

Spectrum-Based Fault Localization in Model Transformations

JAVIER TROYA, SERGIO SEGURA, JOSE ANTONIO PAREJO, and ANTONIO RUIZ-CORTÉS, Universidad de Sevilla, Spain

Model transformations play a cornerstone role in Model-Driven Engineering (MDE), as they provide the essential mechanisms for manipulating and transforming models. The correctness of software built using MDE techniques greatly relies on the correctness of model transformations. However, it is challenging and error prone to debug them, and the situation gets more critical as the size and complexity of model transformations grow, where manual debugging is no longer possible.

Spectrum-Based Fault Localization (SBFL) uses the results of test cases and their corresponding code coverage information to estimate the likelihood of each program component (e.g., statements) of being faulty. In this article we present an approach to apply SBFL for locating the faulty rules in model transformations. We evaluate the feasibility and accuracy of the approach by comparing the effectiveness of 18 different state-of-the-art SBFL techniques at locating faults in model transformations. Evaluation results revealed that the best techniques, namely *Kulcynski2*, *Mountford*, *Ochiai*, and *Zoltar*, lead the debugger to inspect a maximum of three rules to locate the bug in around 74% of the cases. Furthermore, we compare our approach with a static approach for fault localization in model transformations, observing a clear superiority of the proposed SBFL-based method.

CCS Concepts: • Software and its engineering → Model-driven software engineering; Domain specific languages; Software testing and debugging; Functionality; Dynamic analysis;

Additional Key Words and Phrases: Model transformation, spectrum-based, fault localization, debugging, testing

ACM Reference format:

Javier Troya, Sergio Segura, Jose Antonio Parejo, and Antonio Ruiz-Cortés. 2018. Spectrum-Based Fault Localization in Model Transformations. *ACM Trans. Softw. Eng. Methodol.* 27, 3, Article 13 (September 2018), 50 pages.

<https://doi.org/10.1145/3241744>

13

1 INTRODUCTION

In Model-Driven Engineering (MDE), models are the central artifacts that describe complex systems from various viewpoints and at multiple levels of abstraction using appropriate modeling

This work has been partially supported by the European Commission (FEDER) and Spanish Government under CICYT project BELI (TIN2015-70560-R) and the Andalusian Government project COPAS (P12-TIC-1867).

Author's addresses: J. Troya, S. Segura, J. A. Parejo, and A. Ruiz-Cortés, Departamento de Lenguajes y Sistemas Informáticos - ETS de Ingeniería Informática, Avenida Reina Mercedes s/n. 41012 Sevilla; emails: {jtroya, sergiosegura, jparejo, aruiz}@us.es.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2018 Association for Computing Machinery.

1049-331X/2018/09-ART13 \$15.00

<https://doi.org/10.1145/3241744>

formalisms. Model transformations (MTs) are the cornerstone of MDE [28, 71], as they provide the essential mechanisms for manipulating and transforming models. They are an excellent compromise between strong theoretical foundations and applicability to real-world problems [71]. Most MT languages are composed of model transformation rules.¹ Each MT rule deals with the construction of part of the target model. They match input elements from the source model and generate output elements that compose the target model.

The correctness of software built using MDE techniques typically relies on the correctness of the operations executed using MTs. For this reason, it is critical in MDE to maintain and test them as it is done with source code in classical software engineering. However, checking whether the output of a MT is correct is a manual and error-prone task, which suffers the *oracle problem*. The oracle problem refers to, given an input for a system, the challenge of distinguishing the corresponding desired, correct behavior from potentially incorrect behavior [13]. To alleviate this problem in the model transformation domain, the formal specification of MTs has been proposed by the definition and use of *contracts* [14, 18, 21, 104], i.e., assertions that the execution of the MTs must satisfy. These assertions can be specified on the models resulting from the MTs, the models serving as input for the MTs, or both, and they can be tested in a black-box manner. These assertions are typically defined using the Object Constraint Language (OCL) [108].

However, even when using the assertions as oracle to test if MTs are faulty, it is still challenging to debug them and locate what parts of the MTs are wrong. The situation gets more critical as the size and complexity of MTs grow, where manual debugging is no longer possible, such as in aviation, medical data processing [106], the automotive industry [94], or embedded and cyber-physical systems [82]. Therefore, there is an increasing need to count on methods, mechanisms, and tools for debugging them.

Some works propose debugging model transformations by bringing them to a different domain such as Maude [102], DSLTrans [80], or Colored Petri Nets [110], where some specific analysis can be applied. The problem with these approaches is that the user needs to be familiar with such domains; besides, their performance and scalability can be worse than that of the original model transformation [102]. There are a few works that propose the use of contracts to debug model transformations [18, 22, 23]. Among them, the work by Burgueño et al. [18] is the closest to ours. They address the debugging of ATL model transformations based on contracts with a static approach that aims to identify the *guilty* rule, i.e., the faulty rule. It statically extracts the types appearing in the contracts as well as those of the MT rules and decides which rules are more likely to contain a bug. This is a *static* approach, since the transformation is not executed. Despite that, it achieves relatively good results on several case studies [18]. However, the effectiveness of *dynamic* approaches is an open question. Answering this question is one of the goals of this work.

Spectrum-Based Fault Localization (SBFL) is a popular technique used in software debugging for the localization of bugs [3, 115]. It uses the results of test cases and their corresponding code coverage information to estimate the likelihood of each program component (e.g., statements) of being faulty. A program spectrum details the execution information of a program from a certain perspective, such as branch or statement coverage [50]. SBFL entails identifying the part of the program whose activity correlates most with the detection of errors.

This article presents and evaluates in detail the first approach that applies spectrum-based fault localization in model transformations, extending our article with the initial ideas [99]. SBFL being a dynamic approach, our approach takes advantage of the information recovered after MT runs, which may help improve the results over static approaches [18] and at the same time complement them. We follow the approaches in References [14, 18, 21, 104] and use the previously

¹Throughout the article, we may also refer to *model transformation (MT) rules* as *transformation rules* or merely *rules*.

described contracts (assertions) as oracle to determine the correctness of MTs. Given an MT, a set of assertions, and a set of source models, our approach indicates the violated assertions and uses the information of the MT coverage to rank the transformation rules according to their suspiciousness of containing a bug. Also, of the many existing techniques proposed in the literature for computing the suspiciousness values [85, 115, 117], we select 18 of them and compare their effectiveness in the context of MTs.

There is a plethora of frameworks and languages to define MTs. Among them, The ATLAS transformation language (ATL) [61, 84] has come to prominence in the MDE community both in the academic and the industrial arenas, so the testing of ATL transformations is of prime importance. This success is due to ATL's flexibility, support of the main metamodeling standards, usability that relies on strong tool integration within the Eclipse world, and a supportive development community [80]. To implement our approach and achieve automation, we have built a prototype for debugging ATL model transformations. However, we may mention that the proposed approach is applicable to any model transformation language as long as it is able to store the execution of the transformation in traces. Therefore, the approach could be trivially applied to languages such as QVT [45], Maude [26], Kermeta [56], and many more, since in most transformation languages it is possible to define the generation of an extra target model that stores the traces (cf. Section 2.2.3).

We have thoroughly evaluated the approach using the implemented prototype. To do so, we have selected four different case studies that differ regarding the application domains, size of metamodels and transformations, and the number and types of features of ATL used. For instance, the number of rules ranges from 8 to 39 and the lines of code from 53 to 1055. We have defined 117 OCL assertions for the four case studies, many of them taken from Reference [18], and have applied mutation testing by creating 158 mutants using the operators presented in Reference [98], where each mutant is a faulty variation of the original model transformation. Experimental results reveal that the best techniques place the faulty transformation rule among the three most suspicious rules in around 74% of the cases. Looking into each of the four case studies, the best techniques allow the tester to locate the fault by inspecting only 1.59, 2.99, 2.4, and 4.8 rules in each of the case studies, which are composed of 9, 19, 8, and 39 rules, respectively. Furthermore, we compared our approach with a state-of-the-art approach based on the static analysis of transformation rules and assertions, observing a clear superiority of the proposed SBFL-based approach. The conclusions from our experiments serve as a proof of concept of the effectiveness of SBFL techniques to aid in the process of debugging model transformations.

Like ATL, our prototype is compliant with the Eclipse Modeling Framework and is completely automated and executable, dealing with Ecore metamodels and XML Metadata Interchange (XMI) model instances and tailored at iteratively debugging ATL model transformations, although it could be trivially extended to support other transformation languages based on rules.

The remainder of this article is organized as follows. Section 2 presents the basics of our approach, namely it explains metamodeling, model transformations and the ATL language, and spectrum-based fault localization. Then, Section 3 details our approach for applying SBFL in MTs and explains the proposed methodology for debugging model transformations as well as the implemented automation. It is followed by a thorough evaluation in Section 4, for which four case studies have been used. The comparison with the static approach [18] is also presented in this section. Then, Section 5 presents and describes some works related to ours, and the article finishes with the conclusions and some potential lines of future work in Section 6.

2 BACKGROUND

In this section, we present the basics to understand our approach. First, an introduction to metamodeling and an explanation of its most basic concepts are given. Then, we focus on a

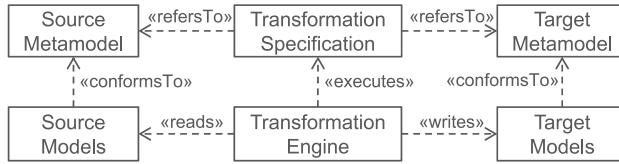


Fig. 1. Model transformation pattern (from [28]).

detailed explanation of model transformations and the ATL transformation language, followed by the introduction of the ATL MT that serves as running example. Finally, we explain the rationale behind spectrum-based fault localization.

2.1 Metamodeling

MDE [29] is a methodology that advocates the use of models as first-class entities throughout the software engineering life cycle. MDE is meant to increase productivity by maximizing compatibility between systems, simplifying the process of design and promoting communication between individuals and teams working on the system, since they can all share a high-level picture of the system.

Metamodels, models, domain-specific languages (DSLs), and model transformations are, among others, key concepts in MDE. A model is an abstraction of a system often used to replace the system under study [66, 72]. Thus, (part of) the complexity of the system that is not necessary in a certain phase of the system development is removed in the model, making it more simple to manage, understand, study, and analyze. Models are also used to share a common vision and facilitate the communication among technical and non-technical stakeholders [29].

Every model must conform to a metamodel. Indeed, a metamodel defines the structure and constraints for a family of models [75]. Like everything in MDE, a metamodel is itself a model, and it is written in the language defined by its meta-metamodel. It specifies the concepts of a language, the relationships between these concepts, the structural rules that restrict the possible elements in the valid models, and those combinations between elements with respect to the domain semantic rules.

A metamodel dictates what kind of models can be defined within a specific domain, i.e., it defines the abstract syntax of a DSL. The concrete syntax of DSLs can be defined in several ways, normally either graphically or textually. To provide a DSL with semantics and behavior, its defining metamodel may not be enough. Therefore, apart from its concrete and abstract syntaxes, also its semantics may need to be defined. For instance, model transformations can be used to give semantics to a DSL by translating it to a different domain where further analysis, simulations, and so on can be performed [103]. This mechanism enables the definition of flexible and reusable DSLs, where several kinds of analysis can be defined [32, 76].

2.2 Model Transformations

Model transformations play a cornerstone role in MDE, since they provide the essential mechanisms for manipulating and transforming models [16, 96]. They allow querying, synthesizing, and transforming models into other models or into code, so they are essential for building systems in MDE. A model transformation is a program executed by a transformation engine that takes one or more input models and produces one or more output models, as illustrated by the model transformation pattern [28] in Figure 1.² Model transformations are developed on the metamodel level, so they are reusable for all valid model instances. Most MT languages are composed of model transformation rules, where each rule deals with the construction of part of the target

²In the article, we use the terms *input/output* models/metamodels and *source/target* models/metamodels indistinctly.

model. They match input elements from the source model and generate output elements that compose the target model.

There is a plethora of frameworks and languages to define MTs, such as Henshin [10], AGG [97], Maude [26], AToM³ [30], e-Motions [86], VIATRA [27], MOMoT [35–37], QVT [45], Kermeta [56], JTL [24], and ATL [62]. In most of these frameworks and languages, model transformations are composed of transformation rules. Among them, we focus in this article on the ATL language due to its importance in both the academic and the industrial arenas.

2.2.1 *ATL* Transformation Language. ATL has come to prominence in the model-driven engineering community due to its flexibility, support of the main meta-modeling standards, usability that relies on strong tool integration with the Eclipse world, and a supportive development community [61, 84].

ATL is a textual rule-based model transformation language that provides both declarative and imperative language concepts. It is thus considered a hybrid model transformation language. An ATL transformation is composed of a set of transformation rules and helpers.³ Each rule describes how certain output model elements should be generated from certain input model elements. Declarative rules are called *matched rules*, while (*unique*) *lazy* and *called* rules are invoked from other rules. Rules are mainly composed of an *input pattern* and an *output pattern*. The input pattern is used to match *input pattern elements* that are relevant for the rule. The output pattern specifies how the *output pattern elements* are created from the input model elements matched by the input pattern. Each output pattern element can have several *bindings* that are used to initialize its attributes and references.

Methods in the ATL context are called *helpers*. There exist two different, although very similar from their syntax, kinds of helpers: the functional and the attribute helpers. Both can be defined in the context of a given data type, and functional helpers can accept parameters, while attribute helpers cannot. Functional helpers make it possible to define factorized ATL code that can then be called from different points of an ATL program. Attribute helpers, in turn, can be viewed as constants.

2.2.2 *Transformation Example*. The *BibTeX2DocBook* model transformation [54], taken from the open-access repository known as *ATL Transformation Zoo* [12], is used throughout this article as running example. It transforms a BibTeXXML model to a DocBook composed document. BibTeXXML⁴ is an XML-based format for the BibTeX bibliographic tool. DocBook [107] is an XML-based format for document composition.

The aim of this transformation is to generate, from a BibTeXXML file, a DocBook document that presents the different entries of the bibliographic file within four different sections. The first and second sections provide the full list of bibliographic entries and the sorted list of the different authors referenced in the bibliography, respectively, while the third and last sections present the titles of the bibliography titled entries (in a sorted way) and the list of referenced journals (in article entries), respectively.

The metamodels of this transformation are displayed in Figure 2. The BibTeXXML metamodel (Figure 2(a)) deals with the mandatory fields of each BibTeX entry (for instance, author, year, title, and journal for an article entry). A bibliography is modeled by a *BibTeXFile* element. This element is composed of *entries* that are each associated with an *id*. All entries inherit, directly or indirectly, from the abstract *BibTeXEntry* element. The abstract classes *AuthoredEntry*, *DatedEntry*, *TitledEntry*, and *BookTitledEntry*, as well as the *Misc* entry, directly inherit from *BibTeXEntry*.

³https://wiki.eclipse.org/ATL/User_Guide_-_The_ATL_Language.

⁴<http://bibtexml.sourceforge.net/>.

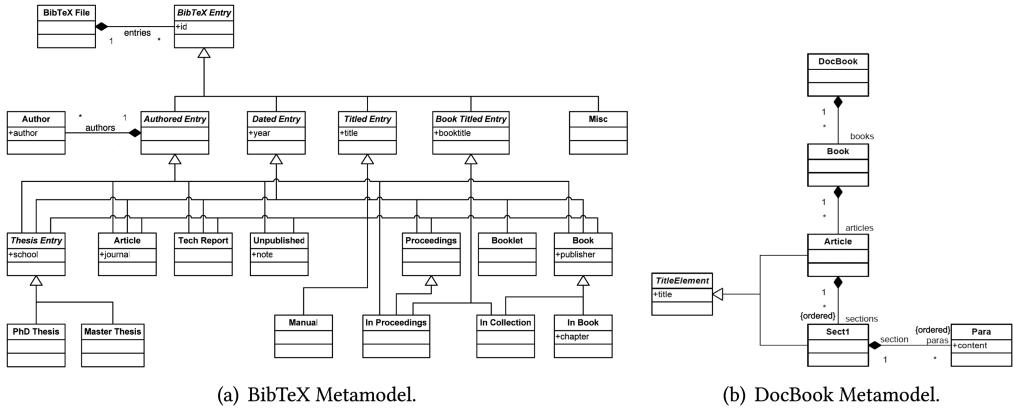


Fig. 2. Metamodels of the BibTeX2DocBook transformation (from Reference [54]).

Concrete BibTeX entries inherit from some of these abstract classes according to their set of mandatory fields. There are 13 possible entry types: *PhDThesis*, *MasterThesis*, *Article*, *TechReport*, *Unpublished*, *Manual*, *InProceedings*, *Proceedings*, *Booklet*, *InCollection*, *Book*, *InBook*, and *Misc*. An authored entry may have several authors.

The DocBook metamodel (Figure 2(b)) represents a limited subset of the DocBook definition. Within this metamodel, a DocBook document is associated with a *DocBook* element. Such an element is composed of several *Books* that, in turn, are composed of several *Articles*. An *Article* is composed of sections (class named *Sect1*) that are ordered. A *Sect1* is composed of paragraphs (class *Para*) that are also ordered within each section. Both *Article* and *Sect1* inherit from the *TitledElement* abstract class.

The *BibTeX2DocBook* model transformation [54] is shown in Listing 1, which contains nine rules. We may mention that the transformation is shown here in a “compressed” way in order not to occupy too much space, since, normally, line breaks are used when, for instance, adding a new binding. The first rule, *Main*, creates the structure of a *DocBook* from a *BibTeXFile* and creates four sections with their corresponding titles. The paragraphs of each section are to be resolved when the remaining rules are executed. This rule uses the helpers *authorSet*, *titledEntrySet*, and *articleSet*. They return, respectively, the sequence of distinct authors (with unique names), *TitledEntries* (with unique titles), and *Articles* (with unique journal names) referenced in the input BibTeX model.

The second rule, *Author*, creates a paragraph for each author and sets as content the author name. The third one creates a paragraph for each untitled entry and uses helper *buildEntryPara* to set its content. This helper builds a string containing all information of a given *BibTeXEntry*. The fourth rule, *TitledEntry_Title_NoArticle*, creates two paragraphs for each *TitledEntry* that is not an article and that is included in the set of *TitledEntry* with unique titles (helper *titledEntrySet*). The next one, *TitledEntry_NoTitle_NoArticle*, creates a paragraph for each *TitledEntry* that is not an article and is not included in the set of *TitledEntry* with unique titles. The next two rules, *Article_Title_Journal* and *Article_NoTitle_Journal*, create paragraphs for those articles whose title is either included in the set of *TitledEntry* with unique titles or not, respectively. Also, the article must be included in the set of *Articles* whose journal name is unique (helper *articleSet*). Finally, the eighth and ninth rules, *Article_Title_NoJournal* and *Article_NoTitle_NoJournal*, create paragraphs for those articles whose title is either included in the set of *TitledEntry* with unique titles or not, respectively. Also, the article must not be included in the set of *Articles* whose journal name is unique (helper *articleSet*). We refer the interested reader to the document explaining the complete model transformation [54].

```

1 module BibTeX2DocBook;
2 create OUT : DocBook from IN : BibTeX;
3
4 helper def: authorSet : Sequence(BibTeX!Author) =
5   BibTeX!Author.allInstances()>iterate(e; ret : Sequence(BibTeX!Author) = Sequence {} | 
6     if ret->collect(e | e.author)->includes(e.author) then ret else ret->including(e)
7     endif)>sortedBy(e | e.author);
8
9 helper def: titledEntrySet : Sequence(BibTeX!TitledEntry) =
10  BibTeX!TitledEntry.allInstances()>iterate(e; ret : Sequence(BibTeX!TitledEntry) = 
11    Sequence {} | 
12      if ret->collect(e | e.title)->includes(e.title) then ret else ret->including(e)
13      endif)>sortedBy(e | e.title);
14
15 helper def: articleSet : Sequence(BibTeX!Article) =
16   BibTeX!Article.allInstances()>iterate(e; ret : Sequence(BibTeX!Article) = Sequence {} | 
17     if ret->collect(e | e.journal)->includes(e.journal) then ret else ret->including(e)
18     endif)>sortedBy(e | e.journal);
19
20 helper context BibTeX!BibTeXEntry def: buildEntryPara() : String =
21 '[' + self.id + ']'
22 + '_'+self.oclType().name
23 + (if self.oclisKindOf(BibTeX!TitledEntry) then '_'+self.title else ''endif)
24 + (if self.oclisKindOf(BibTeX!AuthoredEntry)
25   then self.authors->iterate(e; str : String = '' | str + '_'+e.author) else ''endif)
26 + (if self.oclisKindOf(BibTeX!DatedEntry) then '_'+self.year else ''endif)
27 + (if self.oclisKindOf(BibTeX!BookTitledEntry) then '_'+self.booktitle else ''endif)
28 + (if self.oclisKindOf(BibTeX!ThesisEntry) then '_'+self.school else ''endif)
29 + (if self.oclisKindOf(BibTeX!Article) then '_'+self.journal else ''endif)
30 + (if self.oclisKindOf(BibTeX!Unpublished) then '_'+self.note else ''endif)
31 + (if self.oclisKindOf(BibTeX!Book) then '_'+self.publisher else ''endif)
32 + (if self.oclisKindOf(BibTeX!InBook) then '_'+self.chapter.toString() else''endif);
33
34 rule Main {                                     -- tr1
35   from
36     bib : BibTeX!BibTeXFile
37   to
38     doc : DocBook!DocBook (books <- boo),
39     boo : DocBook!Book (articles <- art),
40     articles : DocBook!Article (title <- 'BibTeXXML_to_DocBook',
41       sections <- Sequence{se1, se2, se3, se4}),
42     se1 : DocBook!Sect1 (title <- 'References_List',
43       paras <- BibTeX!BibTeXEntry.allInstances()>sortedBy(e | e.id)),
44     se2 : DocBook!Sect1 (title <- 'Authors_List',
45       paras <- thisModule.authorSet),
46     se3 : DocBook!Sect1 (title <- 'Titles_List',
47       paras <- thisModule.titledEntrySet->collect(e | thisModule.resolveTemp(e, '
48         title_para'))),
49     se4 : DocBook!Sect1 (title <- 'Journals_List',
50       paras <- thisModule.articleSet->collect(e | thisModule.resolveTemp(e, '
51         journal_para'))))
52 }
53
54 rule Author {                                -- tr2
55   from
56     a : BibTeX!Author (thisModule.authorSet->includes(a))
57   to
58     p1 : DocBook!Para (content <- a.author)
59 }
60
61 rule UntitledEntry {                         -- tr3
62   from
63     e : BibTeX!BibTeXEntry (not e.oclisKindOf(BibTeX!TitledEntry))
64   to
65     p : DocBook!Para (content <- e.buildEntryPara())
66
67 rule TitledEntry_Title_NoArticle { -- tr4
68   from

