stm		
T2	Proc		SC	CF	Stm		Min
T3	Proc		SC	FW	Module		
T4	Proc		SC	Slicing	Stm		Min
T5	Proc		SC	FW	Module		
T6	Proc		SC	Slicing	Stm		
T7	Proc		SC	CF	Stm	Safe	
T8	Proc		SC	Dep	Func	Safe	
T9	OO		IM	FW	Class		
T10	Proc		SC		Stm	Safe	
T11	Proc		SC	CF	Stm		
T12	Proc		SC	CF	Stm		
T13	Proc		SC		Stm		Min
T14	OO		SC	Dep	Func		
T15	OO		IM	CF	Stm	Safe	
T16	OO	Comp	Spec	CF	Stm		
T17	OO	Comp	Spec				
T18	Proc		SC	Dep	Func		
T19	OO		IM	FW + CF	Class + Stm		
T20	Ind	Comp	BIN	Dep	Comp		
T21	OO		IM	FW	Class		
T22	OO	DB	SC	CF	Stm	Safe	
T23	OO	DB	SC	CF	Stm	Safe <sup>a</sup>	
T24	OO	Comp	BIN + Spec	Dep	Stm	Safe	
T25	OO		SC?	FW	Class		
T26	Ind	Comp	BIN	FW	Func		
T27	Ind	Comp	BIN + SC	FW	Func		
T28	Ind	Comp	BIN + SC	FW	Func		
	Proc = Procedural language; Ind = Independent; OO = Object oriented	Comp = Component based; DB = Database driven	SC = Source code; IM = Intermediate code for virtual machines; BIN = Machine code; Spec = Input of a certain format	CF = Control flow; FW = Fire wall; Slicing; Dep = Dependency based	Stm = Statement; Func = Function; Class; Module; Comp = Component	Safe	Min = Minimization

<sup>a</sup> Safe only in DB-state.

ium-sized components, which is observed among the studies, is a substantial step in the right direction towards real-world relevance and applicability.

The empirical quality of the studies varies a lot. In order to obtain a sufficiently large amount of papers, our inclusion criteria regarding quality had to be weak. Included in our analysis was any empirical evaluation of regression test selection techniques if relevant metrics were used and a sufficiently rigorous data collection and analysis could be followed in the report, see Section 2.4 for more details. This was independently assessed by two researchers.

An overview of the empirically studied relations between techniques and studies is shown in Fig. 4. Circles represent techniques and connective lines between the techniques represent comparative studies. CS on the lines refers to the number of case studies conducted in which the techniques are compared, and Exp denotes the number of experimental comparisons. Some techniques have not been compared to any of the other techniques in the diagram:

**Table 10**

Primary studies of different types and sizes.

Type of studies	Size of subjects under study	Number of studies	%
Experiment	Large	1	3
Experiment	Medium	7	19
Experiment	Small	13	36
Case study	Large	4	11
Case study	Medium	5	14
Case study	Small	4	11
Case study	Not reported	2	6
	Total	36	100

T13, T14 and T20. These techniques are still empirically evaluated in at least one study, typically a large scale case study. If no comparison between the proposed techniques is made, the techniques are compared to a reference technique instead, e.g. the retest of all test cases, and in some cases a random selection of a certain percentage of test cases is used as a reference as well. The reference techniques are not shown in Fig. 4 for visibility reasons.

Researchers are more apt to evaluate new techniques or variants of techniques than to replicate studies, which is clearly indicated by that we identified 28 different techniques in 27 papers. This gives rise to clusters of similar techniques compared among themselves and techniques only compared to a reference method such as retest all.

Three clusters of techniques have been evaluated sufficiently to allow for meaningful comparison, see Fig. 4; C1: T2, T7, T8 and T18, C2: T4, T5, T6 and T12, and C3: T3, T9 and T25. Each of these pairs of techniques has been compared in at least two empirical studies. However, not all studies are conducted according to the same evaluation criteria, nor is the quality of the empirical evidence equally high. Therefore we classified the results with respect to empirical quality, as described in Section 2.6, and with respect to evaluation criteria, as described below.