```

```

68     e : BibTeX!TitledEntry (thisModule.titledEntrySet->includes(e) and
69                         not e.oclIsKindOf(BibTeX!Article))
70   to
71     entry_para : DocBook!Para (content <- e.buildEntryPara()),
72     title_para : DocBook!Para (content <- e.title)
73 }
74
75 rule TitledEntry_NoTitle_NoArticle { -- tr5
76   from
77     e : BibTeX!TitledEntry (not thisModule.titledEntrySet->includes(e) and
78                         not e.oclIsKindOf(BibTeX!Article))
79   to
80     entry_para : DocBook!Para (content <- e.buildEntryPara())
81 }
82
83 rule Article_Title_Journal { -- tr6
84   from
85     e : BibTeX!Article (thisModule.titledEntrySet->includes(e) and
86                         thisModule.articleSet->includes(e))
87   to
88     entry_para : DocBook!Para (content <- e.buildEntryPara()),
89     title_para : DocBook!Para (content <- e.title),
90     journal_para : DocBook!Para (content <- e.journal)
91 }
92
93 rule Article_NoTitle_Journal { -- tr7
94   from
95     e : BibTeX!Article (not thisModule.titledEntrySet->includes(e) and
96                         thisModule.articleSet->includes(e))
97   to
98     entry_para : DocBook!Para (content <- e.buildEntryPara()),
99     journal_para : DocBook!Para (content <- e.journal)
100 }
101
102 rule Article_Title_NoJournal { -- tr8
103   from
104     e : BibTeX!Article (thisModule.titledEntrySet->includes(e) and
105                         not thisModule.articleSet->includes(e))
106   to
107     entry_para : DocBook!Para (content <- e.buildEntryPara()),
108     title_para : DocBook!Para (content <- e.title)
109 }
110
111 rule Article_NoTitle_NoJournal { -- tr9
112   from
113     e : BibTeX!Article (not thisModule.titledEntrySet->includes(e) and
114                         not thisModule.articleSet->includes(e))
115   to
116     entry_para : DocBook!Para (content <- e.buildEntryPara())
117 }

```

Listing 1. *BibTeX2DocBook* MT.

2.2.3 ATL Internal Traces Mechanism. The ATL engine works in two steps. First, all elements are created. Second, their features are initialized. This second phase implies to resolve the corresponding references. For instance, in the transformation shown in Listing 1, line 45 initializes the *paras* reference of the *Sect1* element created. This reference will actually point elements that are created in the first phase by rule *Author*, as we explain with an example later.

To resolve these references, ATL uses an internal tracing mechanism. Thereby, every time a rule is executed, it creates a new trace and stores it in the internal trace model. A trace model can be automatically obtained from a transformation execution, e.g., by using Jouault's *TraceAdder* [60], and is composed of a set of traces, one for each rule execution. In our approach, we obtain trace models that conform to the metamodel displayed in Figure 3(a). A trace captures the name of the applied rule and the elements instantiating classes of the source metamodel (*sourceElems* reference) that are used to create new elements that instantiate classes in the target metamodel (*targetElems*

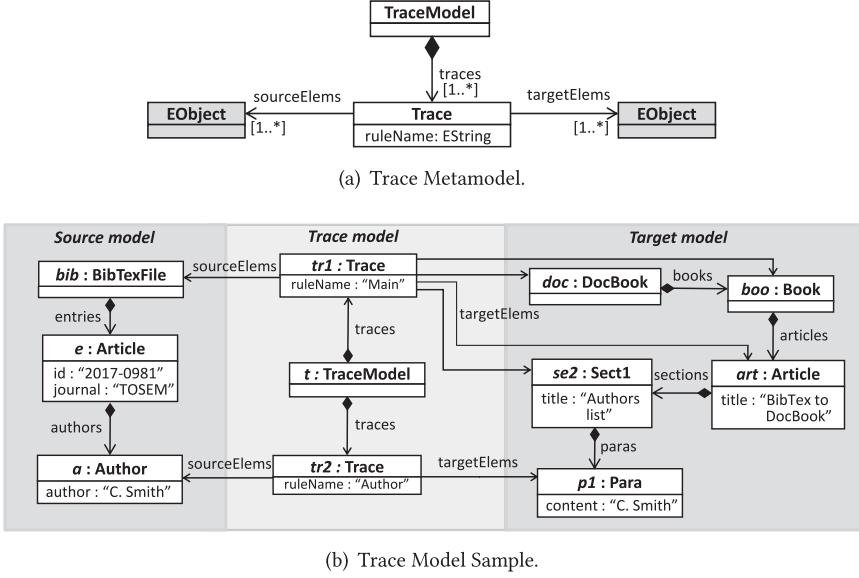


Fig. 3. Traces in model transformations.

reference). The elements pointed by such references are represented with *EObject*, because, when the metamodel is instantiated in a specific trace model, they can be any element of the source and target models, respectively. The execution of both imperative—*(unique)* *lazy* and *called*—and declarative—*matched*—rules are stored in the traces. This means that we have three models (the source model, the target model and the trace model) linked by several so-called inter-model references. Therefore, by navigating the traces, the ATL engine obtains information of which target element(s) have been created from which source element(s) and by which rule.

An example that reflects the information stored in a trace model is displayed in Figure 3(b). In the left-hand side of the figure we can see a sample source model composed of three elements. In the right-hand side we have the target model obtained after transforming elements *bib* and *a*—to keep the figure simple, we do not display the transformation of element *e*. The part in the middle of the figure represents the trace model. Since two different elements, *bib* and *a*, are transformed by two different rules, we have two traces, *tr1* and *tr2*. The first one, created by the execution of rule *Main*, records the generation of elements *doc*, *boo*, *se2* and *art* from element *bib*; while *tr2* stores the generation of *p1* from *a* by rule *Author*.

The interesting aspect in this figure is the *paras* reference between *se2* and *p1* in the target model, created using the traces. The process how ATL resolves such association is the following. As mentioned, after creating all target elements in the first phase, it resolves the references in the second phase. In our example, it means resolving the binding in line 45, where helper *authorSet* returns all *Authors* in the source model. Therefore, such binding is expressing that the *paras* reference from element *se2* in the target model should point all elements of type *Author* in the source model, *a* in our case. Of course, target elements cannot point source elements, so the ATL engine searches in the traces to recover the target elements created from *Author* elements. In our example, it recovers element *p1*, created from *Author* *a*, by inspecting trace *tr2*.

When ATL resolves the references, it recovers the first target element created by the corresponding rule, i.e., the first one that is specified in the rule (right after the *to* part of the rule). For instance, element *entry_para* (line 71 in Listing 1) would be the one recovered when

Table 1. An Example Showing the Suspiciousness Value Computed Using Tarantula (Taken from Reference [115])

Statement	Code	tc_1	tc_2	tc_3	N_{CF}	N_{CS}	Susp	Rank
s_1	input(a)	•	•	•	1	2	0.5	3
s_2	i = 1;	•	•	•	1	2	0.5	3
s_3	sum = 0;	•	•	•	1	2	0.5	3
s_4	product = 1;	•	•	•	1	2	0.5	3
s_5	if (i < a) {	•	•	•	1	2	0.5	3
s_6	sum = a + i;			•	1	0	1	1
s_7	product = a * i; // BUG: 2i → i			•	1	0	1	1
s_8	} else {	•	•		0	2	0	10
s_9	sum = a - i;	•	•		0	2	0	10
s_{10}	product = a / i;	•	•		0	2	0	10
s_{11}	}	•	•		0	2	0	10
s_{12}	print(sum);	•	•	•	1	2	0.5	3
s_{13}	print(product);	•	•	•	1	2	0.5	3
Execution Results		S	S	F				

resolving a binding reference with rule *TitledEntry_Title_NoArticle*. To specify a different element to recover when resolving the references, ATL provides the *resolveTemp* function. For instance, the *resolveTemp* function used in line 47 makes the engine retrieve the target element identified with the string “*title_para*,” i.e., the one created in line 72.

As we explain in Section 3, having this trace information is key in our approach, where we are interested only in the information of the rules that have been fired to apply our SBFL approach.

Please note that the availability of such a simple trace model is useful in many different model transformations languages, what also increases the applicability of our approach beyond ATL. For instance, Falleri et al. [33] propose a simple trace metamodel for Kermeta model transformations, and Anastasakis et al. [8] simply require the link between source and target models in an Alloy transformation. Rose et al. [87] mention the Fujaba and the MOLA (graphical transformation language developed at the University of Latvia) traceability associations, similar to the one we use, and Troya and Vallecillo [102] apply the same trace metamodel to represent model transformations in Maude. Due to the importance of trace models, Jiménez et al. [123] propose a toolkit that allows not only the definition of model transformations but also supports trace generation.

2.3 Spectrum-Based Fault Localization

SBFL uses the results of test cases and their corresponding code coverage information to estimate the likelihood of each program component (e.g., statements) of being faulty. A program spectrum details the execution information of a program from a certain perspective, such as branch or statement coverage [50]. Table 1 depicts an example showing how the technique is applied to a sample program [115]. This programs receives a natural number, a . If it is bigger than 1, then the program must print the result of adding 1 to such number as well as its double. Otherwise, it must print the number minus 1 as well as the number itself.

Having a look at the table, it horizontally shows the code statements of the program, i.e., its components. Note that a bug is seeded in statement s_7 , so that it does not multiply the number by 2. Also note that SBFL considers all lines as statements, so, for instance, the line containing only the character that closes a branch, “},” conforms statement s_{11} . However, statement s_5 includes a condition as well as the opening of a branch with character “{.” Therefore, the way of writing

a program may have an impact in the results returned by SBFL techniques. Vertically, the table shows three test cases of the program. For each test case (i.e., tc_1 , tc_2 , and tc_3), a cell is marked with “•” if the program statement of the row has been exercised by the test case of the column, creating what is known as *coverage matrix* [4]. Additionally, the final row depicts the outcome of each test case, either “Successful” or “Failed,” conforming the so-called *error vector* [4]. Based on this information, it is possible to identify which components were involved in a failure (and which ones were not), narrowing the search for the faulty component that made the execution fail.

Once a coverage matrix and an error vector as those shown in Table 1 are constructed, a number of techniques can be used to rank the program components according to their *suspiciousness*, that is, their probability of containing a fault. For instance, a popular fault localization technique is *Tarantula* [59], which for a program statement is computed as $(N_{CF}/N_F)/(N_{CF}/N_F + N_{CS}/N_S)$, where N_{CF} is the number of failing test cases that cover the statement, N_F is the total number of failing test cases, N_{CS} is the number of successful test cases that cover the statement, and N_S is the total number of successful test cases. The suspiciousness score of each statement is in the range $[0,1]$, i.e., the higher the suspiciousness score of each component, the higher the probability of having a fault. The values of N_{CF} , N_{CS} , and the *Tarantula* suspiciousness value for each statement are given in the sixth, seventh, and eighth columns of Table 1, respectively. Let us focus for instance in the row for statement s_4 . N_{CF} is 1, because only the failing test case tc_3 covers the statement. Then, N_{CS} is 2, because both tc_1 and tc_2 cover the statement, and they are successful test cases. By applying the formula, we get a value of 0.5 for suspiciousness. Finally, the last column indicates the position of the statement in the suspiciousness-based ranking where top-ranked statements are more likely to be faulty. In the example, the faulty statement s_7 is ranked first.

The effectiveness of suspiciousness metrics is usually measured using the EXAM score [117, 120], which is the percentage of statements in a program that has to be examined until the first faulty statement is reached, i.e.,

$$\text{EXAM}_{\text{Score}} = \frac{\text{Number of statements examined}}{\text{Total number of statements}}.$$

It is noteworthy that suspiciousness techniques may often provide the same value for different statements, as they are tied for the same position in the ranking, e.g., statements s_6 and s_7 in Table 1. To break ties, different approaches are applicable, such as measuring the effectiveness in the best-, average-, and worst-case scenarios, using an additional technique to break the tie or using some simple heuristics such as alphabetical ordering [115]. In the best-case scenario, the faulty statement is inspected first in the tie. Conversely, the worst-case scenario is the one where it is inspected last.

In our example, assuming that the statement s_7 is examined in second place (worst-case scenario), the EXAM score of *Tarantula* in the previous example would be $\frac{2}{13} = 0.153$, i.e., 15.3% of the statements must be examined to locate the bug.

The values that the EXAM score can have depend on the number of statements of the program under test, which goes in the denominator of the formula. In the example, the best EXAM score for a statement would be $\frac{1}{13} = 0.0769$. This EXAM score indicates that the buggy statement should be examined first. On the contrary, the worst EXAM value is always 1. In the example, if a statement is to be inspected last, then it has the EXAM score $\frac{13}{13} = 1$. Therefore, the set of values for the EXAM score, from best to worst, is $\{\frac{1}{\text{num_statements}}, \frac{2}{\text{num_statements}}, \dots, \frac{\text{num_statements}}{\text{num_statements}}\}$.

3 SPECTRUM-BASED FAULT LOCALIZATION IN MODEL TRANSFORMATIONS

In this section, we describe our SBFL approach for debugging model transformations. We first describe how the coverage matrix and the error vector are constructed. This is followed by an explanation of the suspiciousness calculation of the different transformation rules and the metric

Table 2. Tarantula [59] Suspiciousness Values for the Simplified *BibTeX2DocBook* MT When OCL_2 Fails

T. Rule	tc_{02}	tc_{12}	tc_{22}	tc_{32}	tc_{42}	tc_{52}	tc_{62}	tc_{72}	tc_{82}	tc_{92}	N_{CF}	N_{UF}	N_{CS}	N_{US}	N_C	N_U	Susp	Rank
tr_1	•	•	•	•	•	•	•	•	•	•	9	0	1	0	10	0	0.5	3
tr_2 (BUG)	•	•	•	•	•	•	•	•	•	•	9	0	0	1	9	1	1	1
tr_3	•	•	•	•	•	•	•	•	•	•	7	2	1	0	8	2	0.44	7
tr_4	•	•	•	•	•	•	•	•	•	•	9	0	1	0	10	0	0.5	3
tr_5	•	•	•	•	•	•	•	•	•	•	8	1	1	0	9	1	0.47	6
tr_6	•	•	•	•	•	•	•	•	•	•	9	0	1	0	10	0	0.5	3
tr_7	•						•				2	7	0	1	2	8	1	1
tr_8	•	•		•		•	•	•			5	4	1	0	6	4	0.36	8
tr_9		•				•	•			2	7	1	0	3	7	0.18	9	
Test Result	F	S	F	F	F	F	F	F	F									

used for measuring the effectiveness of SBFL techniques. Then we describe the methodology to apply our approach. The section ends with an explanation of the implementation and automation of our approach.

3.1 Constructing the Coverage Matrix and Error Vector

The construction of the coverage matrix requires information about the execution of the MT with a set of source models: $S = \{s_1, s_2, \dots, s_n\}$. These models must conform to the source metamodel. The oracle that determines whether the result of a MT is correct or not is a set of OCL assertions: $O = \{ocl_1, ocl_2, \dots, ocl_m\}$. These assertions are defined by specifying the expected properties of the output models of the transformation or properties that the \langle input, output \rangle model pairs must satisfy. As an example, Listing 2 shows three OCL assertions for the model transformation described in Section 2.2.2, where classes of the source and target metamodels have the prefixes *Src* and *Trg*, respectively. We consider a *test case* as a pair composed of a source model and an OCL assertion: $tc_{ij} = \langle s_i, ocl_j \rangle$. Therefore, the test suite is composed of the Cartesian product of source models and OCL assertions: $T = S \times O = \{tc_{11}, tc_{12}, \dots, tc_{nm}\}$. The test case tc_{ij} produces an error if the result of executing the MT with the source model s_i does not satisfy the OCL assertion ocl_j . It is worth noting that OCL assertions must hold for any source model. Therefore, an OCL assertion is not satisfied for a MT when there is, at least, one test case where it is violated. This means, for instance, that for ocl_1 to be satisfied, it must be satisfied in $\{tc_{11}, tc_{21}, \dots, tc_{n1}\}$.

We may recall that this article focuses on debugging and not testing. Thus, we do not impose any constraint on how the source models are generated, either manually or automatically, or on the type of OCL assertions used. However, please note that the effectiveness of SBFL is directly related to the design of the test cases. For instance, having only one test case is useless, since no coverage matrix can be created. Besides, the source models should have a good coverage of the MT and, at the same time, be diverse. This way, different source models will likely fire different parts of the model transformation, and together they will exercise all rules and will produce a useful spectra.

Table 2 depicts a sample coverage matrix constructed using our approach. Horizontally, the table shows the transformation rules in which we aim to locate bugs. In particular, we list the nine transformation rules $\langle tr_1, tr_2, \dots, tr_9 \rangle$ of the *BibTeX2DocBook* example [54], where a bug has been seeded in tr_2 . Vertically, the table shows 10 test cases aiming to check the correctness of constraint OCL_2 in Listing 2, $\langle tc_{02}, tc_{12}, \dots, tc_{92} \rangle$. A cell is marked with “•” if the transformation rule of the row has been exercised by the test case of the column. The information about the rules triggered by a given test case can be collected by inspecting the trace model, e.g., using Jouault’s

```

1 --OCL1. The main Article must be properly created and named
2 TrgBook.allInstances()->forAll(b|b.articles->exists(a|a.title='BibTeXML_to_DocBook'))
3 --OCL2. For each author, there must be a paragraph with the name of the author within a
        section named 'Authors List'
4 SrcAuthor.allInstances()->forAll(a|TrgPara.allInstances()->exists(p|p.content=a.author and
        p.section.title='Authors_List')
5 --OCL3. The titles of all publications must appear in the content of a paragraph of a
        section
6 SrcTitledEntry.allInstances()->forAll(te|TrgSect1.allInstances()->exists(s|s.paras->exists
        (p|p.content=te.title)))
7 --OCL4. There must be the same number of BibTeXFile and DocBook instances
8 SrcBibTeXFile.allInstances()->size()=TrgDocBook.allInstances()->size()

```

Listing 2. Sample OCL assertions for the *BibTeX2DocBook* MT.

TraceAdder [60] (cf. Section 2.2.3). The final row depicts the error vector with the outcome of each test case, either successful (“S”) or failed (“F”). In the example, all test cases fail except tc_{12} , i.e., OCL_2 is violated.

Note that by grouping those test cases using the same OCL assertion we can simplify debugging by providing not only the most suspicious transformation rules but also the specific assertion revealing the failure. This is very useful for the user of our approach, since (s)he can focus on the non-satisfied assertion to repair the model transformation when the faulty rule is found. In practice, this means that our approach needs to generate a coverage matrix and an error vector for each violated OCL assertion, since each of them is dealt with independently from the others.

3.2 Calculating Suspiciousness

The following notation will be used throughout the article and is used in our implementation to compute the suspiciousness of transformation rules from the information collected in the coverage matrix and the error vector. This notation is directly translated from the context of SBFL in software programs[115] by using transformation rules as the components (e.g., instead of statements), namely:

N_{CF}	number of failed test cases that cover a rule
N_{UF}	number of failed test cases that do not cover a rule
N_{CS}	number of successful test cases that cover a rule
N_{US}	number of successful test cases that do not cover a rule
N_C	total number of test cases that cover a rule
N_U	total number of test cases that do not cover a rule
N_S	total number of successful test cases
N_F	total number of failed test cases

Table 2 shows the values of N_{CF} , N_{UF} , N_{CS} , N_{US} , N_C , and N_U for each transformation rule. The values of N_S and N_F are the same for all the rows, 9 and 1, respectively, and are omitted. Based on this information, the suspiciousness of each transformation rule using *Tarantula* is depicted in the column “Susp,” followed by the position of each rule in the suspiciousness-based ranking. In the example, the faulty rule tr_2 is ranked first, tied with tr_7 . Assuming that the faulty rule was inspected in the second place (worst-case scenario), the EXAM score would be calculated as $\frac{2}{9} = 0.222\%$, i.e., 22.2% of the transformation rules need to be examined to locate the bug.

3.3 Methodology

In this section, we describe the proposed methodology to help developers debug model transformations by using our approach based on spectrum-based fault localization. It is graphically depicted in Figure 4.

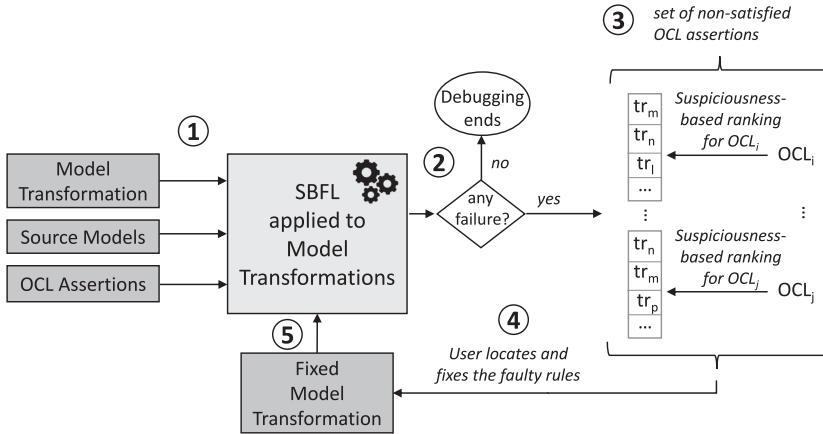


Fig. 4. Debugging of a MT applying our SBFL approach.

- (1) The inputs have to be provided, namely the *Model Transformation* under test as well as the sets of *Source Models* and *OCL Assertions*.
- (2) The approach executes and indicates whether there is *any failure*, ending the process if there is none.
- (3) If there is a failure, then it indicates the *set of non-satisfied OCL assertions*, i.e., those that are violated for at least one test case. As explained in Section 3.1, it constructs a coverage matrix and an error vector for each non-satisfied assertion and returns the *suspiciousness-based rankings* in each case.
- (4) The user picks the ranking of one of the OCL assertions to *locate and fix the faulty rule* that made the assertion fail. As described in Section 4.2.5, we study the effectiveness of 18 SBFL techniques. The idea is to use the ranking of the best techniques, which are discovered in Sections 4.3 and 4.4.
- (5) Now, the user has a *Fixed Model Transformation* that has potentially less bugs than the original *Model Transformation*. The user can decide whether to use it as input for the approach, together with the *Source Models* and *OCL Assertions*, or to keep repairing it according to the suspiciousness rankings obtained for the remaining non-satisfied OCL assertions.
- (6) In the upcoming execution of the approach with the *Fixed Model Transformation*, fewer OCL assertions should be violated, and the user would repeat the process to keep fixing the bugs. This process is repeated iteratively until all bugs have been fixed.

3.4 Implementation and Automation

Our approach is supported by a toolkit. It has been implemented for debugging ATL model transformations. Within one run, it executes the MT with all the input models, checks which assertions are violated, and returns the suspiciousness-based rankings for the violated assertions together with the corresponding coverage matrices and error vectors. Additionally, if we indicate as input the faulty rules, the approach also returns the EXAM score of the results. This is possible thanks to a Java program from which ATL transformations can be triggered, indicating their inputs and doing any post-processing with the outputs. In this section, we describe the implemented tasks used for automating and orchestrating all this process.

The overview of the implementation and automation of our approach is depicted in Figure 5. As we can see, it consists of six steps, which are explained in the following:

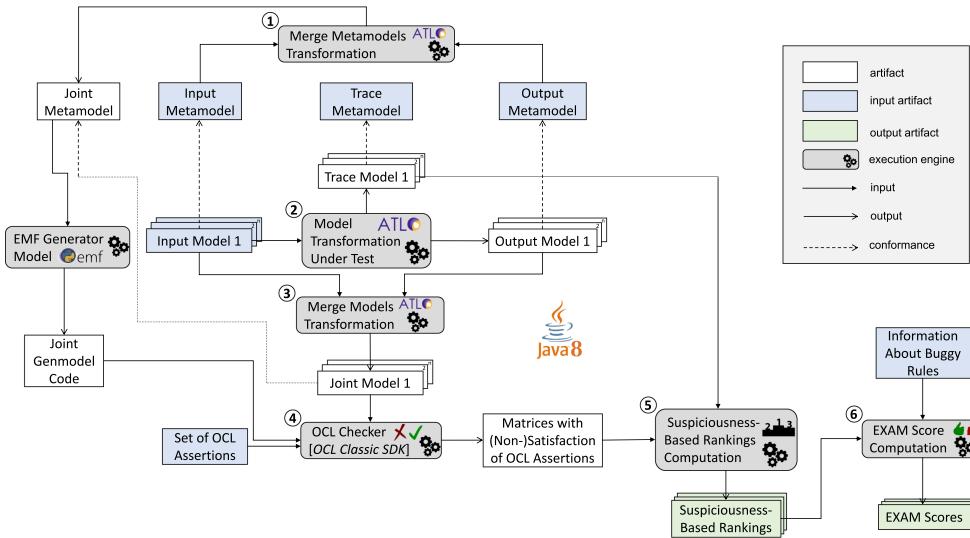


Fig. 5. Implementation and automation of our approach.

- (1) The tool of which we have made use for checking the satisfaction of the OCL assertions is *OCL Classic SDK: Ecore/UML Parsers, Evaluator, Edit*,⁵ which is part of the *Eclipse Modeling Tools*. With this tool, we can check the satisfaction of OCL assertions of a given model conforming to a metamodel. However, in our approach, the OCL assertions are typically defined on both metamodels, namely the input and output metamodels. For this reason, we need to merge both metamodels into one and do the same thing for the <input, output> model pairs. Therefore, the first step in our approach takes the *Input Metamodel* and *Output Metamodel* as input and, with the *Merge Metamodels Transformation*, it creates a *Joint Metamodel*. Due to the possibility of having classes with the same name in the input and output metamodels, this transformation puts the prefixes “Src” and “Trg” in all classes of the input and output metamodels, respectively.

Besides, the *OCL Checker* requires the Java code of the *Joint Metamodel*. This can be generated out of the box by the *EMF Generator Model*, so this is included in this first step.

- (2) The next step in our approach is to run the *ATL Transformation Under Test* with all the *Input Models* to generate the corresponding *Output Models*. Our Java program orchestrates all these model transformation executions.
- (3) For the same reason as explained in the first step, we need to merge the input and output models into the so-called *Joint Models*. These models must conform to the *Joint Metamodel* obtained in the first step. The *Merge Models Transformation* generates all the *Joint Models* for all the <*InputModel_i*, *OutputModel_i*> pairs.
- (4) The next step is to check the *Set of OCL Assertions*. This must be done for all the *Joint Models* constructed after the executions of the model transformation. As explained in the first step, we need for this the Java code obtained from the *Joint Metamodel*. This step produces as output information about the satisfiability of the OCL assertions, captured in the figure as *Matrices with (Non-)Satisfaction of OCL Assertions*. This is different from the coverage matrices explained before, since the purpose now is to identify those OCL

⁵<https://eclipse.org/modeling/mdt/downloads/?project=ocl>.

assertions that fail for at least one test case, so that coverage matrices and error vectors will be then computed for such assertions. This matrix, used internally by the program, has the OCL assertions as rows and the joint models as columns. Cell $\langle i, j \rangle$ is assigned 1 if the i th OCL assertion is not satisfied when executing the model transformation with the j th input model, and 0 otherwise. Therefore, an OCL assertion has failed when there is at least a 1 in its row.

- (5) With the information obtained in the previous step, plus the information of the rules execution stored in the *Trace Models* (cf. Section 2.2.3), this step, namely *Suspiciousness-Based Rankings Computation*, produces the *Suspiciousness-Based Rankings* for all the non-satisfied OCL assertions. In our implementation, we have integrated 18 techniques, so 18 rankings for each non-satisfied assertion are computed. Any other technique can be trivially included in our tool.

To obtain these rankings, we first need to construct the coverage matrices and error vectors. This is done with the information of the (non-)satisfied OCL assertions in the execution of each input model. For the coverage information, we need the *Trace Models*. This means that the coverage matrices are constructed by reading all trace models. As we see for instance in Table 2, the coverage matrices store information of the rules exercised in the execution with each input model. For the creation of the error vectors, we need information of the non-satisfied OCL assertions.

With the information of the coverage matrices and error vectors, we are able to automatically compute the eight values described in Section 3.2 for computing the suspiciousness, namely N_{CF} , N_{UF} , N_{CS} , N_{US} , N_C , N_U , N_S , and N_F . Finally, with these values and the formulae for calculating the suspiciousness with the 18 techniques considered in this study (cf. Section 4.2.5 and Table 6), we obtain the *Suspiciousness-Based Rankings*. These rankings, together with the coverage matrices, error vectors, and values, are returned as comma-separated values (CSV) files by our tool. In particular, it returns a detailed CSV file for each non-satisfied OCL assertion.

- (6) Finally, our tool returns the EXAM scores—in the best-case, worst-case, and average-case scenarios (cf. Section 4.2.6)—for all the 18 techniques and for every non-satisfied OCL assertion. This information is inserted in the CSV files mentioned before. As explained in Section 3.2, this score basically measures the percentage of rules that need to be checked until the faulty rule is found. For this reason, we need as input information of which the buggy rules are represented in the figure as *Information About Buggy Rules*. The automatic computation of the EXAM score has been extremely useful for evaluating our approach, since no manual calculations have been needed. The results of the evaluation are described in the next section.

4 EVALUATION

4.1 Research Questions

The research questions (RQs) that we want to answer in this work are the following:

- **RQ1 - Feasibility.** *Is it possible to automate the process of locating faults in model transformations applying spectrum-based techniques?* Since, at the time of writing, there was no proposal for applying spectrum-based techniques for locating faults in model transformations, we want to answer whether this is feasible. This means that we want to check whether it is possible to automatically obtain for a model transformation a suspiciousness-based ranking, according to SBFL techniques, that indicates which rules should be inspected first in case of failure.

Table 3. Model Transformations Used as Case Studies and Their Characteristics

Transformation Name	# Classes MM Input - Output	# LoC	# Rules M-(U)L-C	# Helpers	Rule inheritance	Imperative part	Conditions	Filters	resolveTemp
UML2ER	4 - 8	53	8-0-0	0	✓	✗	✗	✗	✗
BibTeX2DocBook	21 - 8	393	9-0-0	0	✗	✗	✓	✓	✓
CPL2SPL	33 - 77	503	18-1-0	6	✗	✗	✓	✓	✗
Ecore2Maude	13 - 45	1055	7-7-25	41	✗	✓	✓	✓	✓

- **RQ2 - Effectiveness.** *How effective are state-of-the-art techniques for suspiciousness computation in the localization of faulty rules in model transformations?* Since many techniques have been proposed in the literature in different fields, we want to determine how they behave, comparing among each other, in the context of model transformations. This means we want to study which techniques provide the best suspiciousness-based rankings and which ones provide the worst rankings.
- **RQ3 - Accuracy.** *Is our approach able to accurately locate faulty rules in model transformations?* After studying the 18 techniques and comparing them, we want to conclude whether it is possible to state that applying spectrum-based techniques can accurately help the developer in the debugging of model transformations. This will be answered affirmatively if the techniques that are more effective, according to the answer to the previous RQ, provide accurate suspiciousness-based rankings.
- **RQ4 - Dynamic vs Static.** *How does our approach behave in comparison with a static approach?* As it is approach dynamic, we want to compare its performance with a notable approach for locating bugs in model transformations applying a static approach [18].

4.2 Experimental Setup

4.2.1 *Case Studies.* We have used four case studies to evaluate our approach and developed solution. Two of them have been taken from the open-source ATL Zoo repository [12] and the two others from research projects and tools. They all differ regarding the application domains, size of metamodels, and the number and types of ATL features used. Table 3 summarizes some information regarding the transformations. For instance, the size of the metamodels vary from 4 to 33 classes in the input metamodels and from 8 to 77 classes in the output metamodels. Regarding the size of the transformations, the number of rules ranges from 8 to 39 (in the table, *M* stands for *matched rules*, *(U)L* for *(unique) lazy rules*, and *C* for *called rules*) and the lines of code (LoC) from 53 to 1055. This means that the smaller transformation is approximately 20 times smaller, in terms of LoC, than the biggest one. Furthermore, the transformations contain from no helper to 41 of them. The table includes further information, namely whether rule inheritance, imperative rules, conditions, and filters are used within the transformations. We have slightly tweaked some transformations to increase their variability. For instance, in the *BibTeX2DocBook*, we have integrated the helpers within the rules, since the same transformation with the same behavior can be written with and without helpers [80], or in the *CPL2SPL* we have included some rules to transform features that were not included in the original transformation. All transformations are available on our website [100] and briefly described in the following:

- **UML2ER.** This transformation is taken from the structural modeling domain. It generates Entity Relationship (ER) diagrams from UML Class Diagrams. This transformation is originally taken from References [111], and we have considered the version proposed in Reference [18], which represents an extension. The aspect to highlight in this model transformation is the high use of rule inheritance. If $R_i < R_j$ means that R_i inherits from R_j ,

then we have $R8, R7 < R6; R6, R5 < R4; R4, R3, R2 < R1$. The presence of inheritance may worsen the results of SBFL techniques. Imagine we have, for instance, $R3 < R2 < R1$ in a model transformation, and rule $R3$ is executed. In the trace, it is stored not only the execution of $R3$ but also the execution of $R2$ and $R1$, since the code in the *out* part of these rules is actually executed. Therefore, if we have an error in one of the three rules, then the suspiciousness rankings will not make any difference between the three rules, with the three of them having the same suspiciousness value.

- **BibTeX2DocBook.** This case study is the one used as running example in our article. It is shown in Listing 1, and a complete description is available in Reference [54].
- **CPL2SPL.** This transformation, described in Reference [63], is a relatively complex example available in the ATL Zoo [12]. It handles several aspects of two telephony DSLs, SPL and CPL, and was created by the inventors of ATL.
- **Ecore2Maude.** This is a very large model transformation that is used by a tool called *e-Motions* [86]. It converts models conforming to the Ecore metamodel into models that conform to the Maude [26] metamodel to apply some formal reasoning on them afterwards.

4.2.2 Test Suite. Since this is a dynamic approach, we need input models to trigger the model transformations. For evaluating our work, we have developed a lightweight random model generator that, given any metamodel, produces a user-defined number of random model instances. The rationale behind our model generator is to produce a set of models with a certain variability degree. It creates an instance of the root class of the metamodel and, from such instance, it traverses the metamodel and randomly decides, for each containment relationship, how many instances to create for each contained class, if any. This process is repeated iteratively until the whole metamodel is traversed. After all instances and containment relationships are set, non-containment relationships are created, respecting the multiplicities indicated in the metamodel. Also, attributes are given random values. Alternatively, it is possible to generate models with some predefined structure, by indicating the minimum and maximum number of entities to create. The values to be given to specific attributes can also be preset by the user.

For our evaluation, we have created 100 models conforming to the input metamodel of each of the case studies with our model generator. We may mention that any model generator tool that produces models with a certain degree of variability could be used for generating the models—recall that such variability is necessary so that the input models exercise different parts of the transformation, producing a useful spectra. For instance, the *EMF (pseudo) random instantiator* could be used.⁶ Also, if there were enough models available produced manually, then these could be used, and we would not need to execute any models generator.

In total, we have created 117 OCL assertions for the four case studies, as displayed in the first part of column 3 in Table 4. These assertions are satisfied by the original version of the model transformations. Some of them correspond to the OCL assertions defined in the static approach by Burgueño et al. [18], since we want to compare our approach with this one (cf. Section 4.5)—see the second part of column 3. As indicated in the table, we use 100 input models for evaluating each case study. According to Section 3.1, the total number of test cases is measured as the Cartesian product of input models and OCL assertions: $|T| = |S| \times |O|$. As shown in the table, we have 1,400, 2,700, 3,400, and 4,200 test cases in each of the transformations, having a total of 11,700 test cases.

4.2.3 Mutants. To test the usability and effectiveness of our approach, we apply mutation analysis [58], so that we have produced mutants for all model transformations, where artificial bugs have been seeded. We have used the operators presented in Reference [98] and have applied them

⁶It is described on <http://modeling-languages.com/a-pseudo-random-instance-generator-for-emf-models/>.

Table 4. Case Studies and Artifacts for the Evaluation

Case study	# Input models	# OCL assertions (/ from [18])	# Test suite ($ T = S \times O $)	# Mutants	# OCL assertions violated
UML2ER	100	14/10	1,400	18	90
BibTeX2DocBook	100	27/16	2,700	40	269
CPL2SPL	100	34/15	3,400	50	150
Ecore2Maude	100	42/3	4,200	50	155
Total	400	117/44	11,700	158	664

in the different case studies. The aim of these operators is the same as the ones presented by Mottu et al. [77], i.e., they try to mimic common semantic faults that programmers introduce in model transformations. While Mottu et al. propose operators independent of any transformation language, we use a set of operators specific for ATL [62]. For instance, Mottu et al. [77] present several operators related to model *navigation*, such as *ROCC: relation to another class change*, which in Reference [98] it is materialized as *binding feature change*.

Recall that the aim of our approach is to semantically check the correctness of model transformations against a set of OCL assertions and to help localize the bugs. These OCL assertions are specified on input and output models. This means that, to be able to apply the approach, the model transformation needs to finish, i.e., it must generate output models. For this reason, the model transformation mutants that we have generated do not throw runtime errors for any of the test models created, i.e., they all finish their execution and generate output models. To be able to have such restricting mutants and as many other approaches do [9, 46, 77], we have generated them manually using the operators proposed in Reference [98].

In total, we have manually created 158 mutants, where each mutant is a variation of the original model transformation. For instance, Listing 3 displays the excerpt of a mutant where the operator *binding deletion* is used for deleting the only binding of rule 2. Our set of OCL assertions in the different case studies is complete enough to kill all 158 mutants, i.e., all mutants make at least one OCL assertion fail, indicating there is an error in the MT. Table 4 displays the number of mutants that have been created for each case study (column 5), while Table 5 presents the mutation operators [98] that have been used for creating the mutants, and the number of mutants in which the operators are applied. Please note that fewer mutants have been created for the *UML2ER* case study. This is due to the fact that this model transformation is the smallest one, and it is actually almost four times smaller than the second smaller one in terms of lines of code (cf. Table 3).

We may mention that more than one mutation operators can be used for constructing one mutant. For instance, we can combine *out-pattern element addition* with *binding addition* to create a mutant that adds one more target element and updates one of its features. We have also created mutants with more than one faulty rule. The reason for including higher-order mutants [57, 81] is the definition of realistic mutants, i.e., mutants that produce valid models and reproduce typical mistakes caused by developers. In fact, as presented in the survey on software fault localization by Wong et al. [115], having programs with a single bug (i.e., each faulty program has exactly one bug) is not the case for real-life software, which in general contains multiple bugs. Results of a study [49] based on an analysis of fault and failure data from two large, real-world projects show that individual failures are often triggered by multiple bugs spread throughout the system. Another study [69] also reports a similar finding. The very same reality occurs in model transformations, where it is not common to have isolated faults located in only one rule. Indeed, since some rules have implicit relations among them (cf. Section 2.2.3), it is very common to have errors spread in several rules.

Table 5. Mutation Operators Used and Number of Mutants Where They Are Applied

Mutation Operator (from Reference [98])	UML2ER	BT2DB	CPL2SPL	Ecore2Maude	Total
In-pattern element addition	1	2	5	3	11
In-pattern element class change	0	1	4	0	5
Filter addition	1	0	5	5	11
Filter deletion	0	3	1	0	4
Filter condition change	3	6	1	0	10
Out-pattern element addition	4	5	11	10	30
Out-pattern element deletion	0	3	4	8	15
Out-pattern element class change	2	3	6	0	11
Out-pattern element name change	0	1	0	3	4
Binding addition	2	3	8	0	13
Binding deletion	3	13	17	11	44
Binding value change	3	17	12	15	47
Binding feature change	1	1	5	6	13
Total mutation operators used	20	58	79	61	218

```

1 rule Author {                               -- tr2
2   from
3     a : BibTeX!Author (thisModule.authorSet->includes(a))
4   to
5     p1 : DocBook!Para () --binding deletion
6 }
```

Listing 3. Excerpt of a mutant of *BibTeX2DocBook* MT.

4.2.4 Set of Non-Satisfied OCL Assertions. As described, we have produced 158 mutants that correspond to buggy versions of the model transformations in the different case studies. Each one of them may violate one or more of the OCL assertions defined for the corresponding case study (of course, more than one mutant may violate the same assertion). In total, the 158 mutants make 664 OCL assertions fail, as displayed in the last column of Table 4, so the results of our evaluation are extracted from the 664 suspiciousness-based rankings obtained, one for each violated assertion. These rankings are the results of suspiciousness values calculated with 664 coverage matrices and with the corresponding 664 error vectors. These coverage matrices have different sizes depending on the case study. All of them have 100 columns, since we are using 100 input models, and the number of rows is determined by the number of rules in the model transformation.

4.2.5 Techniques for Suspiciousness Computation. We are interested in studying how different techniques⁷ for computing the suspiciousness of program components behave in the context of model transformations. To this end, we have surveyed articles that apply spectrum-based fault localization techniques in different contexts and have selected the 18 techniques that, together with their corresponding formulae, are displayed in Table 6. *Tarantula* [59] is one of the best-known fault localization techniques. It follows the intuition that statements that are executed primarily by more failed test cases are highly likely to be faulty. Additionally, statements that are executed primarily by more successful test cases are less likely to be faulty. The *Ochiai* similarity coefficient is known from the biology domain, and it has been proven to outperform several other coefficients used in fault localization and data clustering [3]. This can be attributed to the Ochiai coefficient

⁷Throughout the evaluation, we use the terms *techniques* and *metrics* indistinctly.