### 3.5.2. Evaluation criteria

Do and Rothermel proposed a cost model for regression testing evaluation [9]. However, this model requires several data which are not published in the primary studies. Instead, we evaluated the results with respect to each evaluation criterion separately. We identified two main categories of metrics: *cost reduction* and *fault detection effectiveness*. Five different aspects of cost reduction



**Table 11**

Use of evaluation metrics in the studies.

Evaluated metrics	Number	%	Rothermel framework [48]
<i>Cost reduction</i>			
Test suite reduction	29	76	Efficiency
Test execution time	7	18	Efficiency
Test selection time	5	13	Efficiency
Total time	16	42	Efficiency
Precision (omission of non-fault-revealing tests)	1	3	Precision
<i>Fault detection effectiveness</i>			
Test case-related detection effectiveness	5	13	Inclusiveness
Fault-related detection effectiveness	8	21	

and two of fault detection effectiveness have been evaluated in the primary studies. Table 11 gives an overview of the extent to which the different metrics are used in the studies. Size of test suite reduction is the most frequent, evaluated in 76% of the studies. Despite this, it may not be the most important metric. If the cost for performing the selection is too large in relation to that for performing the reduction, no savings are achieved. In 42% of the studies the total time (test selection and execution) is evaluated instead. The effectiveness measures are either related (1) to test cases, i.e. the percentage of fault-revealing test cases selected out of all fault-revealing test cases, or (2) to faults, i.e. the percentage of faults out of all known ones, detected by the selected test cases.

Several of the studies concerning reduction of number of test cases are only compared to retest all (S8, S10, S14–S21, S26, S32–S34) [19,32,42,49,50,65–67] with the only conclusion that a reduction of test cases can be achieved, but nothing on the size of the effect in practice. This is a problem identified in experimental studies in general [26]. Many studies evaluating time reduction are conducted on small programs, and the size of the differences is measured in milliseconds, although there is a positive trend, over time, towards using medium-sized programs. Only 30% of the studies consider both fault detection and cost reduction. Rothermel proposed a framework for evaluation of regression test selection techniques [48] which have been used in some evaluations. This framework defines four key metrics: *inclusiveness*, *precision*, *efficiency*, and *generality*. Inclusiveness and precision correspond to test case-related fault detection effectiveness and precision, respectively, as shown in Table 11. Efficiency is related to space and time requirements and varies with test suite reduction as well as with test execution time and test selection time. Generality is more of a theoretical reasoning, which is not mirrored in the primary studies.

### 3.6. Comparison of techniques (RQ4)

In response to our fourth research question (RQ4) we are analyzing the empirically evaluated relations between the techniques by visualizing the results of the studies. Due to the diversity in evaluation criteria and in empirical quality this visualization cannot give a complete picture. However, it may provide answers to specific questions: e.g. is there any technique applicable in my context proven to reduce testing costs more than the one we use today?

Our taxonomy for analyzing the evidence follows the definitions in Table 2. Grey arrows indicate *light weight* empirical result and black arrows indicate *medium weight* result. A connection without arrows in the figures means that the studies have similar effect, while where there is a difference, the arrow points to the technique that is better with respect to the chosen criterion. A connection with thicker line represents more studies. In Section 3.6.1, we report our findings regarding cost reduction and in Section 3.6.2, we report our findings regarding fault detection. Note that

the numbers on the arrows indicate the number of collected metrics, which may be more than one per study.

#### 3.6.1. Cost reduction

Fig. 5 reports the empirically evaluated relations between the techniques regarding the cost reduction, including evaluations of execution time as well as of test suite reduction and precision.

The strongest evidence can be found in cluster C1, where T2 provides most reduction of execution costs. T7, T8 and T18 reduce the test suites less than T2, and T8 among those reduces execution cost less than T18. All techniques, however, reduce test execution cost compared to REF1 (retest all), which is a natural criterion for a regression test selection technique.

In cluster C2, there is strong evidence that T6 and T12 have similar cost for test execution. On the other hand, there is a study with weaker empirical evidence, indicating that T12 reduces execution cost more than T6.

The rest of the studies show rather weak empirical evidence, showing that the evaluated techniques reduce test execution cost better than retest all.

One component of the cost for regression test selection is the analysis time needed to select which test cases to re-execute. The selection time is reported separately for a small subset of the studies, as shown in Fig. 6.