Table 6. Techniques Applied for Suspiciousness Computation

Technique	Formula
Arithmetic Mean [117]	$\frac{2(N_{CF} \times N_{US} - N_{UF} \times N_{CS})}{(N_{CF} + N_{CS}) \times (N_{US} + N_{UF}) + (N_{CF} + N_{UF}) \times (N_{CS} + N_{US})}$
Barinel [2]	$1 - \frac{N_{CS}}{N_{CS} + N_{CF}}$
Baroni-Urbani & Buser [114]	$\frac{\sqrt{N_{CF} \times N_{US} + N_{CF}}}{\sqrt{N_{CF} \times N_{US} + N_{CF} + N_{CS} + N_{UF}}}$
Braun-Banquet [115]	$\frac{N_{CF}}{\max(N_{CF} + N_{CS}, N_{CF} + N_{UF})}$
Cohen [78]	$\frac{2 \times (N_{CF} \times N_{US} - N_{UF} \times N_{CS})}{(N_{CF} + N_{CS}) \times (N_{US} + N_{CS}) + (N_{CF} + N_{UF}) \times (N_{UF} + N_{US})}$
D* [113]	$\frac{(N_{CF})^*}{N_{CS} + N_F + N_{CF}}$
Kulczynski2 [78]	$\frac{1}{2} \times \left(\frac{N_{CF}}{N_{CF} + N_{UF}} + \frac{N_{CF}}{N_{CF} + N_{CS}} \right)$
Mountford [114]	$\frac{N_{CF}}{0.5 \times ((N_{CF} \times N_{CS}) + (N_{CF} \times N_{UF})) + (N_{CS} \times N_{UF})}$
Ochiai [3]	$\frac{N_{CF}}{\sqrt{N_F \times (N_{CF} + N_{CS})}}$
Ochiai2 [11]	$\frac{N_{CF} \times N_{US}}{\sqrt{(N_{CF} + N_{CS}) \times (N_{US} + N_{UF}) \times (N_{CF} + N_{UF}) \times (N_{CS} + N_{US})}}$
Op2 [78]	$N_{CF} - \frac{N_{CS}}{N_S + 1}$
Phi [74]	$\frac{N_{CF} \times N_{US} - N_{UF} \times N_{CS}}{\sqrt{(N_{CF} + N_{CS}) \times (N_{CF} + N_{UF}) \times (N_{CS} + N_{US}) \times (N_{UF} + N_{US})}}$
Pierce [115]	$\frac{(N_{CF} \times N_{UF}) + (N_{UF} \times N_{CS})}{(N_{CF} \times N_{UF}) + (2 \times N_{UF} \times N_{US}) + (N_{CS} \times N_{US})}$
Rogers & Tanimoto [73]	$\frac{N_{CF} + N_{US}}{N_{CF} + N_{US} + 2(N_{UF} + N_{CS})}$
Russel-Rao [85]	$\frac{N_{CF}}{N_{CF} + N_{UF} + N_{CS} + N_{US}}$
Simple Matching [115]	$\frac{N_{CF} + N_{US}}{N_{CF} + N_{CS} + N_{US} + N_{UF}}$
Tarantula [59]	$\frac{\frac{N_{CF}}{N_F}}{\frac{N_{CF}}{N_F} + \frac{N_{CS}}{N_S}}$
Zoltar [55]	$\frac{N_{CF}}{N_{CF} + N_{UF} + N_{CS} + \frac{10000 \times N_{UF} \times N_{CS}}{N_{CF}}}$

being more sensitive to activity of potential fault locations in failed runs than to activity in passed runs. *Ochiai2* is an extension of such technique [11, 78, 115]. *Kulczynski2*, taken from the field of artificial intelligence, and *Cohen* have showed promising results in preliminary experiments [78, 117]. *Russel-Rao* has shown different results in previous experiments [85, 116, 117], while *Simple Matching* has been used in clustering [78]. *Reogers and Tanimoto* presented a high similarity with *Simple Matching* when ranking in the study performed in Reference [78]. The framework called *Barinel* [70] combines spectrum-based fault localization and model-based debugging to localize single and multiple bugs in programs. *Zoltar* [55] is also a tool set for fault localization. *Arithmetic Mean*, *Phi* (Geometric Mean), *Op2*, and *Pierce* have been considered in some comparative studies with other metrics [78, 115, 117]. *Mountford* behaves as the second best technique, among 17 of them, for a specific program in a study performed in Reference [114], where *Baroni-Urbani and*

Buser is also studied. As for D^* , its numerator, $(N_{CF})^*$, depends on the value of “ \ast ” selected. This technique resulted the best technique in the study performed in Reference [113], where “ \ast ” was assigned a value of 2. We have followed the same approach, so we have $(N_{CF})^2$ in the numerator of the formula.

Note that the computation of these formulae may result in having zero in a denominator. Different approaches mention how to deal with such cases [79, 118, 119]. Following the guidelines of these works, if a denominator is zero and the numerator is also zero, then our computation returns zero. However, if the numerator is not 0, then it returns 1 [119].

4.2.6 Evaluation Metric. To compare the effectiveness of the different SBFL techniques, we apply the EXAM score described in Section 2.3. In the context of this work, this score indicates the percentage of transformation rules that need to be examined until the faulty rule is reached. Its value is in the range $[1/(\text{num rules}), 1]$, and the higher its value, the worse.

Since there can be ties in the rankings obtained from the suspiciousness values, we compute the EXAM scores in the best-, worst-, and average-case scenarios. If the faulty rule is ranked in the same position as several other rules, then the best-case scenario assumes that the faulty rule is inspected first. In this sense, if the faulty rule is tied with many other rules, then the EXAM score is likely to be low. On the contrary, the worst-case scenario assumes that the faulty rule is inspected last. For this reason, if the faulty rule is tied with many other rules, then the EXAM score is likely to be high. In between we have the average-case scenario, which considers that the faulty rule is located in the middle place of a tie. Therefore, if it is in a tie with $(n - 1)$ other rules, then it will be inspected in the $(n/2)^{\text{th}}$ position.

4.2.7 Execution Environment. All the runs have been executed in a PC running the 64-bit OS Windows 10 Pro with processor Intel Core i7-4770 @ 3.40GHz and 16GB of RAM. We have used Eclipse Modeling Tools version Mars Release 2 (4.5.2), and we had to install the plugins ATL (we have used version 3.6.0) and ATL/EMFTVM (version 3.8.0). Finally, Java 8 is required.

4.3 Experimental Results

Table 7 shows the descriptive statistics of the EXAM score for each suspiciousness computation technique when applied to each case study on the three evaluation scenarios (average-, best-, and worst-case scenarios)—ignore for now the rows with the numbers for *Burgueño’15*, as those numbers are commented in Section 4.5.3. We may recall that the EXAM score, in the range $(0, 1]$, indicates the percentage of transformation rules that need to be inspected to locate the faulty rule. This score is never 0, since the inspection of the faulty rule counts. For this reason, since the MTs under test for each case study contain a different number of rules (9 in *Bibtex2DocBook*, 19 in *CPL2SPL*, 39 in *Ecore2Maude*, and 8 in *URML2ER*), the best possible values (the case where the faulty rule is ranked first in the suspiciousness rank) for the score are $\frac{1}{9} = 0.\overline{1}$, $\frac{1}{19} = 0.052631$, $\frac{1}{39} = 0.\overline{025641}$, and $\frac{1}{8} = 0.125$, respectively. Conversely, the worst value is always $1 = \frac{9}{9} = \frac{19}{19} = \frac{39}{39} = \frac{8}{8}$ (the faulty rule is ranked last). The table also shows, in the last two columns, the average mean and standard deviation values considering all case studies.

Having a look at the average EXAM scores in the average-case scenario, we observe there are eight techniques where less than 25% of the rules need to be inspected to locate the faulty rule, i.e., their EXAM score is below 0.25. These are, ordered from lower to higher percentage, *Mountford*, *Kulcynski2*, *Ochiai*, *Zoltar*, *Phi*, *Arithmetic Mean*, *Braun-Banquet*, and *Op2*. If we have a look at these eight techniques in the best-case scenario, then we see that *Phi* and *Arithmetic Mean* have the lowest, and therefore best, numbers. However, their numbers are the worst among these techniques in the worst-case scenario, implying that these techniques produced quite a large number of ties.

Table 7. Descriptive Statistics of the EXAM Score per Scenario and Case Study and Overall Values

	Technique	Bibtex2DocBook			CPL2SPL			Ecore2Maude			UML2ER			AVERAGE	
		mdn	mean	sd	mdn	mean	sd	mdn	mean	sd	mdn	mean	sd	mean	sd
AC	Arithmetic Mean	.111	.284	.240	.105	.166	.171	.077	.175	.181	.188	.313	.229	.234	.205
	Barinel	.333	.391	.216	.184	.245	.146	.192	.237	.172	.406	.363	.176	.309	.178
	Braun-Banquet	.222	.277	.182	.105	.192	.199	.090	.160	.179	.188	.337	.319	.242	.219
	B-U & Buser	.333	.365	.236	.105	.172	.189	.077	.134	.142	.188	.336	.319	.252	.221
	Cohen	.333	.343	.229	.105	.169	.170	.077	.176	.182	.188	.313	.229	.250	.202
	Dstar	.222	.326	.310	.263	.293	.213	.231	.423	.344	.469	.537	.298	.395	.292
	Kulcynski2	.111	.177	.142	.105	.185	.203	.077	.133	.139	.188	.331	.313	.207	.199
	Mountford	.111	.204	.151	.053	.157	.194	.077	.123	.111	.188	.337	.316	.205	.193
	Ochiai	.111	.188	.147	.105	.185	.193	.077	.134	.143	.188	.333	.318	.210	.200
	Ochiai2	.444	.443	.270	.105	.180	.191	.077	.175	.182	.188	.313	.229	.278	.218
	Op2	.111	.182	.149	.105	.226	.217	.154	.245	.228	.188	.331	.313	.246	.227
	Phi	.111	.268	.237	.105	.166	.172	.077	.172	.179	.188	.313	.229	.230	.204
	Pierce	.833	.682	.285	.737	.636	.283	.667	.596	.203	.688	.611	.301	.631	.268
	Rogers & Tani.	.556	.454	.277	.053	.206	.235	.077	.132	.140	.188	.302	.289	.273	.235
	Russel Rao	.222	.255	.121	.105	.240	.222	.333	.367	.182	.375	.438	.262	.325	.197
	Simple Matching	.556	.454	.277	.053	.206	.235	.077	.132	.140	.188	.302	.289	.273	.235
	Tarantula	.333	.398	.221	.092	.164	.191	.167	.211	.172	.438	.499	.259	.318	.211
	Zoltar	.111	.177	.142	.105	.182	.198	.154	.197	.185	.188	.331	.313	.222	.209
	Burgueño'15	.388	.436	.245	.105	.239	.224	.167	.317	.312	.375	.476	.297	.367	.269
BC	Arithmetic Mean	.111	.260	.233	.105	.161	.165	.026	.073	.112	.125	.196	.173	.173	.171
	Barinel	.333	.342	.235	.158	.229	.141	.051	.095	.112	.125	.168	.139	.208	.157
	Braun-Banquet	.222	.277	.182	.105	.180	.178	.026	.107	.175	.125	.308	.325	.218	.215
	B-U & Buser	.333	.365	.236	.105	.163	.171	.026	.081	.130	.125	.307	.326	.229	.216
	Cohen	.333	.320	.228	.105	.164	.163	.026	.075	.116	.125	.196	.173	.189	.170
	Dstar	.222	.325	.309	.263	.284	.202	.205	.372	.335	.438	.494	.310	.369	.289
	Kulcynski2	.111	.177	.142	.105	.176	.185	.026	.080	.132	.125	.301	.320	.184	.195
	Mountford	.111	.203	.151	.053	.148	.175	.026	.069	.097	.125	.304	.325	.181	.187
	Ochiai	.111	.188	.147	.105	.176	.176	.026	.081	.135	.125	.304	.325	.187	.196
	Ochiai2	.444	.416	.272	.105	.171	.174	.026	.072	.106	.125	.196	.173	.214	.181
	Op2	.111	.182	.149	.105	.221	.210	.026	.193	.241	.125	.301	.320	.225	.230
	Phi	.111	.245	.228	.105	.161	.165	.026	.070	.108	.125	.196	.173	.168	.169
	Pierce	.667	.587	.260	.658	.605	.262	.359	.410	.169	.375	.461	.308	.516	.250
	Rogers & Tani.	.556	.450	.277	.053	.195	.229	.026	.080	.131	.125	.274	.292	.250	.232
	Russel Rao	.111	.141	.122	.053	.196	.205	.026	.171	.247	.125	.261	.319	.192	.223
	Simple Matching	.556	.450	.277	.053	.195	.229	.026	.080	.131	.125	.274	.292	.250	.232
	Tarantula	.333	.349	.241	.053	.146	.176	.026	.068	.112	.125	.304	.330	.217	.215
	Zoltar	.111	.177	.142	.105	.173	.180	.026	.144	.192	.125	.301	.320	.199	.208
	Burgueño'15	.333	.342	.219	.0526	.106	.086	.154	.270	.279	.312	.458	.296	.253	.172

(Continued)

Table 7. Continued

	Technique	Bibtex2DocBook			CPL2SPL			Ecore2Maude			UML2ER			AVERAGE	
		mdn	mean	sd	mdn	mean	sd	mdn	mean	sd	mdn	mean	sd	mean	sd
WC	Arithmetic Mean	.111	.307	.283	.105	.171	.179	.128	.276	.317	.250	.429	.359	.296	.285
	Barinel	.444	.441	.261	.211	.260	.155	.256	.380	.303	.625	.557	.299	.409	.255
	Braun-Banquet	.222	.277	.182	.105	.205	.224	.154	.213	.190	.250	.365	.315	.265	.228
	B-U & Buser	.333	.365	.236	.105	.181	.211	.128	.187	.161	.250	.365	.315	.275	.231
	Cohen	.333	.367	.269	.105	.174	.177	.128	.278	.318	.250	.429	.359	.312	.281
	Dstar	.222	.326	.311	.263	.302	.229	.282	.474	.357	.500	.579	.290	.420	.297
	Kulcynski2	.111	.177	.142	.105	.194	.224	.128	.186	.155	.250	.360	.308	.229	.208
	Mountford	.111	.204	.151	.053	.167	.217	.128	.178	.134	.250	.369	.310	.230	.203
	Ochiai	.111	.188	.147	.105	.194	.215	.128	.188	.159	.250	.363	.314	.233	.209
	Ochiai2	.444	.470	.303	.105	.189	.213	.128	.279	.322	.250	.429	.359	.342	.299
	Op2	.111	.182	.149	.105	.231	.224	.231	.298	.219	.250	.360	.308	.268	.225
	Phi	.111	.292	.282	.105	.171	.179	.128	.273	.317	.250	.429	.359	.291	.284
	Pierce	1,000	.777	.335	.737	.667	.311	1,000	.781	.297	1,000	.761	.350	.746	.323
	Rogers & Tani.	.556	.458	.277	.053	.217	.244	.128	.184	.156	.250	.331	.290	.297	.242
	Russel Rao	.333	.369	.156	.105	.284	.249	.641	.563	.161	.625	.615	.260	.458	.206
	Simple Matching	.556	.458	.277	.053	.217	.244	.128	.184	.156	.250	.331	.290	.297	.242
	Tarantula	.444	.448	.264	.105	.182	.214	.231	.354	.303	.750	.693	.269	.419	.262
	Zoltar	.111	.177	.142	.105	.191	.220	.231	.250	.185	.250	.360	.308	.244	.214
	Burgueño'15	.444	.529	.317	.131	.372	.413	.179	.365	.371	0.5	.494	.303	.440	.351

Observing the eight techniques in the worst-case scenario, we see that *Mountford* and *Kulcynski2* are able to locate the faulty rule by inspecting less than 23% of the rules, so they seem to be the best techniques. In particular, *Kulcynski2* is able to locate the faulty rule first in the rank in 45% of the cases in the worst-case scenario, and it ranks the guilty rule in the top 3—i.e., only up to three rules need to be inspected to locate the fault—in 70% of the cases. The EXAM scores for these two techniques in the worst-case scenario are similar to the ones in the best- and average-case scenarios, concluding that there are not many ties. These two techniques are closely followed by *Ochiai* and *Zoltar*, techniques that do not produce many ties either and that are able to locate the faulty rule by inspecting less than 25% of the rules in the worst-case scenario. We can conclude that the four techniques with best results are, in this order, *Kulcynski2*, *Mountford*, *Ochiai*, and *Zoltar*. However, this ordering is not strict, since they behave slightly differently among them depending on the case study. In particular, to locate the fault, these techniques lead the debugger to inspect between 1.59 and 1.84 (of 9) rules in *BibTex2DocBook*, 2.98 and 3.5 (of 19) rules in *CPL2SPL*, between 4.78 and 7.68 (of 39) rules in *Ecore2Maude*, and between 2.65 and 2.69 (of 8) rules in *UML2ER* in the average-case scenario. The average standard deviation in all scenarios is around 0.2 for these four techniques, meaning that the results they provide are quite stable.

Going back to the average-case scenario and looking for techniques that give bad results, there are five techniques that need to inspect more than 30% of the rules to locate the faulty one, namely *Barinel*, *Russel Rao*, *Tarantula*, *Dstar*, and *Pierce*. Interestingly, the worst technique, so-called *Pierce* [115], needs to inspect more than 63% of the rules. This means that it performs even worse than random testing, and this is true in all case studies. Regarding the other four, *Dstar* needs to inspect almost 37% of the rules even in the best-case scenario, what is not a good result either. If we go to the worst-case scenario, then all these techniques need to inspect more than 40% of the

rules, so we can conclude that they do not behave good, and therefore we do not recommend to use them when applying SBFL in the MT domain.

The distributions of the results of each technique are graphically depicted in the box-plots of Figure 6, so they are useful to analyze each case study separately and see if the conclusions drawn so far are confirmed. The figure contains one box-plot per scenario (average-case labelled as AC, best-case labelled as BC, and worst-case labelled as WC) and case study, where the Y and X axis indicate the EXAM score and technique, respectively. These box-plots gather the results of the EXAM scores obtained with all mutants and depict them in vertical boxes—ignore for now the boxes for *Burgueño'15*, since they are commented in Section 4.5.3. The dots outside the boxes are known as outliers.

Having a look at the average-case scenarios in Figure 6, we can appreciate how techniques are categorized in two groups. On the one hand, techniques that perform well are represented by small boxes located at the bottom of each box-plot. We refer to these techniques as *good-performers*. However, the boxes of techniques with bad performance are stretched and typically located around the middle of the plot. We will name this group of techniques *bad-performers*. It is worth mentioning that among the group of good-performers, the most reliable ones are those with smaller vertical line segments above the box. This means that in the cases where faulty rules are difficult to locate, they provide lower EXAM scores than other good-performers.

For instance, in the *UML2ER* case study for the average-case scenario, the set of most-reliable good-performers comprises of *Kulcynski2*, *Mountford*, *Zoltar*, and *Ochiai*. In fact, the boxes of these four techniques are quite stable and similar in all scenarios of all case studies. Some other techniques, such as *Op2*, seem to provide similar performance, since for instance its boxes in the *UML2ER* case study are similar to those of these four techniques. However, the boxes are clearly worse in the *Ecore2Maude* and *CPL2SPL* case studies. At the same time, *Mountford* shows slightly better performance than the other three good-performers in some box-plots, such as in the *CPL2SPL* case study. Regarding the five techniques mentioned before that give bad results in the table of descriptive statistics, they fit in the profile of bad-performers. We can observe that their boxes are not uniform when comparing box-plots, having some boxes even located in the top of the charts.

Finally, it is worth noting that the box-plots in the *UML2ER* case study present the highest disparity among the three scenarios and that most techniques seem to behave worse in this case study than in the other three, showing larger boxes. This suggests that it is more challenging for the techniques to properly rank the faulty rule in this case study than in the other three case studies. This may be due to the high use of rules inheritance in this case study, what might complicate the location of the fault as explained in Section 4.2.1. Please also note, as commented in Section 4.2.3, that fewer mutants have been created for this case study compared with the other three. This could also have an impact in the results.

To study the data from a different perspective, namely the average values, we have constructed Figure 7. The figure presents a matrix where the suspiciousness computation techniques are represented by rows, and the mutants of the different case studies are represented by columns. Each cell is therefore colored according to its $EXAM_{m,t}^{Average}$, where m represents the mutant (X-axis) and t the technique (Y-axis).