The left group primarily tells that T19 has less selection time than T15, and in C1, T8 has less analysis time than T7.

The results from cluster C2 show mixed messages. T4 has in most cases the shortest selection time, although it is in one study more time consuming than T6. The selection time is hence dependent on the subject programs, test cases, and types of changes done.

In Fig. 7, the total time for analysis and execution together is shown for those studies where it is reported. It is worth noting that some regression test selection techniques actually can be more time consuming than retest all (T7, T8, T10). Again, this is case dependent, but it is interesting to observe that this situation actually arises under certain conditions.

Other relations are a natural consequence of the expansion of certain techniques. T9 (object-oriented firewall) is less time consuming than T25 (extended OO firewall with data-paths). Here an additional analysis is conducted in the regression test selection.

#### 3.6.2. Fault detection effectiveness

In addition to saving costs, regression test selection techniques should detect as many as possible of the faults found by the original test suite. Evaluations of test case-related as well as fault-related detection effectiveness are presented in Fig. 8.

Some techniques are proven to be *safe*, i.e. guarantee that the fault detection effectiveness is 100% compared to the original test suite (see Section 3.4). This property is stated to hold for seven techniques: T7, T8, T10, T15, T22, T23 and T24.

T7 and T8 within C2 are also those that can be found superior or equal from Fig. 8, which is in line with the *safe* property. T4 in C2 also tends to be better or equal to all its reference techniques. However, for the rest, the picture is not clear.

## 4. Discussion

### 4.1. The reviewed studies

The overall goal of the study was to identify regression test selection techniques and systematically assess the empirical evidence collected about those techniques. As the selection of a specific technique is dependent on many factors, the outcomes of empirical studies also depend on those factors. However only few factors are specifically addressed in the empirical studies and hence it is not possible to draw very precise conclusions. Nor is it possible to draw general conclusions. Instead we have conducted mostly qualitative assessments of the empirical studies. From those we try to aggregate recommendations of which regression test selection techniques to use.

A comparison of the techniques in cluster C1 indicates that the minimization technique, T2, is the most efficient in reducing time and/or number of test cases to run. However this is an unsafe technique (see Section 3.4) and all but one of six studies report on significant losses in fault detection. When it comes to safe techniques, T7 is shown to be the most efficient in reducing test cases. However analysis time for T7 is shown to be too long (it exceeds the time for rerunning all test cases) in early experiments, while in later experiments, it is shown to be good. Hence, there is a trade-off between cost reduction and defect detection ability. This is the case in all test selections, and none of the evaluated techniques seems to have done any major breakthrough in solving this trade-off.

It is interesting to notice that the technique T7 is not changed between the studies that show different results on selection time, but the subject programs on which the experiments are conducted are changed. The subject programs are one factor that heavily impacts the performance of some techniques. This emphasizes the importance of the regression testing context in empirical studies, and may also imply that specific studies have to be conducted when selecting a technique for a specific environment.

As mentioned before, many techniques are incremental improvements of existing techniques, which are demonstrated to perform better. For example, T25 is an extension of T9, with better fault detection at the cost of total time. This is a pattern shown in many of the studies: improvements may be reached, but always at a price for something else.

### 4.2. Implications for future studies

The standards for conducting empirical studies, and which measures to evaluate differ greatly across the studies. Rothermel and Harrold proposed a framework to constitute the basis for comparison [48], but it is not used to any significant level in later research. Hence, it is not possible to conduct very strict aggregation of research results, e.g. through meta-analysis. It is however not necessarily the ultimate goal to compare specific techniques. More general concepts would be more relevant to analyze, rather than detailed implementation issues.

Examples of such concepts to evaluate are indicated in the headings of Table 9. *Applicability*: are different techniques better suited for different languages or programming concepts, or for certain types of software? *Method*: are some selection approaches better suited to find faults, independently of details in their implementation? Which level of granularity for the analysis is

effective – statement, class, component, or even specification level? Other concepts are related to process, product and resources factors [53]. *Process*: How frequent should the regression testing cycles be? At which testing level is the regression testing most efficient: unit, function, system? *Product*: Is regression testing different for different types and sizes of products? *Resources*: Is the regression testing different with different skills and knowledge among the testers?