As we can see in the color key, cells with lower values are lightly colored, while cells with higher values are darkly colored. The lighter the shade of cell $\langle i, j \rangle$, the better has performed technique j in mutant i on average. Observing the four techniques with good performance mentioned before, namely *Kulcynski2*, *Mountford*, *Ochiai*, and *Zoltar*, we see that their rows are lighter than the others in most cases. Similarly, a dark column in the matrix points out a MT with high EXAM values, meaning that it is hard to identify the faulty rule for such MT. This allows us to identify that

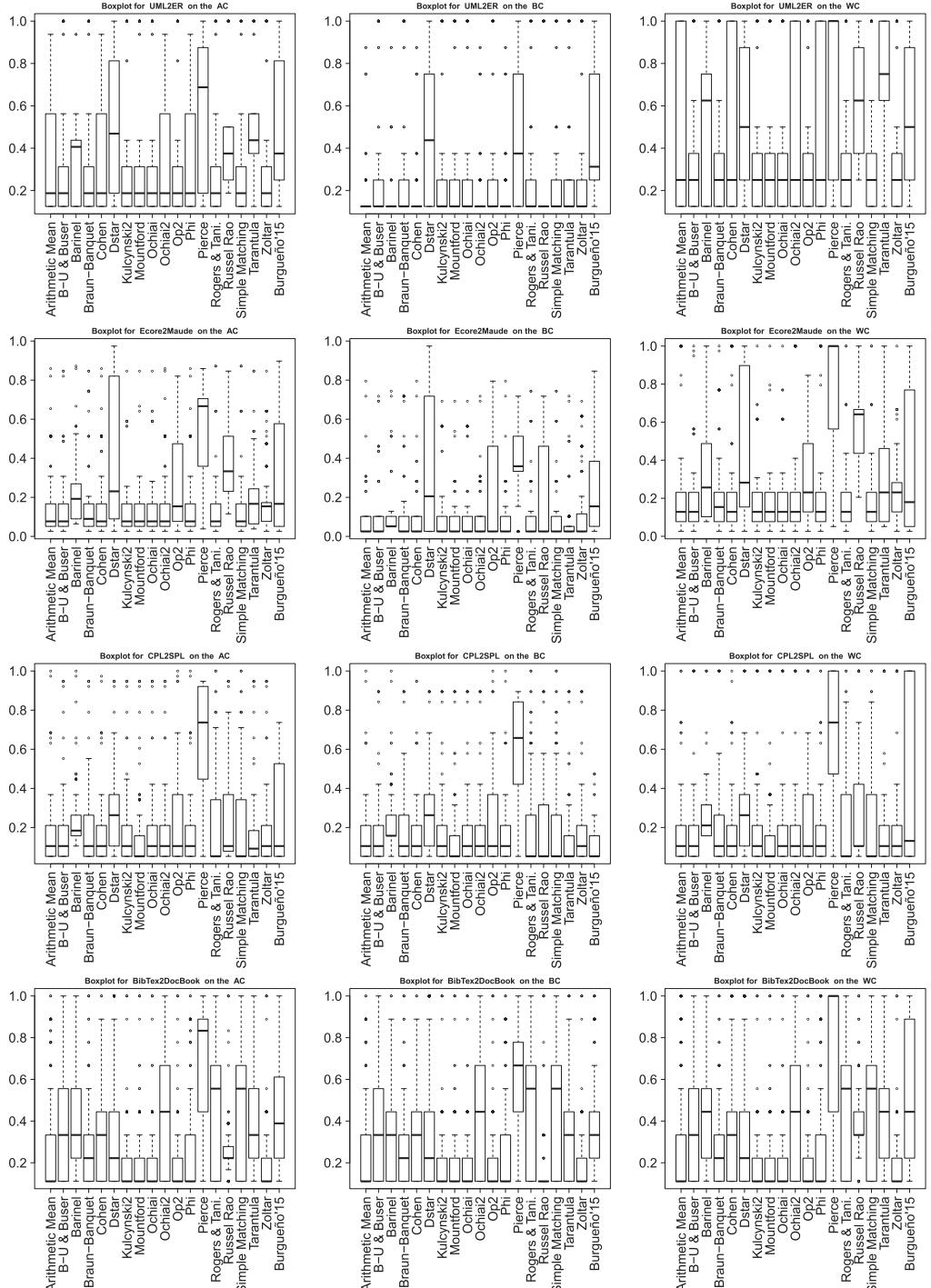


Fig. 6. Box-plot of the EXAM score of each technique per scenario and case study.

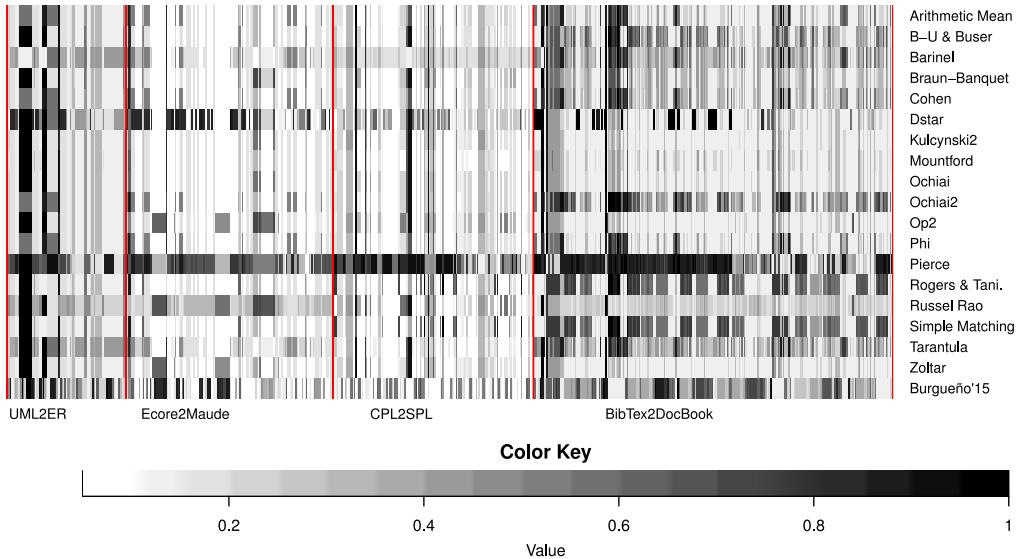


Fig. 7. Value per technique and case study.

the hardest case study in the study is *UML2ER*, with a significant amount of dark columns, what supports the conclusion drawn before. As mentioned earlier, we hypothesize that the reason for the techniques to behave worse in this case study than in the others is the high use of rules inheritance, since part of the behavior of the children rules is encoded in the parent rules, what may jeopardize the precision in the localization of the buggy rules.

Regarding performance in terms of run time, each run of our approach has taken between 4 and 75s (per mutant) on all the case studies. Please note that this is the time taken to execute the MT with all the source models, print in the console the violated OCL assertions, and compute and save in CSV files all the coverage matrices, error vectors, and suspiciousness rankings for all 18 techniques together with the automatically computed EXAM score for each violated assertion.

4.4 Statistical Results

The mutants and input models used in the evaluation were randomly generated, and thus a statistical analysis of the data was performed to study whether the differences observed among techniques are due to chance or not. Since the differences observed among the best-, average-, and worst-case scenarios are not disquieting, and to keep this article at a reasonable size, we focus on the analysis of results obtained in the average scenario, as it provides a better approximation to the accuracy of the technique in real settings.

4.4.1 Null Hypothesis Tests. The null hypothesis (H_0) states that there is not a statistically significant difference between the results obtained by different suspiciousness computation techniques, while the alternative hypothesis (H_1) states that at least for one pair of techniques such difference is statistically significant. Statistical tests provide the probability (named *p-value*) of getting the actual observed results based on the assumption that the null hypothesis is true. *p*-Values range in $[0, 1]$, for which researchers have established by convention that *p*-values under 0.05 represent so-called statistically significant values, and are sufficient to reject the null hypothesis. The results of the study do not follow a normal distribution, as confirmed by Shapiro-Wilk normality tests, thus the Friedman test was used for the analysis [40]. The resulting *p*-values were $< 1^{-10}$

for the results of the four case studies, leading us to reject H_0 for all of them. To find the specific techniques with statistically significant differences, pairwise comparisons were performed using Conover-Iman's Test [52]. More specifically, we compared all the possible pairs of techniques, out of 18 techniques under study, i.e., $\binom{18}{2} = \frac{18!}{2!(18-2)!} = 153$ pairwise comparison per case study. Additionally, we applied a correction of the p-values using the Holm post-hoc procedure [53], as recommended in Reference [31]. The results of the corrected p-values for the pairwise comparisons of all techniques are available on the project's website [100]. In summary, the percentage of pairwise comparisons revealing statistically significant differences was 96% in *Bibtex2DocBook*, 82% in *CPL2SPL*, 78% in *Ecore2Maude*, and 49% in *UML2ER*. Again, these data highlight that the latter case study is the one giving worse and more unstable results, which is consistent with the conclusion drawn in the analysis of the results in Section 4.3.

4.4.2 Effect-Size Estimation. To further investigate the differences between the different suspiciousness computation techniques, Vargha and Delaney's \widehat{A}_{12} statistic [105] was used to evaluate the effect size, i.e., determine which technique performs better and to what extent. Table 8 shows the effect size statistic for every pair of techniques. Each cell shows the \widehat{A}_{12} value obtained when comparing the suspiciousness computation technique in the column against the technique in the row. Vargha and Delaney [105] suggested thresholds for interpreting the effect size: 0.5 means no difference at all; values over 0.5 indicates a small (0.5–0.56), medium (0.57–0.64), large (0.65–0.71), or very large (0.72–1) difference in favour of the technique in the column; values below 0.5 indicate a small (0.5–0.44), medium (0.43–0.36), large (0.36–0.29), or very large (0.29–0.0) difference in favour of the technique in the row. Cells indicating medium, large, and very large differences in favor of the column are shaded in light grey, grey, and dark grey, respectively. The values in boldface are those where hypothesis test revealed statistical differences (p -value < 0.05). As expected, there is not a clear winner technique for all the case studies. However, the results confirm the superiority of *Mountford*, *Kulcynski2*, *Ochiai*, and *Zoltar*, showing from medium to large differences in 35, 30, 29, and 28 (of 72) pairwise comparisons. Analogously, the results support the bad performance of *Pierce*—outperformed by all of other techniques—*Barinel* and *Tarantula*.

4.5 Comparison Study

To answer RQ4 and to study whether our approach performs well in the location of faults in model transformations, we want to compare its effectiveness with a state-of-the-art approach based on the static analysis of transformation rules and assertions that obtained good results [18]. We believe the comparison of our approach with this one is fair and adequate for several reasons. First, the model transformation language used as proof of concept in both approaches is ATL. Second, OCL assertions are used in both approaches as oracle, i.e., to determine whether a model transformation is correct or not. Third, both approaches determine an order in which the rules must be examined to locate the faulty rules. Fourth, we are using in the evaluation of our approach the same four case studies proposed in Reference [18]. Fifth, we are able to use the mutants developed for evaluating our approach to evaluate the approach in Reference [18], and we are also able to compute the EXAM values (cf. Section 4.5.2) for such approach, so we can fairly compare both approaches. Finally, in the set of OCL assertions that we have created for each case study in this work, we have included all the OCL assertions the authors in Reference [18] proposed for evaluating their approach⁸ (cf. third column of Table 4). We have used those assertions, as well as some more that we have defined, for evaluating our approach. Since the tools developed for the static approach [18] are available (cf. Section 4.5.1), we have been able to run the approach with them. We

⁸The OCL assertions used in Reference [18] are available at http://atenea.lcc.uma.es/index.php/Main_Page/Resources/MTB.

have used the complete set of OCL assertions for comparing both approaches. To demonstrate that the results are not biased due to the new defined OCL assertions, we have also made a comparison considering only the OCL assertions defined in Reference [18]. This comparison is presented in the appendix of this article, where the results presented and conclusions drawn are very similar.

In the following, we first summarize how the approach by Burgueño et al. works and computes the rankings. Then, we explain how we are able to obtain EXAM values for such approach. Finally, we present and discuss the results of the comparison.

4.5.1 Static Fault Localization in Model Transformations. The article by Burgueño et al. [18] proposes a static approach for the localization of faults in model transformations. As in our approach, ATL is the language used as proof of concept and the assertions that the transformations must satisfy are also defined in OCL. Therefore, it follows the same methodology as proposed in this article (cf. Section 3.3). Also, like our approach, theirs is backed up by a tool. However, for determining if any OCL assertion fails (step 2 in the methodology), their approach relies on an external tool, namely *TractsTool* [19, 109]. This means that the user also needs to get familiarized with this other tool.

When executed, this static approach computes, for all OCL assertions, the order in which rules should be inspected to locate the bug. To do so, it computes, for each pair <assertion, rule>, the probability that the assertion failure comes from the rule making use of the common denominator that both have, namely the structural elements belonging to the metamodels. The approach builds on the following steps:

- (1) *Footprint Extraction.* The *structural elements*, referred to as *footprints*, of both model transformation and assertions are extracted.
- (2) *Footprint Matching.* The footprints extracted are compared for each rule and assertion.
- (3) *Matching Tables Calculation.* The percentage of types overlapping, so-called *alignment*, for each transformation rule and assertion is calculated. This information is used to produce the matching tables.
- (4) *Matching Tables Interpretation.* The resulting tables are analyzed for identifying the order in which rules should be inspected in case any OCL assertion fails.

To apply this approach, three tools need to be executed, two of which are proposed and implemented in Reference [18]. First, as mentioned, the *TractsTool* is executed to check which OCL assertions fail. Then, the *ATL Transformation Types Extractor* is executed to generate a model with the footprints of the ATL transformation. Finally, the *Matching Tables Calculator* uses, among others, such model as input and generates the matching tables, indicating also the order in which rules should be inspected in case of failure.⁹

Three matching tables are generated by this approach. They are matrices that have the OCL assertions as rows and transformation rules as columns. Two of them need to be inspected to determine the order in which rules should be inspected in case of failure—for a detailed explanation, the interested reader is referred to Reference [18].

4.5.2 EXAM Values Computation. To obtain the suspiciousness-based rankings for the approach in Reference[18], we have obtained the matching tables of all 158 mutants, for which we have made use of the available tools mentioned before, namely *ATL Transformation Types Extractor* and *Matching Tables Calculator*. We have developed a program that, for each case study, takes the matching tables as input together with a CSV file that contains information of the buggy rules

⁹The *ATL Transformation Types Extractor* and *Matching Tables Calculator* tools are available at http://atenea.lcc.uma.es/index.php/Main_Page/Resources/MTB/MTB.

in each mutant and the OCL assertions that fail in such mutant. With those inputs, this program computes the order in which the rules should be inspected for each of the OCL assertions that fail, which is the same concept as the suspiciousness-based ranking proposed in spectrum-based fault localization. Therefore, with this rules ordering, and since we have as input information of the buggy rule of each mutant, we are able to easily compute the EXAM score. As output, our program generates a CSV file indicating, for each mutant and each OCL assertion that fails, the EXAM score in the best-, worst-, and average-case scenarios (cf. Section 4.2.6).

All the artifacts used for the comparison, namely the 158 mutants and 117 OCL assertions, together with all the matching tables generated for all case studies are available on our project’s website [100].

4.5.3 Static Approach vs. Dynamic Approach: Results. As described, Table 7 provides the descriptive statistics of the EXAM score, where *Mountford*, *Kulcynski2*, *Ochiai*, and *Zoltar* show the best numbers. Regarding the static technique proposed by Burgueño et al. [18] (*Burgueño’15* in the table), it performs worse than these techniques. In the average-case scenario, the static approach needs to inspect around 37% of the rules to locate the fault, which is much more than the 20% that needs to be inspected by the best techniques. In particular, for each case study in the average-case scenario, the static technique needs to inspect 2.3 (of 9) more rules in *BibTex2DocBook* (25.9% of the MT), 1.387 (of 19) more rules in *CPL2SPL* (7.3% of the MT), 7.18 (of 39) more rules in *Ecore2Maude* (18.4% of the MT), and 1.3 (of 8) more rules in *UML2ER* (16.2% of the MT) compared with the best techniques in each case. Regarding the number of ties, there is no uniform behavior. For instance, in *BibTex2DocBook* and *CPL2SPL* there are clearly more ties in the static technique compared to the best dynamic techniques, since the difference in the EXAM score in the best- and worst-case scenarios is bigger. As for *Ecore2Maude* and *UML2ER*, the number of ties seems to be similar among both approaches.

Looking at the worst dynamic techniques, the static approach seems to behave better than some of them. Having a look at the average mean (penultimate column), it behaves much better than *Pierce* in the average-case scenario, since the latter technique needs to inspect more than 63% of the rules to locate the fault. It also performs better than *Dstar* in this scenario, since this technique needs to inspect more than 39% of the rules. Finally, the static technique by Burgueño et al. performs worse than *Russel Rao* in the average-case scenario but a bit better in the worst-case scenario. Therefore, for now we can conclude that the static technique may behave better than 2 dynamic techniques and clearly behaves worse than other 15 techniques, but let us delve deeper into the results.

We further analyze the results by looking at each case study in the box-plots of Figure 6. In general, we notice that the results of the static approach are typically similar among the three scenarios, although the boxes are larger than those of most dynamic techniques, indicating a worse performance. We can appreciate that the static approach behaves normally better than *Pierce*, confirming our previous finding. As for *Dstar*, its boxes are in many plots larger than the ones of the static approach. However, in other plots its boxes are smaller, so we cannot confirm the superiority of the static technique with regards to *Dstar*. For instance, in the *BibTex2DocBook* and *CPL2SPL* case studies, the shape of the box-plots for *Dstar* seems to be normally better. Finally, regarding *Russel Rao*, its boxes have in most cases better shapes than those of the static approach.

The effect-size estimations of the statistical analysis for the static approach are displayed in Table 8. To begin with, we can see in the *BibTex2DocBook* case study that the four best SBFL techniques are clearly better than the static approach by Burgueño et al. [18], since the values in the row of the static approach are above 0.78 for the corresponding cells, indicating a very large difference in favor of the technique in the column. Also, the technique that seemed to be similar

to the static approach, namely *Dstar*, is proved to be better in this case study. In general, the color of the row shows that most techniques behave better than the static one.

In fact, looking at the four case studies, the numbers in the cells of the rows of the static approach and the columns with the best SBFL techniques—*Kulcynski2*, *Mountford*, *Zoltar*, and *Ochiai*—are always above 0.55, leaving no doubt that the static approach behaves worse. Besides, all these cells reveal statistical differences (p -value <0.05 , displayed in bold in the table).

The superiority of the static approach regarding *Pierce* is confirmed in all case studies. However, it cannot be concluded that it is better than any other of the techniques, since the rows of the static technique do not present a value <0.5 in more than two case studies for any of the other techniques. Finally, we see that in the *UML2ER* case study the static approach behaves generally much worse than most techniques. An explanation can be that the static approach, based on types matching, does not behave well in the presence of rule inheritance.

In summary, we can confirm that all SBFL techniques have a better performance when locating the faulty rule than the static technique, except for *Pierce*, where the static technique behaves clearly better. Besides, the static approach normally presents more ties than the best dynamic techniques.

Regarding performance in terms of runtime, static approaches are typically faster since they do not need to execute the program under test. This is the same in our case, where the static approach is faster. In any case, it requires to perform footprints extractions—both in OCL assertions and ATL transformation rules—and footprints matching, that also requires some resources. Altogether, the static approach takes from less than 1s (in *UML2ER*) to 42s (in *Ecore2Maude*) per mutant, less than required by our dynamic approach (from 4 to 75se, cf. Section 4.3).

4.6 Discussion

The results of the exhaustive experiments described in the previous sections allow us to answer the research questions formulated in Section 4.1.

4.6.1 RQ1 - Feasibility. The first research question, related to the feasibility of the approach, “*RQ1: Is it possible to automate the process of locating faults in model transformations applying spectrum-based techniques?*,” can be answered affirmatively. Indeed, we have automated the process of locating the faulty rules in model transformations by means of a Java program¹⁰ that orchestrates ATL model transformations and uses the information stored in the traces to compute the suspiciousness-based rankings based on the program spectra. This automation is explained in Section 3.4. All the artifacts used as input and generated as output are available on our project’s website [100].

Furthermore, even though our program has been implemented for ATL model transformations, we are confident that it can be adapted for any transformation language that is able to store in a trace model the result of the execution. In fact, the trace model is nothing but an output model. Therefore, any model transformation language that is able to produce more than one output model can generate a trace model as output.

4.6.2 RQ2 - Effectiveness. The second research question has to do with the comparison of the techniques evaluated with our automated approach: “*RQ2: How effective are state-of-the-art techniques for suspiciousness computation in the localization of faulty rules in model transformations?*” This question has to do with how well the different techniques are able to position the faulty rule in the suspiciousness-based ranking. According to the results presented in Sections 4.3 and 4.4, we can conclude that the top four most-effective techniques are *Kulcynski2*, *Mountford*, *Ochiai*, and

¹⁰ Available on Github: https://github.com/javitroya/SBFL_MT.

Zoltar, the first two presenting slightly better overall results. At the other end, we have *Pierce* as the least-effective technique. The top three of non-effective techniques is completed by *Barinel* and *Tarantula*.

4.6.3 RQ3 - Accuracy. The third research question is related to the accuracy of the approach: “*RQ3: Is our approach able to accurately locate faulty rules in model transformations?*” The answer to this question is related to the previous one, since depending on the effectiveness of the techniques we will conclude whether the approach is accurate or not. In particular, we need to look at the most effective techniques. Evaluation results revealed that the best techniques place the faulty transformation rule among the three most suspicious rules in around 74% of the cases. Looking into each of the four case studies, the best techniques allow the tester to locate the fault by inspecting, on average, only 1.59 rules (of 9) in *BibTex2DocBook*, 2.99 rules (of 19) in *CPL2SPL*, 4.8 rules (of 39) in *Ecore2Maude*, and 2.4 rules (of 8) in *UML2ER*. According to these numbers, we can conclude that the application of spectrum-based fault localization is accurate in the context of model transformations if techniques such as *Mountford*, *Kulcynski2*, *Zoltar*, and *Ochiai* are applied, so we shall recommend to apply this approach to debug model transformations. These conclusions are supported by the evaluation of more case studies available on our project’s website [100].

4.6.4 RQ4 - Dynamic vs. Static. Our last research question has to do with the comparison of our dynamic approach with a notable static approach [18]: “*RQ4: How does our approach behave in comparison with a static approach?*” In summary, we can conclude that most dynamic techniques based on spectrum computation are better than the static approach for the localization of faults in model transformations. This was expected, since dynamic techniques execute the model transformation—from which they extract a lot of information—and the static approach does not. However, the static approach is still clearly better than one dynamic technique, namely *Pierce*. Furthermore, it also behaves better than other techniques, such as *Dstar*, *Tarantula*, *Simple Matching*, *Rogers and Tanimoto*, and *Barinel* in some case studies. It is also noteworthy that both approaches are complementary, and so it should be possible to define heuristics for the selection of the best technique on each application scenario, or even combine them. For example, it is better to apply the static approach in environments with low resources or when the transformations are very expensive to execute [65], for instance, in the case of transforming very large models [15, 25], and when it is not possible to get model instances of the source metamodel at the time of developing the model transformation.

4.7 Threats to Validity

According to Wohlin et al. [112], there are four basic types of validity threats that can affect the validity of our study. We cover each of these in the following paragraphs.

4.7.1 Conclusion Validity Threats. Threats to the conclusion validity are concerned with the issues that affect the ability to draw correct conclusions from the data obtained from the experiments. To mitigate these threats, we have applied statical analysis to confirm the conclusions drawn from the means and figures, and we have used the specific statistical tests and effect size measures recommended by the guidelines on empirical methodology. Furthermore, all the assumptions required for the application of the tests were checked, and the raw data and scripts for replication are available in the companion materials of this article [100].

4.7.2 Construct Validity Threat. It is concerned with the relationship between theory and what is observed. A possible construct validity threat (known as the mono-method bias) is related to the use of one single metric, the so-called EXAM score, to evaluate the performance of the approach and the suspiciousness computation techniques compared. Other metrics have been

proposed [115], such as the T-score [68], P-score [122], and N-score [43]. However, EXAM score is an accepted metric for measuring the quality of spectrum-based fault localization techniques, and has been used in a variety of works [115]. Moreover, we have decided to obtain the EXAM score as it is directly applicable in the context of model transformations. By considering the transformation rules as units of examination, the EXAM score is easy to understand, since it is directly proportional to the amount of rules to be examined, rather than to an indirect measurement in terms of the amount of code that does not need to be examined, as proposed by other scores.

Another possible construct validity threat is the mono-operation bias, which is related to the use of a single treatment or technique that could bias our conclusions. Since we compare the approach with a static alternative [18] and have used up to 18 suspiciousness-computation techniques and 4 use cases in our experiments, we consider that this threat is neutralized.

4.7.3 Internal Validity Threats. These threats are related to those factors that might affect the results of our evaluation. First, we may remark that this is a debugging approach, not a testing approach. Therefore, the objective of this work is not to generate high-quality test models, something addressed in related articles [6, 39, 44, 48, 95], but to localize the faults that triggered test failures. In fact, a key point in favor of our approach is that it can be used in conjunction with any method for test model generation, either random or guided. For evaluating our work, we have developed a lightweight random model generator that, given any metamodel, produces a user-defined number of random model instances, as explained in Section 4.2.2. With this generator we have obtained a set composed of 100 source models in the test suite of each case study, so a total of 400 models have been generated. These models have achieved full coverage—all rules and lines of code have been exercised—in all case studies. However, using more complex input model generation approaches [6, 48] may be required in those cases where random generation is not enough to achieve a sufficient coverage.

A second threat is that we have used in total 117 OCL assertions in the first study and 44 in the dynamic-vs.-static comparison study (cf. Table 4). We have tried to minimize this threat by constructing a set of OCL assertions that cover much of the specifications of the transformations. Besides, for the comparison study to be fair, we have taken the OCL assertions proposed in Reference [18]. Third, we have tried to create a large set of mutants, composed of 158 of them, and we have aimed at maximizing the variation of semantic faults and mutation operators used. Having used more or fewer mutants could have had an impact in the results. For instance, we recall that fewer mutants have been created for the UML2ER case study than for any of the other case studies, as commented in Section 4.2.3. Having used different mutation operators could have also had an impact in the results. For instance, some approaches propose mutation operators that yield run-time errors, such as the work by Sánchez-Cuadrado et al. [91], which proposes a powerful approach that relies on static analysis and type inference to locate, among others, run-time errors. However, please bear in mind that the purpose of the approach presented in this article is to localize semantic faults, i.e., it needs the model transformation to finish and produce output models, so that their satisfaction can be checked against the set of OCL assertions available. That is why we have used a subset [98] of the operators defined in Reference [91] and that mimic semantic faults likely to be made by programmers [77], as explained in Section 4.2.3. In any case, our approach is complementary to those aiming at spotting bugs that produce runtime errors [91].

As a final threat to internal validity, we may mention a weakness of SBFL, and generally of all fault localization techniques [115], which is the incapability of locating bugs resulting from missing code [121]. Same thing happens with our approach, it is likely to produce bad results if there are missing rules. For this reason, and as commented above, our approach can be complemented with Sánchez-Cuadrado et al.'s [91] approach, which identifies rules absence with a static analysis.

Indeed, the target elements created in a transformation rule typically reference or are referenced by target elements created in other transformation rules, so static analysis is a good technique for identifying the absence of rules that should create referencing or referenced target elements. For instance, in the model transformation shown in Listing 1, the target elements created by rule *Main* reference the target elements created by all the other rules. Likewise, the target elements created by all the other rules are referenced by those created by rule *Main*. Therefore, the absence of any of these rules can be detected with a static analysis tailored at examining that there will be no dangling references among the target elements created. Finally, this threat can also be mitigated with the definition of proper OCL assertions. For instance, in the transformation of Listing 1, the specification should dictate that there must be an element of class *DocBook* created for each element of class *BibTexFile*, so that the number of instances of both classes must be the same after executing the model transformation. This can be expressed with assertion *OCL4* in Listing 2. Therefore, even if we do not count on approaches like the one by Sánchez-Cuadrado et al. [91], the non-satisfaction of assertions such as *OCL4* can help the developer realize a rule is missing.

4.7.4 External Validity. These threats have to do with the extent to which it is possible to generalize the findings of the experiments. The first threat is that the results of our experiments have been obtained with four case studies, which externally threatens the generalizability of our results. To mitigate this threat, we have tried to select a set of model transformations that considers all ATL constructs and where the model transformations differ in their domains, size of metamodels and transformation, and variability of features used within the transformations, as reflected in Table 3. Furthermore, we have selected the same case studies as those used in the related article compared to our approach in Section 4.5, published in 2015 in the *IEEE Transactions on Software Engineering* journal [18]. Second, we have analyzed a set of 18 techniques for the computation of the suspiciousness-based rankings. Although it is a large set, result of doing a thorough literature review, we might have left aside some techniques that could give better results than the ones obtained with the best techniques of our study. Also, we have implemented our approach for ATL due to its importance both in industry and academia, so it would be interesting to test it with other transformation languages. However, we do believe our approach would produce similar results for any model transformation language based in rules as long as the result of their executions can be stored in traces (cf. Section 2.2.3), that allows to construct the coverage matrix and error vector and, therefore, apply SBFL techniques.

There are two other threats related to the external validity of the results that have to do with the program spectra creation in our implementation. In particular, we have used in our prototype the ATLAS transformation language and have considered the rules, of any type, as unit of examination and therefore as the components to be considered for constructing the spectra. Should we also have considered helpers in the spectra, the results of techniques effectiveness could have been different. This decision has been made considering related works that also check (ATL) transformations correctness against OCL assertions. While some approaches only check whether an assertion is violated or not by a model transformation [7, 19, 42, 80, 109], others propose to locate the fault when an assertion is not satisfied [7, 22, 23], but none of them inspect the helpers—they remain at the rule level. Crucial for our decision has been the static approach for locating faults proposed by Burgueño et al. [18], which does not consider helpers either and only locate faults in ATL rules. Should we have considered them, the thorough comparison with this approach presented in Section 4.5 would have been unfair. After having proved the effectiveness of SBFL techniques in the model transformation domain according to the extensive evaluation presented in this article, a natural evolution of this work is to perform a thorough study considering helpers to check if these techniques remain effective. In any case, if our current approach determines that a rule is faulty,

and it is calling a helper, then the user of the approach would inspect the rule and, if (s)he sees no fault, (s)he would proceed by inspecting the helper, so this threat is reduced.

Finally, the other threat is that the components considered in our approach might be too coarse-grained: our approach works at rule level. This means that the user would need, for example, to put more effort in locating a bug in a big rule than when doing it in a small rule. However, the complexity of transformation rules and model transformations is inherent to the bridges they try to build among different semantic domains, and different types of model transformations can be written depending on the problem to be solved [28, 64]. For instance, the creator of ATL recommends to use declarative code as much as possible.¹¹ Besides, some approaches exist for modularizing model transformations, so that they become as easy-to-understand and reusable as possible [34, 89]. Like with the threat before, another reason that led us to work at rule level in this approach is that related works that aim to locate bugs in model transformations against OCL satisfaction also propose approaches at rule level [22, 23], and especially the work with which we do an extensive comparison [18].

5 RELATED WORK

Due to the lack of oracles and formal semantics in model transformation languages, some approaches propose to translate the transformation specifications to other domains where formal treatment is possible. For instance, Troya and Vallecillo propose to translate ATL to the rewriting logic framework Maude [102], where some formal analysis can be performed, although the translation is not fully automated. Anastasakis et al. propose to translate QVT model transformations to Alloy to verify if some properties hold for the transformation, and there are also approaches for verifying contracts for ATL transformations based on the Coq proof assistant [20, 83]. Oakes et al. propose to translate the declarative part of ATL to the visual graph-based model transformation engine DSLTrans [80]. Visual contracts similar to our OCL assertions but less expressive can be then tested for satisfaction in DSLTrans. Similar visual contracts, using a visual language with formal semantics called PaMoMo, are used by Guerra et al. [47], but in this case their approach compiles such contracts into QVT and their satisfiability is checked with the PACO-Checker tool. A big difference of our approach with these is that we do not need to leave the model transformation development environment to check for the correctness of the MTs, so our approach stays within the Eclipse Modeling Framework dealing with Ecore metamodels and XMI models and the user does not need to be familiar with any other domain-specific language such as Maude, DSLTrans, Alloy, or Coq. Furthermore, our approach helps locate the faulty rules, that is not addressed in these approaches.

As in our approach, Cheng et al. [22] propose to verify if ATL transformations satisfy OCL assertions. However, to prove the correctness of the ATL transformation, they encode both the OCL assertions and the ATL transformation specification into the Boogie language [88]. Boogie is a procedure-oriented language that is based on Hoare-logic. Then, their developed VeriATL verification system indicates whether the ATL specification satisfies the specified OCL assertions or not. However, this approach does not report useful feedback to help the transformation developers fix the fault, which is the main objective of our approach. Cheng and Tisi [23] then build on this approach and tool (VeriATL) with the goal of localizing the fault by applying natural deduction and program slicing. However, instead of offering the developer with a rules ranking according to their chance to contain a bug, their approach determines scenarios, which are slices of the model transformation under test, where a certain OCL assertion is not satisfied together with the proof

¹¹<http://www.idi.ntnu.no/emner/tdt4250/Slides/M2M-atl-intria1117.pdf>.

tree. This is achieved by deriving sub-goals from the OCL assertions. Since this approach aims at locating a fault from a different perspective than ours, they can complement each other.

Burgueño et al. [18] propose an approach with a similar purpose as ours, but their approach is static. They also count on ATL model transformations and OCL assertions that must be checked for correctness and try to locate the faulty rule without translating the OCL assertions nor ATL transformations to any formal language. Their approach proposes to locate the faulty rules based on matching functions that automatically establish alignments among the metamodels footprints appearing in the transformation rules and those present in the OCL assertions. A comparison of this approach and the one we present in this article has been done in Section 4.5, where we have seen that most techniques for the spectrum-based localization of model transformations give better results than this static approach. Furthermore, this approach does not check if an OCL assertion is satisfied, but it resorts to the Tracts tool [19]. Contrarily, our approach does not need any input from external tools. A good aspect of the static approach is that it does not need any input model, since actually the transformation is not executed, and it requires shorter runtimes. This aspect makes this approach very useful in several situations. For example, it is better to apply the static approach in environments with low resources or when the transformations are very expensive to execute [65], for instance in the case of transforming very large models [15, 25], and when it is not possible to get model instances of the source metamodel at the time of developing the model transformation. Both approaches are therefore tangential.

There are other approaches that propose static analysis for debugging model transformations. Sánchez-Cuadrado et al. [91] combine static analysis and constraint solving to discover errors in ATL model transformations such as navigation errors (like invalid collection operations and operators), disconformities between the types used in the transformation and those declared in its source/target metamodels, integrity constraints regarding the semantics of ATL, problems related to dependencies between transformation rules and, in summary, any error that the current syntactic checker of ATL is not able to identify. They even provide possible suitable quick fixes based on speculative analysis [90]. These approaches have meant an important milestone in the evolution of ATL. Our approach is orthogonal to these and, consequently, can serve to complement them. Finally, Sánchez-Cuadrado et al. [92] have built, on top of their so-called *anATLyzer* tool described in the previous cited articles, an approach for checking contracts specified in the target language. Their approach translates these target contracts into source contracts by using the model transformation, so that they can predict, without the need to execute the model transformation, whether any specific input model will yield (in)correct target models. Since they use the model transformation for generating source contracts, it has to be correct. Therefore, different from our approach, this is not targeted to debugging model transformations, but to statically check target constraints in a lightweight manner.

The approach we present in this work is perfectly in line with the approach we presented in Reference [101] with the aim of locating bugs in three application scenarios of model transformations, namely regression testing, incremental transformations and migrations among transformation languages. Thus, the approach in Reference [101] proposes to automatically derive OCL assertions from a given ATL model transformation, which are satisfied by the transformation. The approach applies a technique known as metamorphic testing [93]. By applying metamorphic testing, and after identifying a set of patterns that normally takes place in the trace information stored after the execution of a model transformation, it is able to automatically derive so-called likely metamorphic relations, which can be seen as precisely the OCL assertions used in the current work. In this way, (i) in regression testing, (ii) when an original transformation is migrated to a different transformation language, or (iii) an incremental transformation is developed with the same behavior of the original transformation, the approach presented in this article can be used to check whether

the OCL assertions obtained for the original transformation by the approach in Reference [101] are satisfied in the latter evolved or modified transformations. Metamorphic testing has also been applied by He et al. [51], in this case for bidirectional model transformation testing.

We recall that this article focuses on debugging and not testing. Thus, we do not impose any constraint on how the source models are generated, either manually or automatically. In any case, some proposals for generating models have been proposed in the literature, where some of them require input by the tester. For instance, the model generator in Reference [17] requires the tester to provide metamodel fragments as input, or the one in Reference [95] requires input from the MMCC external tool [38] to provide model fragments. Other approaches propose the generation of models in different formats such as the Human Usable Textual Notation [41], so they need to be transformed prior to their use as input for model transformation languages integrated in the Eclipse Modeling Framework such as ATL. Some other more sophisticated model generators try to derive a set of input models from model transformations [44], what is not desired in our case, because we may be debugging erroneous transformations, and from OCL constraints [6, 48]. Most of these approaches can be used for generating test models for our approach. However, as explained in Section 4.2.2, we have used a lightweight model generator that, given a metamodel, it returns a set of random models conforming to such metamodel, where the models present certain variability among them with respect to the classes of the metamodel. Since none of the case studies contain complex graph constraints as preconditions that are difficult to cover with random graph generation, the models we have generated have obtained full coverage in all case studies, i.e., they have exercised all rules. However, in other cases, obtaining models with full coverage may require the use of more complex and computationally expensive methods and tools [5, 6, 44, 48, 95].

6 CONCLUSION

In this article, we have presented the first approach for debugging model transformations following an SBFL approach. We have implemented and automated it for the ATLAS transformation language due to its importance in both industry and academia. However, we are confident that the approach can be extensible to any model transformation language as long as it can store the output of its execution in a trace model. The implemented automation has allowed us to perform a thorough evaluation.

Taking as input the model transformation under test and a set of source models and OCL assertions that serve as oracle, our approach determines which assertions are not satisfied and, for each of them, it ranks the transformation rules according to their suspiciousness of being the faulty rule causing the failure. We have compared the effectiveness of 18 state-of-the-art techniques proposed in the literature for the suspiciousness computation of program components (e.g., statements) in the context of model transformations. The evaluation has been carried out using four case studies that differ regarding the application domains, size of metamodels, and the number and types of ATL features used. Our experiments conclude that the best techniques place the faulty transformation rule among the three most suspicious rules in around 74% of the cases. These conclusions are supported by more case studies, other than the four presented in this article, whose evaluation is available on our project's website [100].

We have also evaluated our approach by comparing it with a static approach that presented notable results [18]. The conclusion is that applying dynamic techniques based on spectra computation allows to identify the faulty rule more quickly. However, the runtime of the static technique is shorter, and it does not need any input model, since the model transformation is not executed. Therefore, both approaches are tangential and can complement each other.

Summarizing, we have proved the effectiveness in the context of model transformations of SBFL, a technique never applied before for localizing faults in this domain. We have proved it is feasible to

automate such technique in this domain, offering novel ways of debugging model transformations. Despite we have obtained good effectiveness results, further experiments can be performed as future work. For instance, we can consider helpers in the program spectra, and even each line of code could be considered as a component. In both cases, the trace model used has to be extended. Also, to break ties in the suspiciousness rankings, we could use the rules execution frequency, as some works have proposed for procedural programs [1, 67].

VERIFIABILITY

For the sake of verifiability, our prototype as well as all artifacts of the experiments are available on our project’s website [100]. For each case study, it is available the transformation and its metamodels, the OCL assertions defined, the transformation mutants together with information of the mutation operators applied and the OCL assertions that fail with each mutant, as well as the CSV files with the results generated by our program for all mutants and all OCL assertions. For the comparison study, it is available for each case study the subset of mutants used together with the matching tables generated with the approach in Reference [18] for each mutant, and the subset of OCL assertions obtained from Reference [18]. Several files with statistical results and raw data and scripts for replication are also available. Finally, the implemented prototype is available on Github: https://github.com/javitroya/SBFL_MT.

APPENDIX

A STATIC-VS-DYNAMIC COMPARISON WITH REDUCED SET OF OCL ASSERTIONS

The comparison study presented in Section 4.5 has compared our approach with the static approach by Burgueño et al. [18]. In that comparison, we have used all OCL assertions: those taken from Reference [18] and several others defined for evaluating this work. This appendix is devoted to presenting the figures and results for the comparison using only the OCL assertions defined in Reference [18]. This way we show that the new OCL assertions defined for evaluating our approach are not tailored to defeat the approach by Burgueño et al.

As shown in the second part of the third column in Table 9, 44 OCL assertions, of the total of 117 assertions created for the four case studies, have been taken from the static approach we want to compare our approach with [18]. First, of the 158 mutants we have created for the four case studies, we select those that make any of the 44 OCL assertions fail. They are a total of 104 mutants, so they are the ones to be considered in this comparison. The second part of the fifth and third columns of Table 9 display the number of mutants and OCL assertions considered in each case study for the comparison study, respectively. All the artifacts used for the comparison, namely the 104 mutants and 44 OCL assertions, together with all the matching tables generated for all case studies are available on our project’s website [100].

The approach by Burgueño et al. as well as the way we compute the EXAM values are explained in Sections 4.5.1 and 4.5.2, respectively. The descriptive statistics of the EXAM score provided by the techniques when applied to the 104 MT mutants are shown in Table 10.

First, it is worth noting that the conclusions drawn from the experiments considering all OCL assertions and mutants (cf. Sections 4.3 and 4.4) hold for this study with the 104 MT mutants, i.e., *Mountford*, *Kulcynski2*, *Ochiai*, and *Zoltar* have again the best numbers. Regarding the static technique proposed by Burgueño et al. [18], it performs worse than these techniques. In the average-case scenario, the static approach needs to inspect around 35% of the rules to locate the fault, which is much more than the 20% that needs to be inspected by the best techniques. In particular, for each case study in the average-case scenario, the static technique needs to inspect 2.17 (of 9)

Table 9. Case Studies and Artifacts for the Comparison

Case study	# Input models	# OCL assertions (/ from [18])	# Test suite ($ T = S \times O $)	# Mutants (/ comparison study)	# OCL assertions violated
UML2ER	100	14/10	1,400	18/16	90
BibTeX2DocBook	100	27/16	2,700	40/40	269
CPL2SPL	100	34/15	3,400	50/39	150
Ecore2Maude	100	42/3	4,200	50/9	155
Total	400	117/44	11,700	158/104	664

more rules in *BibTeX2DocBook* (24.1% of the MT), 0.916 (of 19) more rules in *CPL2SPL* (4.82% of the MT), 5 (of 39) more rules in *Ecore2Maude* (12.8% of the MT), and 1.58 (of 8) more rules in *UML2ER* (19.75% of the MT) compared with the best techniques in each case. Regarding the number of ties, there is not a uniform behavior. For instance, in *BibTeX2DocBook* and *CPL2SPL* there are clearly more ties in the static technique compared to the best dynamic techniques, since the difference in the EXAM score in the best- and worst-case scenarios is bigger. As for *Ecore2Maude* and *UML2ER*, the number of ties seems to be similar among both approaches.

Looking at the worst dynamic techniques, the static approach seems to behave better than some of them. Having a look at the average mean (penultimate column), it behaves much better than *Pierce* in the average-case scenario, since the latter technique needs to inspect more than 65% of the rules to locate the fault. It also performs better than *Dstar* in this scenario, since this technique needs to inspect more than 44% of the rules. Finally, the static technique by Burgueño et al. performs slightly worse than *Tarantula* in the average-case scenario but a bit better in the worst-case scenario. Therefore, for now we can conclude that the static technique may behave better than 3 dynamic techniques and clearly behaves worse than other 15 techniques, but let us delve deeper into the results.