In the reviewed studies, some of these aspects are addressed: e.g. the size aspect, scaling up from small programs to medium-sized [50], the level of granularity of the change analysis [3], as well as testing frequency [27] and the effect of changes [11]. However, this has to be conducted more systematically by the research community.

Since the outcomes of the studies depend on many different factors, replication of studies with an attempt to keep as many factors stable as possible is a means to achieve a better empirical foundation for evaluation of concepts and techniques. The use of benchmarking software and test suites is one way of keeping factors stable between studies [8]. However, in general, the strive for novelty in each research contribution tends to lead to a lack of replications and thus a lack of deeper understanding of earlier proposed techniques.

A major issue in this review is to find the relevant information to compare techniques. Hence, for the future, a more standardized documentation scheme would be helpful, as proposed by Jedlitschka and Pfahl [24] for experiments and by Runeson and Höst [52] for case studies. To allow enough detail despite page restrictions, complementary technical reports could be published on the empirical studies.

## 5. Conclusions and future work

In this paper we present results from a systematic review of empirical evaluations of regression test selection techniques. Related to our research questions we have identified that:

- (RQ1) There are 28 empirically evaluated techniques on regression test selection published.
- (RQ2) These techniques might be classified according to: applicability on type of software and type of language; details regarding the method such as which input is required, which approach is taken and on which level of granularity are changes considered; and properties such as classification in safe/unsafe or minimizing/not minimizing.
- (RQ3) The empirical evidence for differences between the techniques is not very strong, and sometimes contradictory, and
- (RQ4) hence there is no basis for selecting one superior technique. Instead techniques have to be tailored to specific situations, e.g. initially based on the classification of techniques.

We have identified some basic problems in the regression testing field which hinder a systematic review of the studies. Firstly, there is a great variance in the uniqueness of the techniques identified. Some techniques may be presented as novel at the time of their publications and others may be regarded as variants of already existing techniques. Combined with a tendency to consider replications as second class research, the case for cooperative learning on regression testing techniques is not good. In addition to this, some techniques are presented in a rather general manner, e.g. claimed to handle object-oriented programs, which gives much space for different interpretations on how they may be implemented due to, e.g. different programming language constructs existing in different programming languages. This may lead to dif-

ferent (but similar) implementations of a specific technique in different studies depending on, e.g. the programming languages used in the studies.

As mentioned in Section 1, to be able to select a strategy for regression testing, relevant empirical comparisons between different methods are required. Where such empirical comparisons exist, the quality of the evaluations must be considered. One goal of this study was to determine whether the literature on regression test selection techniques provides such uniform and rigorous base of empirical evidence on the topic that makes it possible to use it as a base for selecting a regression test selection method for a given software system.

Our study shows that most of the presented techniques are not evaluated sufficiently for a practitioner to make decisions based on research alone. In many studies, only one aspect of the problem is evaluated and the context is too specific to be easily applied directly by software developers. Few studies are replicated, and thus the possibility to draw conclusions based on variations in test context is limited. Of course even a limited evidence base could be used as guidance. In order for a practitioner to make use of these results, the study context must be considered and compared to the actual environment into which a technique is supposed to be applied.

Future work for the research community is to (1) focus more on general regression testing concepts rather than on variants of specific techniques; (2) encourage systematic replications of studies in different contexts, preferably with a focus on gradually scaling up to more complex environments; (3) define how empirical evaluations of regression test selection techniques should be reported, which variation factors in the study context are important.

## Acknowledgments

The authors acknowledge Dr. Carina Andersson for her contribution to the first two stages of the study. The authors are thankful to librarian Maria Johnsson for excellent support in the search procedures. The authors also appreciate review comments from Prof. Sebastian Elbaum and the anonymous reviewers, which substantially have improved the paper. This work is partly funded by the Swedish Governmental Agency for Innovation Systems under Grant 2005-02483 for the UPPREPA project, and partly funded by the Swedish Research Council under Grant 622-2004-552 for a senior researcher position in software engineering.

## References

- [1] H. Agrawal, J.R. Horgan, E.W. Krauser, S.A. London, Incremental regression testing, in: Proceedings of the Conference on Software Maintenance 1993. CSM-93 (Cat. No. 93CH3360-5), IEEE Comput. Soc. Press, 1993, pp. 348–357.
- [2] G. Baradhi, N. Mansour, A comparative study of five regression testing algorithms, in: Proceedings of the Software Engineering Conference, 1997, Australian, 1997, pp. 174–182.
- [3] J. Bible, G. Rothermel, D.S. Rosenblum, A comparative study of coarse- and fine-grained safe regression test-selection techniques, ACM Transactions on Software Engineering and Methodology 10 (2) (2001) 149–183.
- [4] D. Binkley, The application of program slicing to regression testing, Information and Software Technology 40 (11–12) (1998) 583–594.
- [5] P. Brereton, B.A. Kitchenham, D. Budgen, M. Turner, M. Khalil, Lessons from applying the systematic literature review process within the software engineering domain, Journal of Systems and Software 80 (4) (2007) 571–583.
- [6] Y.-F. Chen, D.S. Rosenblum, K.-P. Vo, Test tube: a system for selective regression testing, in: Proceedings of the International Conference on Software Engineering, IEEE, Los Alamitos, CA, USA, 1994, pp. 211–220.
- [7] O. Dieste, A. Grimán, N. Juristo, Developing search strategies for detecting relevant experiments, Empirical Software Engineering (2008).
- [8] H. Do, S. Elbaum, G. Rothermel, Supporting controlled experimentation with testing techniques: an infrastructure and its potential impact, Empirical Software Engineering, An International Journal 10 (4) (2005).
- [9] H. Do, G. Rothermel, An empirical study of regression testing techniques incorporating context and lifecycle factors and improved cost-benefit models, in: Proceedings of the ACM SIGSOFT Symposium on Foundations of Software Engineering, November 2006, pp. 141–151.
- [10] T. Dybå, T. Dingsöyr, G.K. Hanssen, Applying systematic reviews to diverse study types: an experience report, in: First International Symposium on Empirical Software Engineering and Measurement, 2007 (ESEM 2007), 2007, pp. 225–234.
- [11] S. Elbaum, P. Kallakuri, A. Malishevsky, G. Rothermel, S. Kanduri, Understanding the effects of changes on the cost-effectiveness of regression testing techniques, Software Testing, Verification and Reliability 13 (2) (2003) 65–83.
- [12] E. Engström, Mats Skoglund, Per Runeson, Empirical evaluations of regression test selection techniques: a systematic review, in: ESEM 08, 2008.
- [13] K. Fischer, F. Raji, A. Chruscicki, A methodology for retesting modified software, in: NTC'81, IEEE 1981 National Telecommunications Conference, Innovative Telecommunications – Key to the Future, IEEE, 1981, p. 6–3.
- [14] P.G. Frankl, G. Rothermel, K. Sayre, F.I. Vokolos, An empirical comparison of two safe regression test selection techniques, in: Proceedings of the ISESE International Symposium on Empirical Software Engineering 2003, 2003, pp. 195–204.
- [15] T.L. Graves, M.J. Harrold, J.M. Kim, A. Porter, G. Rothermel, An empirical study of regression test selection techniques, ACM Transactions on Software Engineering and Methodology 10 (2) (2001) 184–208.
- [16] R. Gupta, M.J. Harrold, M.L. Soffa, An approach to regression testing using slicing, in: Conference on Software Maintenance 1992 (Cat. No. 92CH3206-0), IEEE Comput. Soc. Press, 1992, pp. 299–308.
- [17] R. Gupta, M.J. Harrold, M.L. Soffa, Program slicing-based regression testing techniques, Software Testing, Verification and Reliability 6 (2) (1996) 83–111.
- [18] F. Haftmann, D. Kossmann, E. Lo, A framework for efficient regression tests on database applications, VLDB Journal 16 (1) (2007) 145–164.