We can further analyze the results by looking at each case study in the box-plots of Figure 8. In general, we notice that the results of the static approach are typically similar among the three scenarios, although the boxes are larger than those of most dynamic techniques, indicating a worse performance. We can appreciate that the static approach behaves normally better than *Pierce*, confirming our previous finding. As for *Dstar* and *Tarantula*, their boxes are in many plots similar to the ones of the static approach, each of them presenting slightly better results than the others in certain scenarios, so we cannot confirm the superiority of the static technique with regards to these two techniques. Indeed, for instance, in the *BibTeX2DocBook* case study, the shape of the box-plots for *Dstar* seem to be clearly better.

We have performed a statistical analysis for the comparison study, whose effect-size estimations are displayed in Table 11. We apply the same coloring as the one described in Section 4.4 for Table 8. To begin with, we can see in the *BibTeX2DocBook* case study that the four best SBFL techniques are clearly better than the static approach by Burgueño et al. [18], since the values in the row of the static approach are above 0.78 for the corresponding cells, indicating a very large difference in favor of the technique in the column. Also, the technique that seemed to be similar to the static approach, namely *Dstar*, is proved to be much better in this case study. In general, the color of the row shows that most techniques behave better than the static one.

In fact, looking at the four case studies, the numbers in the cells of the rows of the static approach and the columns with the best SBFL techniques—*Kulcynski2*, *Mountford*, *Zoltar*, and *Ochiai*—are always above 0.55, leaving no doubt that the static approach behaves worse. Besides, all these cells reveal statistical differences (p-value <0.05, displayed in boldface in the table), except for the

Table 10. Descriptive Statistics of the EXAM Score per Scenario and Case Study in the Comparison Study

	Technique	Bibtex2DocBook			CPL2SPL			Ecore2Maude			UML2ER			Average	
		mdn	mean	sd	mdn	mean	sd	mdn	mean	sd	mdn	mean	sd	mean	sd
A	Arithmetic Mean	.222	.272	.221	.066	.169	.194	.192	.267	.225	.188	.314	.239	.256	.220
	Barinel	.333	.397	.198	.184	.254	.179	.269	.315	.197	.438	.334	.164	.325	.185
	Braun-Banquet	.333	.300	.173	.079	.192	.206	.115	.167	.208	.188	.347	.349	.252	.234
	B-U & Buser	.444	.404	.226	.079	.169	.195	.115	.123	.081	.188	.345	.349	.260	.213
	Cohen	.333	.348	.207	.079	.168	.190	.192	.267	.225	.188	.314	.239	.274	.215
	Dstar	.111	.265	.258	.263	.296	.212	.788	.654	.270	.500	.550	.304	.441	.261
	Kulcynski2	.111	.173	.139	.079	.178	.202	.115	.151	.161	.188	.345	.349	.212	.213
	Mountford	.111	.203	.150	.053	.156	.199	.115	.126	.078	.188	.350	.347	.209	.194
	Ochiai	.111	.185	.143	.092	.179	.196	.115	.151	.161	.188	.345	.349	.225	.228
	Ochiai2	.444	.462	.244	.079	.175	.195	.192	.267	.225	.188	.314	.239	.304	.225
	Op2	.111	.175	.142	.105	.213	.214	.115	.167	.208	.188	.345	.349	.225	.228
	Phi	.111	.253	.218	.079	.167	.194	.192	.267	.225	.188	.314	.239	.250	.219
	Pierce	.833	.696	.274	.737	.641	.293	.731	.674	.198	.719	.625	.322	.659	.272
	Russel Rao	.222	.260	.112	.105	.226	.221	.231	.272	.180	.313	.456	.301	.304	.204
	Rogers & Tani.	.556	.518	.256	.053	.221	.274	.115	.123	.081	.125	.287	.301	.287	.228
	Simple Matching	.556	.518	.256	.053	.221	.274	.115	.123	.081	.125	.287	.301	.287	.228
	Tarantula	.333	.402	.202	.079	.160	.196	.244	.290	.197	.438	.521	.280	.343	.219
	Zoltar	.111	.173	.138	.079	.176	.199	.115	.128	.094	.188	.345	.349	.206	.195
	Burgueño'15	.389	.414	.230	.105	.204	.197	.141	.254	.251	.500	.542	.312	.354	.248
B	Arithmetic Mean	.111	.256	.216	.053	.165	.191	.026	.051	.081	.125	.181	.163	.163	.163
	Barinel	.333	.360	.215	.158	.240	.173	.051	.051	.000	.125	.132	.050	.196	.110
	Braun-Banquet	.333	.299	.173	.053	.175	.177	.026	.095	.219	.125	.326	.357	.224	.232
	B-U & Buser	.444	.404	.226	.053	.156	.169	.026	.051	.081	.125	.324	.358	.234	.209
	Cohen	.333	.332	.206	.053	.165	.188	.026	.051	.081	.125	.181	.163	.182	.160
	Dstar	.111	.265	.258	.263	.283	.196	.692	.590	.258	.500	.513	.320	.413	.258
	Kulcynski2	.111	.173	.139	.053	.165	.177	.026	.079	.170	.125	.324	.358	.185	.211
	Mountford	.111	.202	.150	.053	.144	.173	.026	.051	.081	.125	.326	.357	.181	.190
	Ochiai	.111	.185	.143	.079	.166	.171	.026	.079	.170	.125	.324	.358	.189	.211
	Ochiai2	.444	.444	.248	.053	.162	.170	.026	.051	.081	.125	.181	.163	.210	.166
	Op2	.111	.175	.142	.105	.209	.211	.026	.095	.219	.125	.324	.358	.201	.233
	Phi	.111	.237	.211	.053	.163	.192	.026	.051	.081	.125	.181	.163	.158	.162
	Pierce	.667	.592	.244	.605	.601	.268	.538	.487	.222	.438	.493	.332	.543	.267
	Rogers & Tani.	.556	.513	.258	.053	.214	.272	.026	.051	.081	.125	.266	.306	.261	.229
	Russel Rao	.111	.132	.109	.053	.174	.197	.026	.095	.219	.125	.313	.362	.179	.222
	Simple Matching	.556	.513	.258	.053	.214	.272	.026	.051	.081	.125	.266	.306	.261	.229
	Tarantula	.333	.365	.219	.053	.140	.173	.026	.026	.000	.125	.319	.362	.213	.189
	Zoltar	.111	.173	.138	.053	.163	.174	.026	.056	.097	.125	.324	.358	.179	.192
	Burgueño'15	.333	.342	.208	.105	.133	.111	.103	.126	.079	.500	.522	.315	.281	.178

(Continued)

Table 10. Continued

	Technique	Bibtex2DocBook			CPL2SPL			Ecore2Maude			UML2ER			Average	
		mdn	mean	sd	mdn	mean	sd	mdn	mean	sd	mdn	mean	sd	mean	sd
W	Arithmetic Mean	.222	.283	.250	.079	.173	.197	.231	.482	.452	.250	.446	.392	.346	.323
	Braun-Banquet	.333	.300	.174	.105	.209	.242	.205	.238	.210	.250	.368	.342	.279	.242
	B-U & Buser	.444	.404	.226	.105	.181	.228	.205	.195	.108	.250	.366	.343	.287	.226
	Cohen	.333	.359	.233	.105	.172	.193	.231	.482	.452	.250	.446	.392	.365	.318
	Dstar	.111	.265	.257	.263	.308	.237	.833	.718	.288	.500	.587	.291	.470	.268
	Kulcynski2	.111	.173	.139	.105	.190	.235	.205	.223	.168	.250	.366	.343	.238	.221
	Mountford	.111	.203	.150	.053	.168	.232	.205	.200	.100	.250	.375	.339	.237	.205
	Ochiai	.111	.185	.143	.105	.191	.229	.205	.223	.168	.250	.366	.343	.241	.221
	Ochiai2	.444	.475	.266	.105	.187	.228	.231	.482	.452	.250	.446	.392	.398	.335
	Op2	.111	.175	.142	.105	.217	.218	.205	.238	.210	.250	.366	.343	.249	.228
	Phi	.111	.264	.248	.105	.171	.197	.231	.482	.452	.250	.446	.392	.341	.322
	Pierce	1,000	.799	.327	.737	.682	.327	1,000	.862	.237	1,000	.757	.358	.775	.312
	Rogers & Tani.	.556	.524	.255	.053	.228	.277	.205	.195	.108	.125	.308	.299	.314	.235
	Russel Rao	.333	.391	.149	.105	.277	.260	.436	.449	.179	.500	.600	.280	.429	.217
	Simple Matching	.556	.524	.255	.053	.228	.277	.205	.195	.108	.125	.308	.299	.314	.235
	Tarantula	.444	.433	.234	.105	.180	.229	.462	.554	.393	.750	.723	.281	.473	.284
	Zoltar	.111	.173	.138	.105	.188	.232	.205	.200	.116	.250	.366	.343	.232	.207
	Burgueño'15	.444	.485	.284	.105	.275	.332	.179	.382	.433	.500	.563	.318	.426	.342

Ecore2Maude case study. The latter is due to the fact that the results in *Ecore2Maude* have been taken from only 9 mutants (cf. second part of the fifth column in Table 9), which are the ones that make the 3 OCL assertions considered in this case study fail, since only these assertions are defined in the evaluation of the static approach by Burgueño et al. (cf. Reference [18]—second part of third column in Table 9). Indeed, in the comparison with the complete set of OCL assertions (cf. Section 4.5.3), the cells of the *Ecore2Maude* also reveal statistical differences, since 42, instead of 3, OCL assertions are considered. Please note that the conclusions of both comparisons is the same.

The superiority of the static approach regarding *Pierce* is confirmed in the other three case studies. However, it cannot be concluded that it is better than any other of the techniques, since the rows of the static technique do not present a value <0.5 in more than one case study for any of the other techniques. Finally, we see that in the *UML2ER* case study the static approach behaves generally much worse than most techniques. An explanation can be that the static approach, based on types matching, does not behave well in the presence of rule inheritance.

In summary, we can confirm that all SBFL techniques have a better performance when locating the faulty rule than the static technique, except for *Pierce*, where the static technique behaves clearly better. Besides, the static approach normally presents more ties than the best dynamic techniques.

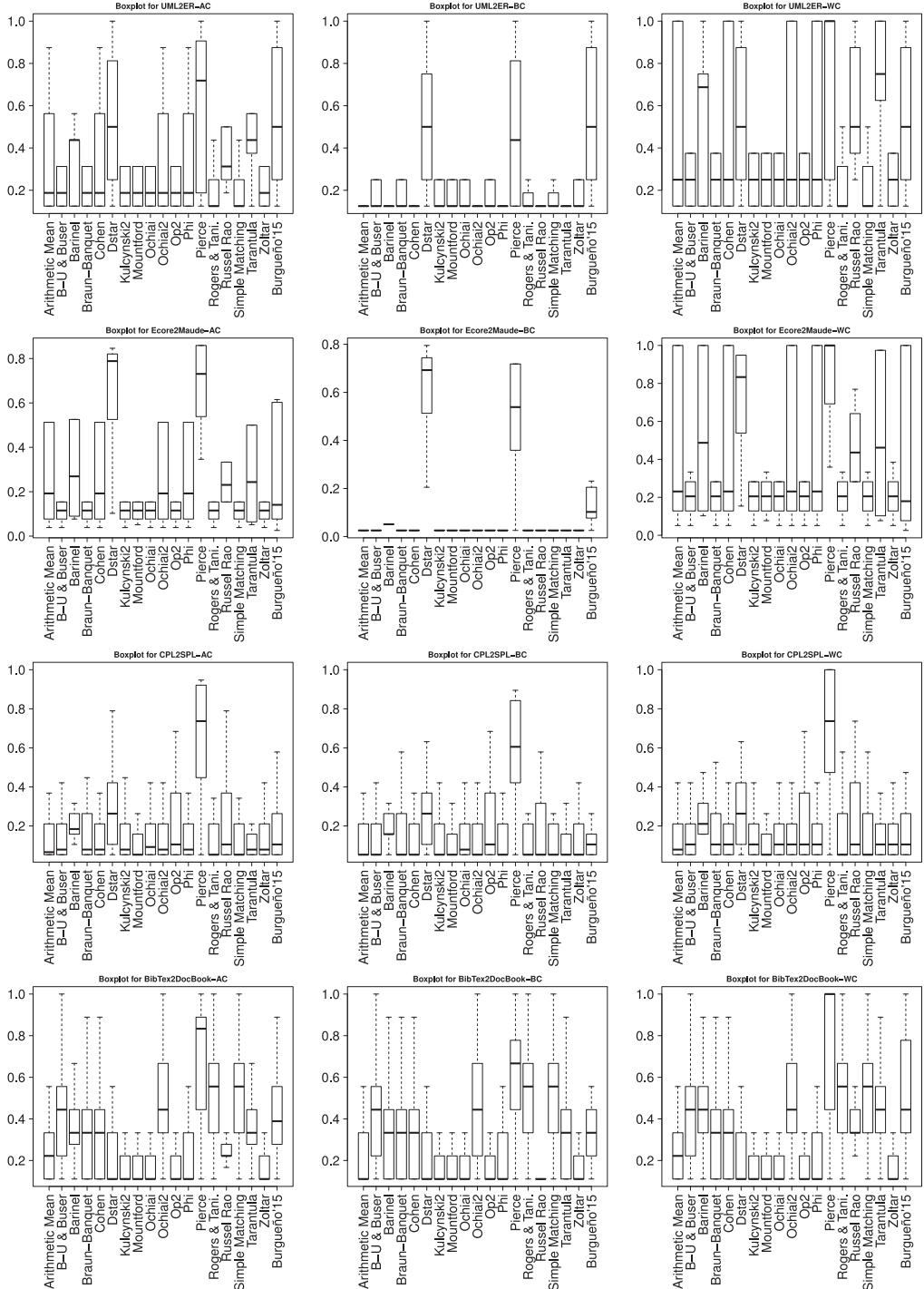


Fig. 8. Box-plot of the EXAM score of each technique per scenario and case study including Reference [18] (Burgueño'15).

REFERENCES

- [1] Rui Abreu, Alberto Gonzalez-Sanchez, and Arjan J. C. van Gemund. 2010. Exploiting count spectra for Bayesian fault localization. In *Proceedings of the 6th International Conference on Predictive Models in Software Engineering (PROMISE'10)*. ACM, Article 12, 10 pages. DOI : <http://dx.doi.org/10.1145/1868328.1868347>
- [2] Rui Abreu, Peter Zoeteweij, and Arjan J. C. van Gemund. 2009. Spectrum-based multiple fault localization. In *Proceedings of IEEE/ACM International Conference on Automated Software Engineering (ASE '09)*. IEEE Computer Society, Los Alamitos, CA, 88–99. DOI : <http://dx.doi.org/10.1109/ASE.2009.25>
- [3] Rui Abreu, Peter Zoeteweij, Rob Golsteijn, and Arjan J. C. van Gemund. 2009. A practical evaluation of spectrum-based fault localization. *J. Syst. Softw.* 82, 11 (2009), 1780–1792. DOI : <http://dx.doi.org/10.1016/j.jss.2009.06.035>
- [4] R. Abreu, P. Zoeteweij, and A. J. C. van Gemund. 2007. On the accuracy of spectrum-based fault localization. In *Proceedings of the Testing: Academic and Industrial Conference Practice and Research Techniques-MUTATION (TAICPART-MUTATION'07)*. IEEE Computer Society, Washington, DC, USA, 89–98. DOI : <http://dx.doi.org/10.1109/TAIC.PART.2007.13>
- [5] Shaukat Ali, Muhammad Zohaib Iqbal, and Andrea Arcuri. 2014. Improved heuristics for solving OCL constraints using search algorithms. In *Proceedings of the 2014 Annual Conference on Genetic and Evolutionary Computation (GECCO'14)*. ACM, New York, NY, 1231–1238. DOI : <http://dx.doi.org/10.1145/2576768.2598308>
- [6] S. Ali, M. Zohaib Iqbal, A. Arcuri, and L. C. Briand. 2013. Generating test data from OCL constraints with search techniques. *IEEE Trans. Softw. Eng.* 39, 10 (2013), 1376–1402. DOI : <http://dx.doi.org/10.1109/TSE.2013.17>
- [7] Jesús Manuel Almendros-Jiménez and Antonio Becerra-Terón. 2016. Automatic generation of ecore models for testing ATL transformations. In *Proceedings of the 6th International Conference on Model and Data Engineering (MEDI'16)*. LNCS, Vol. 9893. Springer, 16–30. DOI : http://dx.doi.org/10.1007/978-3-319-45547-1_2
- [8] Kyriakos Anastasakis, Behzad Bordbar, Geri Georg, and Indrakshi Ray. 2008. On challenges of model transformation from UML to alloy. *Softw. Syst. Model.* 9, 1 (2008). DOI : <http://dx.doi.org/10.1007/s10270-008-0110-3>
- [9] Vincent Aranega, Jean-Marie Mottu, Anne Etien, Thomas Degueule, Benoit Baudry, and Jean-Luc Dekeyser. 2015. Towards an automation of the mutation analysis dedicated to model transformation. *Softw. Test. Verif. Reliabil.* 25, 5–7 (2015), 653–683. DOI : <http://dx.doi.org/10.1002/stvr.1532>
- [10] Thorsten Arendt, Enrico Biermann, Stefan Jurack, Christian Krause, and Gabriele Taentzer. 2010. Henshin: Advanced concepts and tools for in-place EMF model transformations. In *Proceedings of ACM/IEEE 21st International Conference on Model Driven Engineering Languages and Systems (MODELS'10)*, Vol. 6394. 121–135.
- [11] Fatmah Yousef Assiri and James M. Bieman. 2017. Fault localization for automated program repair: Effectiveness, performance, repair correctness. *Softw. Qual. J.* 25, 1 (2017), 171–199. DOI : <http://dx.doi.org/10.1007/s11219-016-9312-z>
- [12] ATL. 2006. ATL Zoo. Retrieved from <http://www.eclipse.org/atl/atlTransformations>.
- [13] E. T. Barr, M. Harman, P. McMinn, M. Shahbaz, and S. Yoo. 2015. The oracle problem in software testing: A survey. *IEEE Trans. Softw. Eng.* 41, 5 (May 2015), 507–525. DOI : <http://dx.doi.org/10.1109/TSE.2014.2372785>
- [14] Benoit Baudry, Trung Dinh-Trong, Jean-Marie Mottu, Devon Simmonds, Robert France, Sudipto Ghosh, Franck Fleurey, and Yves Le Traon. 2006. Model transformation testing challenges. In *Proceedings of the ECMDA Workshop on Integration of Model Driven Development and Model Driven Testing*. <http://www.cs.colostate.edu/france/publications/TransTestin>
- [15] Amine Benelallam, Abel Gómez, Massimo Tisi, and Jordi Cabot. 2015. Distributed model-to-model transformation with ATL on MapReduce. In *Proceedings of the 2015 ACM SIGPLAN International Conference on Software Language Engineering (SLE'15)*. ACM, New York, NY, 37–48. DOI : <http://dx.doi.org/10.1145/2814251.2814258>
- [16] Marco Brambilla, Jordi Cabot, and Manuel Wimmer. 2012. *Model-Driven Software Engineering in Practice*. Morgan & Claypool.
- [17] Erwan Brottier, Franck Fleurey, Jim Steel, Benoit Baudry, and Yves Le Traon. 2006. Metamodel-based test generation for model transformations: An algorithm and a tool. In *Proceedings of Annual IEEE International Symposium on Software Reliability Engineering (ISSRE'06)*. 85–94.
- [18] L. Burgueño, J. Troya, M. Wimmer, and A. Vallecillo. 2015. Static fault localization in model transformations. *IEEE Trans. Softw. Eng.* 41, 5 (May 2015), 490–506.
- [19] Loli Burgueño, Manuel Wimmer, Javier Troya, and Antonio Vallecillo. 2013. TractsTool: Testing MTs based on contracts. In *Invited Talks, Demonstration Session, Poster Session, and ACM Student Research Competition (MODELS'13)*. CEUR.
- [20] Daniel Calegari, Carlos Luna, Nora Szasz, and Alvaro Tasistro. 2010. A type-theoretic framework for certified model transformations. In *Proceedings of Brazilian Symposium on Formal Methods (SBMF'10)*. 112–127. DOI : http://dx.doi.org/10.1107/978-3-642-19829-8_8
- [21] Eric Cariou, Raphaël Marvie, Lionel Seinturier, and Laurence Duchien. 2004. OCL for the specification of model transformation contracts. In *Proceedings of the OCL and Model Driven Engineering Workshop*. <http://web.univ-pau.fr/ecariou/papers/workshop-ocl-mde-uml2004-paper>.