- [19] M.J. Harrold, J.A. Jones, L. Tongyu, L. Donglin, A. Orso, M. Pennings, S. Sinha, S.A. Spoon, A. Gujarathi, Regression test selection for Java software, in: SIGPLAN Not. (USA), ACM, 2001, pp. 312–326.
- [20] M.J. Harrold, M.L. Souffa, An incremental approach to unit testing during maintenance, in: Proceedings of the Conference on Software Maintenance – 1988 (IEEE Cat. No. 88CH2615-3), IEEE Comput. Soc. Press, 1988, pp. 362–367.
- [21] J. Hartmann, D.J. Robson, Approaches to regression testing, in: Proceedings of the Conference on Software Maintenance – 1988 (IEEE Cat. No. 88CH2615-3), IEEE Comput. Soc. Press, 1988, pp. 368–372.
- [22] J. Hartmann, D.J. Robson, Techniques for selective revalidation, IEEE Software 7 (1) (1990) 31–36.
- [23] M. Hutchins, H. Foster, T. Goradia, T. Ostrand, Experiments on the effectiveness of dataflow- and control-flow-based test adequacy criteria, in: ICSE-16, 16th International Conference on Software Engineering (Cat. No. 94CH3409-0), IEEE Comput. Soc. Press, 1994, pp. 191–200.
- [24] A. Jedlitschka, D. Pfahl, Reporting guidelines for controlled experiments in software engineering, in: Proceedings of ACM/ IEEE International Symposium on Empirical Software Engineering, 2005, pp. 95–104.
- [25] N. Juristo, A.M. Moreno, S. Vegas, M. Solari, In search of what we experimentally know about unit testing [software testing], IEEE Software 23 (6) (2006) 72–80.
- [26] B. Kampenes Vigdis, T. Dybå, E. Hannay Jo, I.K. Sjöberg Dag, A systematic review of effect size in software engineering experiments, Information and Software Technology 49 (11–12) (2007) 1073.
- [27] J.-M. Kim, A. Porter, G. Rothermel, An empirical study of regression test application frequency, Software Testing, Verification, and Reliability 15 (4) (2005) 257–279.
- [28] B. Kitchenham, O.P. Brereton, D. Budgen, M. Turner, J. Bailey, S. Linkman, Systematic literature reviews in software engineering – a systematic literature review, Information and Software Technology 51 (1) (2009) 7–15.
- [29] B.A. Kitchenham, Guidelines for performing Systematic Literature reviews in Software Engineering Version 2.3. Technical Report S.o.C.S.a.M. Software Engineering Group, Keele University and Department of Computer Science University of Durham, 2007.
- [30] B.A. Kitchenham, E. Mendes, G.H. Travassos, Cross versus within-company cost estimation studies: a systematic review, IEEE Transactions on Software Engineering 33 (5) (2007) 316–329.
- [31] R.R. Klosch, P.W. Glaser, R.J. Truschneegg, A testing approach for large system portfolios in industrial environments, Journal of Systems and Software 62 (1) (2002) 11–20.
- [32] T. Koju, S. Takada, N. Doi, Regression test selection based on intermediate code for virtual machines, in: Conference on Software Maintenance, Institute of Electrical and Electronics Engineers Inc., 2003, pp. 420–429.
- [33] J.R. Landis, G.K. Gary, The measurement of observer agreement for categorical data, Biometrics 33 (1) (1977) 159–174.
- [34] H.K.N. Leung, L. White, Insights into testing and regression testing global variables, Journal of Software Maintenance. Research and Practice 2 (4) (1990) 209–222.
- [35] H.K.N. Leung, L. White, A study of integration testing and software regression at the integration level, in: In Proceedings of the Conference on Software Maintenance 1990 (Cat. No.90CH2921-5), IEEE Comput. Soc. Press, 1990, pp. 290–301.
- [36] N. Mansour, R. Bahsoon, G. Baradhi, Empirical comparison of regression test selection algorithms, The Journal of Systems and Software 57 (1) (2001) 79–90.
- [37] N. Mansour, K. El-Fakih, Natural optimization algorithms for optimal regression testing, in: Proceedings of the IEEE Computer Society's International Computer Software and Applications Conference, IEEE, Los Alamitos, CA, USA, 1997, pp. 511–514.

- [38] C. Mao, Y. Lu, Regression testing for component-based software systems by enhancing change information, in: Proceedings of the 12th Asia-Pacific Software Engineering Conference, IEEE Computer Society, 2005. 8 pp.
- [39] A.M. Memon, Using tasks to automate regression testing of GUIs, in: IASTED International Conference on Artificial Intelligence and Applications (AIA 2004), ACTA Press, 2004, pp. 477–482.
- [40] A. Orso, M.J. Harrold, D. Rosenblum, G. Rothermel, M.L. Soffa, H. Do, Using component metacontent to support the regression testing of component-based software, in: Proceedings of the IEEE International Conference on Software Maintenance (ICSM 2001), IEEE Comput. Soc., 2001, pp. 716–725.
- [41] A. Orso, S. Nanjuan, M.J. Harrold, Scaling regression testing to large software systems, in: Softw. Eng. Notes (USA), ACM, 2004, pp. 241–251.
- [42] A. Pasala, A. Bhowmick, An approach for test suite selection to validate applications on deployment of COTS upgrades, in: Proceedings of the Asia-Pacific Software Engineering Conference, APSEC, IEEE Computer Society, Los Alamitos, CA 90720-1314, United States, 2005, pp. 401–407.
- [43] H. Pei, L. Xiaolin, D.C. Kung, H. Chih-Tung, L. Liang, Y. Toyoshima, C. Chen, A technique for the selective revalidation of OO software, Journal of Software Maintenance, Research and Practice 9 (4) (1997) 217–233.
- [44] X. Ren, F. Shah, F. Tip, B.G. Ryder, O. Chesley, Chianti: a tool for change impact analysis of java programs, in: 19th Annual ACM Conference on Object-Oriented Programming, Systems, Languages, and Applications, OOPSLA'04, Association for Computing Machinery, New York, NY 10036-5701, United States, 2004, pp. 432–448.
- [45] G. Rothermel, S. Elbaum, A. Malishevsky, P. Kallakuri, B. Davia, The impact of test suite granularity on the cost-effectiveness of regression testing, in: Proceedings of the International Conference on Software Engineering, Institute of Electrical and Electronics Engineers Computer Society, 2002, pp. 130–140.
- [46] G. Rothermel, S. Elbaum, A.G. Malishevsky, P. Kallakuri, Q. Xuemei, On test suite composition and cost-effective regression testing, ACM Transactions on Software Engineering and Methodology 13 (3) (2004) 227–331.
- [47] G. Rothermel, M.J. Harrold, A safe, efficient algorithm for regression test selection, Proceedings of the Conference on Software Maintenance, 1993 (CSM-93), 1993, pp. 358–367.
- [48] G. Rothermel, M.J. Harrold, Analyzing regression test selection techniques, IEEE Transactions on Software Engineering 22 (8) (1996) 529–551.
- [49] G. Rothermel, M.J. Harrold, A safe, efficient regression test selection technique, ACM Transactions on Software Engineering and Methodology 6 (2) (1997) 173–210.
- [50] G. Rothermel, M.J. Harrold, Empirical studies of a safe regression test selection technique, IEEE Transactions on Software Engineering 24 (6) (1998) 401–419.
- [51] G. Rothermel, M.J. Harrold, J. Dedhia, Regression test selection for C++ software, Journal of Software Testing Verification and Reliability 10 (2) (2000) 77–109.
- [52] P. Runeson, M. Höst, Guidelines for conducting and reporting case study research in software engineering, Empirical Software Engineering 14 (2) (2009) 131–164.
- [53] P. Runeson, M. Skoglund, E. Engström, Test Benchmarks – what is the question? in: TestBench Workshop at International Conference on Software Testing, Verification and Validation, Lillehammer, Norway, April 2008.
- [54] A.S.M. Sajeev, B. Wibowo, Regression test selection based on version changes of components, in: Tenth Asia-Pacific Software Engineering Conference, IEEE Comput. Soc., 2003, pp. 78–85.
- [55] T. Shadish, T. Cook, D. Campbell, Experimental and Quasi-Experimental Designs – for Generalized Causal Inference, second ed., Houghton Mifflin Company, Boston, 2002. p. 623.
- [56] M. Skoglund, P. Runeson, A case study of the class firewall regression test selection technique on a large scale distributed software system, in: 2005 International Symposium on Empirical Software Engineering (IEEE Cat. No. 05EX1213), IEEE, 2005. 10 pp.
- [57] M. Staples, M. Niazi, Experiences using systematic review guidelines, The Journal of Systems and Software 80 (9) (2007) 1425–1437.
- [58] K. Toshihiko, T. Shingo, D. Norihisa, Regression test selection based on intermediate code for virtual machines, in: Proceedings International Conference on Software Maintenance ICSM 2003, IEEE Comput. Soc., 2003, pp. 420–429.
- [59] L. White, K. Abdullah, A firewall approach for the regression testing of object-oriented software, Software Quality Week (1997).
- [60] L. White, K. Jaber, B. Robinson, Utilization of extended firewall for object-oriented regression testing, in: IEEE International Conference on Software Maintenance, ICSM, IEEE Computer Society, Los Alamitos, CA 90720-1314, United States, 2005, pp. 695–698.
- [61] L. White, B. Robinson, Industrial real-time regression testing and analysis using firewalls, in: Proceedings of the 20th IEEE International Conference on Software Maintenance, 2004, pp. 18–27.
- [62] L.J. White, H.K.N. Leung, A firewall concept for both control-flow and data-flow in regression integration testing, in: Conference on Software Maintenance 1992 (Cat. No. 92CH3206-0), IEEE Comput. Soc. Press, 1992, pp. 262–271.
- [63] D. Willmor, S.M. Embury, A safe regression test selection technique for database-driven applications, in: Proceedings of the 21st IEEE International Conference on Software Maintenance, IEEE Comput. Soc. Press, 2005, pp. 421–430.
- [64] F.I. Vokolos, P.G. Frankl, Pythia: a regression test selection tool based on textual differencing, in: 3rd International Conference on Reliability, Quality and Safety of Software-Intensive Systems, IFIP TC5 WG5.4, Chapman & Hall, 1997, pp. 3–21.
- [65] F.I. Vokolos, P.G. Frankl, Empirical evaluation of the textual differencing regression testing technique, in: Proceedings of the International Conference on Software Maintenance (Cat. No. 98CB36272), IEEE Comput. Soc., 1998, pp. 44–53.
- [66] W.E. Wong, J.R. Horgan, S. London, H. Agrawal, A study of effective regression testing in practice, in: Proceedings of the Eighth International Symposium on Software Reliability Engineering (Cat. No. 97TB100170), IEEE Comput. Soc., 1997, pp. 264–274.
- [67] Y. Wu, M.-H. Chen, H.M. Kao, Regression testing on object-oriented programs, in: Proceedings 10th International Symposium on Software Reliability Engineering (Cat. No. PR00443), IEEE Comput. Soc., 1999, pp. 270–279.
- [68] C. Yanping, L.P. Robert, D.P. Sims, Specification-based regression test selection with risk analysis, in: Proceedings of the 2002 Conference of the Centre for Advanced Studies on Collaborative research, IBM Press, 2002.
- [69] R.K. Yin, in: D.J.R. Leonard Bickman (Ed.), Case Study Research – Design and Methods Applied Social Research Methods Series, vol. 5, Sage Publications, London, 2003.
- [70] J. Zheng, In regression testing selection when source code is not available, in: Proceedings of the 20th IEEE/ACM international Conference on Automated Software Engineering, ACM, 2005.
- [71] J. Zheng, B. Robinson, L. Williams, K. Smiley, An initial study of a lightweight process for change identification and regression test selection when source code is not available, in: Proceedings of the International Symposium on Software Reliability Engineering, ISSRE, IEEE Computer Society, 2005, pp. 225–234.
- [72] J. Zheng, B. Robinson, L. Williams, K. Smiley, Applying regression test selection for COTS-based applications, in: Proceedings of the International Conference on Software Engineering, Institute of Electrical and Electronics Engineers Computer Society, Piscataway, NJ 08855-1331, United States, 2006, pp. 512–521.
- [73] J. Zheng, B. Robinson, L. Williams, K. Smiley, A lightweight process for change identification and regression test selection in using COTS components, in: Proceedings of the Fifth International Conference on Commercial-off-the-Shelf (COTS)-Based Software Systems, Institute of Electrical and Electronics Engineers Computer Society, Piscataway, NJ 08855-1331, United States, 2006, pp. 137–143.