- [22] Zheng Cheng, Rosemary Monahan, and James F. Power. 2015. *A Sound Execution Semantics for ATL via Translation Validation*. Springer International Publishing, Cham, 133–148. DOI : http://dx.doi.org/10.1007/978-3-319-21155-8_11
- [23] Zheng Cheng and Massimo Tisi. 2017. *A Deductive Approach for Fault Localization in ATL Model Transformations*. Springer, Berlin, 300–317. DOI : http://dx.doi.org/10.1007/978-3-662-54494-5_17
- [24] Antonio Cicchetti, Davide Di Ruscio, Romina Eramo, and Alfonso Pierantonio. 2011. JTL: A bidirectional and change propagating transformation language. In *Software Language Engineering*. LNCS, Vol. 6563. Springer, 183–202. DOI : http://dx.doi.org/10.1007/978-3-642-19440-5_11
- [25] Cauê Clasen, Marcos Didonet Del Fabro, and Massimo Tisi. 2012. Transforming very large models in the cloud: A research roadmap. In *Proceedings of the 1st International Workshop on Model-Driven Engineering on and for the Cloud*. Springer, Copenhagen, Denmark. <https://hal.inria.fr/hal-00711524>
- [26] Manuel Clavel, Francisco Durán, Steven Eker, Patrick Lincoln, Narciso Martí-Oliet, José Meseguer, and Carolyn Talcott. 2007. *All About Maude—A High-Performance Logical Framework*. LNCS, Vol. 4350. Springer. DOI : <http://dx.doi.org/10.1007/978-3-540-71999-1>
- [27] György Csértán, Gábor Huszér, István Majzik, Zsigmond Pap, András Pataricza, and Dániel Varró. 2002. VIATRA—visual automated transformations for formal verification and validation of UML models. In *Proceedings of the 17th International Conference on Automated Software Engineering (ASE'02)*. IEEE/ACM, 267–270. DOI : <http://dx.doi.org/10.1109/ASE.2002.1115027>
- [28] Krzysztof Czarnecki and Simon Helsen. 2006. Feature-based survey of model transformation approaches. *IBM Syst. J.* 45, 3 (2006), 621–646.
- [29] Alberto Rodrigues da Silva. 2015. Model-driven engineering: A survey supported by the unified conceptual model. *Comput. Lang. Syst. Struct.* 43 (2015), 139–155. DOI : <http://dx.doi.org/10.1016/j.cl.2015.06.001>
- [30] Juan de Lara and Hans Vangheluwe. 2002. AToM3: A tool for multi-formalism and meta-modelling. In *Proceedings of the 5th International Conference on Fundamental Approaches to Software Engineering (FASE'02)*. LNCS, Vol. 2306. Springer, 174–188. DOI : http://dx.doi.org/10.1007/3-540-45923-5_12
- [31] Joaquín Derrac, Salvador García, Daniel Molina, and Francisco Herrera. 2011. A practical tutorial on the use of nonparametric statistical tests as a methodology for comparing evolutionary and swarm intelligence algorithms. *Swarm Evol. Comput.* 1, 1 (2011), 3–18. DOI : <http://dx.doi.org/10.1016/j.swevo.2011.02.002>
- [32] Francisco Durán, Steffen Zschaler, and Javier Troya. 2013. On the reusable specification of non-functional properties in DSLs. In *Proceedings of the 5th International Conference on Software Language Engineering (SLE'12). Revised Selected Papers* (LNCS). Springer, 332–351. DOI : http://dx.doi.org/10.1007/978-3-642-36089-3_19
- [33] Jean-Rémy Falleri, Marianne Huchard, and Clémentine Nebut. 2006. Towards a traceability framework for model transformations in Kermeta. In *Proceedings of the European Conference on Model Driven Architecture—Traceability Workshop (ECMDA-TW'06)*. 31–40. <https://hal-lirmm.ccsd.cnrs.fr/lirmm-00102855>
- [34] M. Fleck, J. Troya, M. Kessentini, M. Wimmer, and B. Alkhazi. 2017. Model transformation modularization as a many-objective optimization problem. *IEEE Trans. Softw. Eng.* 43, 11 (2017), 1009–1032. DOI : <http://dx.doi.org/10.1109/TSE.2017.2654255>
- [35] Martin Fleck, Javier Troya, and Manuel Wimmer. 2015. Marrying search-based optimization and model transformation technology. In *Proceedings of First North American Search Based Software Engineering Symposium (NasBASE'15)*. 1–16.
- [36] Martin Fleck, Javier Troya, and Manuel Wimmer. 2016. Search-based model transformations. *J. Softw. Evol. Process* 28, 12 (2016), 1081–1117. DOI : <http://dx.doi.org/10.1002/smri.1804>
- [37] Martin Fleck, Javier Troya, and Manuel Wimmer. 2016. Search-based model transformations with MOMoT. In *Proceedings of Theory and Practice of Model Transformations (ICMT'16)*. Springer, 79–87.
- [38] Franck Fleurey, Benoit Baudry, Pierre-Alain Muller, and Yves Le Traon. 2009. Qualifying input test data for model transformations. *Softw. Syst. Model.* 8, 2 (2009), 185–203. DOI : <http://dx.doi.org/10.1007/s10270-007-0074-8>
- [39] Franck Fleurey, Benoit Baudry, Pierre-Alain Muller, and Yves Le Traon. 2009. Qualifying input test data for model transformations. *Softw. Syst. Model.* 8, 2 (2009), 185–203. DOI : <http://dx.doi.org/10.1007/s10270-007-0074-8>
- [40] Milton Friedman. 1940. A comparison of alternative tests of significance for the problem of m rankings. *Ann. Math. Stat.* 11, 1 (1940), 86–92. <http://www.jstor.org/stable/2235971>
- [41] Pau Giner and Vicente Pelechano. 2009. Test-driven development of model transformations. In *Proceedings of ACM/IEEE 21st International Conference on Model Driven Engineering Languages and Systems (MODELS'09)*. LNCS, Vol. 5795. Springer, 748–752.
- [42] Martin Gogolla and Antonio Vallecillo. 2011. *Tractable Model Transformation Testing*. Springer, Berlin, 221–235. DOI : http://dx.doi.org/10.1007/978-3-642-21470-7_16
- [43] C. Gong, Z. Zheng, W. Li, and P. Hao. 2012. Effects of class imbalance in test suites: An empirical study of spectrum-based fault localization. In *Proceedings of the IEEE 36th Annual Computer Software and Applications Conference Workshops*. 470–475. DOI : <http://dx.doi.org/10.1109/COMPSACW.2012.89>

- [44] Carlos A. González and Jordi Cabot. 2012. *ATLTest: A White-Box Test Generation Approach for ATL Transformations*. Springer, 449–464. DOI : http://dx.doi.org/10.1007/978-3-642-33666-9_29
- [45] Joel Greenyer and Ekkart Kindler. 2010. Comparing relational model transformation technologies: Implementing query/view/transformation with triple graph grammars. *Softw. Syst. Model.* 9, 1 (2010), 21–46.
- [46] Esther Guerra. 2012. Specification-driven test generation for model transformations. In *Theory and Practice of Model Transformations*, Zhenjiang Hu and Juan de Lara (Eds.). Springer, Berlin, 40–55.
- [47] Esther Guerra, Juan de Lara, Manuel Wimmer, Gerti Kappel, Angelika Kusel, Werner Retschitzegger, Johannes Schönböck, and Wieland Schwinger. 2013. Automated verification of model transformations based on visual contracts. *Automat. Softw. Eng.* 20, 1 (2013), 5–46. DOI : <http://dx.doi.org/10.1007/s10515-012-0102-y>
- [48] Esther Guerra and Mathias Soeken. 2015. Specification-driven model transformation testing. *Softw. Syst. Model.* 14, 2 (2015), 623–644. DOI : <http://dx.doi.org/10.1007/s10270-013-0369-x>
- [49] M. Hamill and K. Goseva-Popstojanova. 2009. Common trends in software fault and failure data. *IEEE Trans. Softw. Eng.* 35, 4 (2009), 484–496. DOI : <http://dx.doi.org/10.1109/TSE.2009.3>
- [50] Mary Jean Harrold, Gregg Rothenmel, Kent Sayre, Rui Wu, and Liu Yi. 2000. An empirical investigation of the relationship between spectra differences and regression faults. *Softw., Test., Verif., Reliabil.* 10, 3 (2000), 171–194. DOI : [http://dx.doi.org/10.1002/1099-1689\(200009\)10:3<171::AID-STVR209>3.0.CO;2-J](http://dx.doi.org/10.1002/1099-1689(200009)10:3<171::AID-STVR209>3.0.CO;2-J)
- [51] Xiao He, Xing Chen, Sibo Cai, Ying Zhang, and Gang Huang. 2018. Testing bidirectional model transformation using metamorphic testing. *Information and Software Technology* (2018). DOI : <http://dx.doi.org/10.1016/j.infsof.2018.07.010>
- [52] Myles Hollander, Douglas A. Wolfe, and Eric Chicken. 2013. *Nonparametric Statistical Methods*. John Wiley & Sons.
- [53] Sture Holm. 1979. A simple sequentially rejective multiple test procedure. *Scand. J. Statist.* 6, 2 (1979), 65–70. SJSADG
- [54] INRIA. 2005. ATL Transformation Example: BibTeXML to DocBook. Retrieved from [https://www.eclipse.org/atl/atlTransformations/BibTeXML2DocBook/ExampleBibTeXML2DocBook\[v00.01\].pdf](https://www.eclipse.org/atl/atlTransformations/BibTeXML2DocBook/ExampleBibTeXML2DocBook[v00.01].pdf).
- [55] Tom Janssen, Rui Abreu, and Arjan J. C. van Gemund. 2009. Zoltar: A spectrum-based fault localization tool. In *Proceedings of the 2009 ESEC/FSE Workshop on Software Integration and Evolution @ Runtime (SINTER'09)*. ACM, New York, NY, 23–30. DOI : <http://dx.doi.org/10.1145/1596495.1596502>
- [56] Jean-Marc Jézéquel, Olivier Barais, and Franck Fleurey. 2011. Model driven language engineering with Kermeta. In *Generative and Transformational Techniques in Software Engineering III*. LNCS, Vol. 6491. Springer, 201–221. DOI : http://dx.doi.org/10.1007/978-3-642-18023-1_5
- [57] Yue Jia and Mark Harman. 2009. Higher order mutation testing. *Inf. Softw. Technol.* 51, 10 (2009), 1379–1393. DOI : <http://dx.doi.org/10.1016/j.infsof.2009.04.016>
- [58] Yue Jia and Mark Harman. 2011. An analysis and survey of the development of mutation testing. *IEEE Trans. Softw. Eng.* 37, 5 (2011), 649–678. DOI : <http://dx.doi.org/10.1109/TSE.2010.62>
- [59] James A. Jones and Mary Jean Harrold. 2005. Empirical evaluation of the tarantula automatic fault-localization technique. In *Proceedings of the 20th IEEE/ACM International Conference on Automated Software Engineering (ASE'05)*. ACM, New York, NY, 273–282. DOI : <http://dx.doi.org/10.1145/1101908.1101949>
- [60] Frédéric Jouault. 2005. Loosely coupled traceability for ATL. In *Workshop Proceedings of European Conference on Model Driven Architecture—Traceability Workshop (ECMDA'05)*.
- [61] Frédéric Jouault, Freddy Allilaire, Jean Bézivin, and Ivan Kurtev. 2008. ATL: A model transformation tool. *Sci. Comput. Program.* 72, 1-2 (2008), 31–39.
- [62] Frédéric Jouault, Freddy Allilaire, Jean Bézivin, Ivan Kurtev, and Patrick Valduriez. 2006. ATL: A QVT-like transformation language. In *Companion to the 21st ACM SIGPLAN Symposium on Object-oriented Programming Systems, Languages, and Applications (OOPSLA'06)*. ACM, 719–720. DOI : <http://dx.doi.org/10.1145/1176617.1176691>
- [63] Frédéric Jouault, Jean Bézivin, Charles Consel, Ivan Kurtev, and Fabien Latry. 2006. Building DSLs with AMMA/ATL, a case study on SPL and CPL telephony languages. In *Proceedings of the ECOOP Workshop on Domain-Specific Program Development*. Nantes, France. <https://hal.inria.fr/inria-00353580>
- [64] Begashaw Gezu Kirsie. 2010. *Guideline and Evaluation of Model Transformation Engineering Approaches*. Master's thesis. KTH Industrial Engineering and Management, Sweden.
- [65] Dimitrios S. Kolovos, Louis M. Rose, Nicholas Matragkas, Richard F. Paige, Esther Guerra, Jesús Sánchez Cuadrado, Juan De Lara, István Ráth, Dániel Varró, Massimo Tisi, and Jordi Cabot. 2013. A research roadmap towards achieving scalability in model driven engineering. In *Proceedings of the Workshop on Scalability in Model Driven Engineering (BigMDE'13)*. ACM, New York, NY, Article 2, 10 pages. DOI : <http://dx.doi.org/10.1145/2487766.2487768>
- [66] Thomas Kühne. 2006. Matters of (meta-) modeling. *Softw. Syst. Model.* 5, 4 (2006), 369–385. DOI : <http://dx.doi.org/10.1007/s10270-006-0017-9>
- [67] H. J. Lee, L. Naish, and K. Ramamohanrao. 2010. Effective software bug localization using spectral frequency weighting function. In *Proceedings of the IEEE 34th Annual Computer Software and Applications Conference*. 218–227. DOI : <http://dx.doi.org/10.1109/COMPSAC.2010.26>

- [68] Chao Liu, Long Fei, Xifeng Yan, Jiawei Han, and Samuel P. Midkiff. 2006. Statistical debugging: A hypothesis testing-based approach. *IEEE Trans. Softw. Eng.* 32, 10 (Oct. 2006), 831–848. DOI : <http://dx.doi.org/10.1109/TSE.2006.105>
- [69] Lucia, F. Thung, D. Lo, and L. Jiang. 2012. Are faults localizable? In *Proceedings of the 2012 9th IEEE Working Conference on Mining Software Repositories (MSR'12)*. 74–77. DOI : <http://dx.doi.org/10.1109/MSR.2012.6224302>
- [70] Lucia Lucia, David Lo, Lingxiao Jiang, Ferdinand Thung, and Aditya Budi. 2014. Extended comprehensive study of association measures for fault localization. *J. Softw. Evol. Process* 26, 2 (2014), 172–219. DOI : <http://dx.doi.org/10.1002/smri.1616>
- [71] Levi Lúcio, Moussa Amrani, Jürgen Dingel, Leen Lambers, Rick Salay, Gehan Selim, Eugene Syriani, and Manuel Wimmer. 2016. Model transformation intents and their properties. *Softw. Syst. Model.* 15, 3 (2016), 647–684.
- [72] Jochen Ludewig. 2003. Models in software engineering – an introduction. *Softw. Syst. Model.* 2, 1 (2003), 5–14. DOI : <http://dx.doi.org/10.1007/s10270-003-0020-3>
- [73] Xiaoguang Mao, Yan Lei, Ziyi Dai, Yuhua Qi, and Chengsong Wang. 2014. Slice-based statistical fault localization. *J. Syst. Softw.* 89 (Mar. 2014), 51–62. DOI : <http://dx.doi.org/10.1016/j.jss.2013.08.031>
- [74] A. E. Maxwell and A. E. G. Pilliner. 1968. Deriving coefficients of reliability and agreement for ratings. *Br. J. Math. Statist. Psych.* 21, 1 (1968), 105–116. DOI : <http://dx.doi.org/10.1111/j.2044-8317.1968.tb00401.x>
- [75] S. J. Mellor, K. Scott, A. Uhl, D. Weise, and R. M. Soley. 2004. *MDA Distilled: Principles of Model-driven Architecture*. Vol. 88. Addison-Wesley.
- [76] Antonio Moreno-Delgado, Francisco Durán, Steffen Zschaler, and Javier Troya. 2014. Modular DSLs for flexible analysis: An e-motions reimplementation of palladio. In *Proceedings of the 10th European Conference on Modelling Foundations and Applications (ECMFA'14)*. LNCS. Springer, 132–147.
- [77] Jean-Marie Mottu, Benoit Baudry, and Yves Le Traon. 2006. *Mutation Analysis Testing for Model Transformations*. Springer, Berlin., 376–390. DOI : http://dx.doi.org/10.1007/11787044_28
- [78] Lee Naish, Hua Jie Lee, and Kotagiri Ramamohanarao. 2011. A model for spectra-based software diagnosis. *ACM Trans. Softw. Eng. Methodol.* 20, 3, Article 11 (Aug. 2011), 32 pages. DOI : <http://dx.doi.org/10.1145/2000791.2000795>
- [79] L. Naish, Neelofer, and K. Ramamohanarao. 2015. Multiple bug spectral fault localization using genetic programming. In *Proceedings of the 2015 24th Australasian Software Engineering Conference*. 11–17. DOI : <http://dx.doi.org/10.1109/ASWEC.2015.12>
- [80] Bentley James Oakes, Javier Troya, Levi Lúcio, and Manuel Wimmer. 2018. Full contract verification for ATL using symbolic execution. *J. Softw. Syst. Model.* 17, 3 (2018), 815–849. DOI : <http://dx.doi.org/10.1007/s10270-016-0548-7>
- [81] Elmahdi Omar, Sudipto Ghosh, and Darrell Whitley. 2017. Subtle higher order mutants. *Inf. Softw. Technol.* 81 (2017), 3–18. DOI : <http://dx.doi.org/10.1016/j.infsof.2016.01.016>
- [82] Magnus Persson, Martin Törngren, Ahsan Qamar, Jonas Westman, Matthias Biehl, Stavros Tripakis, Hans Vangheluwe, and Joachim Denil. 2013. A characterization of integrated multi-view modeling in the context of embedded and cyber-physical systems. In *Proceedings of the 11th ACM International Conference on Embedded Software (EMSOFT'13)*. IEEE Press, Los Alamitos, CA, Article 10, 10 pages. <http://dl.acm.org/citation.cfm?id=2555754.2555764>
- [83] Iman Poernomo and Jeffrey Terrell. 2010. Correct-by-construction model transformations from partially ordered specifications in Coq. In *Proceedings of International Conference on Formal Engineering Methods (ICFEM'10)*. 56–73. DOI : http://dx.doi.org/10.1007/978-3-642-16901-4_6
- [84] Eclipse Modeling Project. 2015. Atlas Transformation Language—ATL. Retrieved from <http://eclipse.org/atl>.
- [85] Yuhua Qi, Xiaoguang Mao, Yan Lei, and Chengsong Wang. 2013. Using automated program repair for evaluating the effectiveness of fault localization techniques. In *Proceedings of the 2013 International Symposium on Software Testing and Analysis (ISSTA'13)*. ACM, New York, NY, 191–201. DOI : <http://dx.doi.org/10.1145/2483760.2483785>
- [86] Jose E. Rivera, Francisco Duran, and Antonio Valleccillo. 2009. A graphical approach for modeling time-dependent behavior of DSLs. In *Proceedings of the IEEE Symposium on Visual Languages and Human-Centric Computing (VL/HCC'09)*. IEEE, 51–55. DOI : <http://dx.doi.org/10.1109/VLHCC.2009.5295300>
- [87] Louis M. Rose, Markus Herrmannsdoerfer, Steffen Mazanek, Pieter Van Gorp, Sebastian Buchwald, Tassilo Horn, Elina Kalnina, Andreas Koch, Kevin Lano, Bernhard Schätz, and Manuel Wimmer. 2014. Graph and model transformation tools for model migration. *Softw. Syst. Model.* 13, 1 (2014), 323–359. DOI : <http://dx.doi.org/10.1007/s10270-012-0245-0>
- [88] K. Rustan and M. Leino. 2008. *This Is Boogie 2*. Technical Report. Manuscript KRML 178.
- [89] Jesús Sánchez-Cuadrado and Jesús García-Molina. 2009. Modularization of model transformations through a phasing mechanism. *Softw. Syst. Model.* 8, 3 (Jul 2009), 325–345. DOI : <http://dx.doi.org/10.1007/s10270-008-0093-0>
- [90] Jesús Sánchez-Cuadrado, Esther Guerra, and Juan de Lara. 2018. Quick fixing ATL transformations with speculative analysis. *Softw. Syst. Model.* 17, 3 (2018), 779–813. DOI : <http://dx.doi.org/10.1007/s10270-016-0541-1>
- [91] Jesús Sánchez-Cuadrado, Esther Guerra, and Juan de Lara. 2017. Static analysis of model transformations. *IEEE Trans. Softw. Eng.* 43, 9 (2017), 868–897. DOI : <http://dx.doi.org/10.1109/TSE.2016.2635137>

- [92] J. Sánchez-Cuadrado, E. Guerra, J. de Lara, R. Clarisó, and J. Cabot. 2017. Translating target to source constraints in model-to-model transformations. In *Proceedings of the ACM/IEEE 20th International Conference on Model Driven Engineering Languages and Systems (MODELS'17)*. 12–22. DOI : <http://dx.doi.org/10.1109/MODELS.2017.12>
- [93] S. Segura, G. Fraser, A. Sanchez, and A. Ruiz-Cortes. 2016. A survey on metamorphic testing. *IEEE Trans. Softw. Eng.* 42, 9 (2016), 805–824. DOI : <http://dx.doi.org/10.1109/TSE.2016.2532875>
- [94] Gehan M. K. Selim, Shige Wang, James R. Cordy, and Juergen Dingel. 2012. *Model Transformations for Migrating Legacy Models: An Industrial Case Study*. Springer, Berlin, 90–101. DOI : http://dx.doi.org/10.1007/978-3-642-31491-9_9
- [95] Sagar Sen, Benoit Baudry, and Jean-Marie Mottu. 2009. Automatic model generation strategies for model transformation testing. In *Proceedings of the International Conference on Manufacturing Technologies (ICMT'09)*. LNCS, Vol. 5563. Springer, 148–164.
- [96] Shane Sendall and Wojtek Kozaczynski. 2003. Model transformation: The heart and soul of model-driven software development. *IEEE Softw.* 20, 5 (2003), 42–45.
- [97] Gabriele Taentzer. 2003. AGG: A graph transformation environment for modeling and validation of software. In *Proceedings of the 2nd Workshop on Applications of Graph Transformations with Industrial Relevance (AGTIVE'03)*. LNCS, Vol. 3062. Springer, 446–453. DOI : http://dx.doi.org/10.1007/978-3-540-25959-6_35
- [98] Javier Troya, Alexander Bergmayr, Loli Burgueno, and Manuel Wimmer. 2015. Towards systematic mutations for and with ATL model transformations. In *Proceedings of the IEEE 8th International Conference on Software Testing, Verification and Validation Workshops (ICSTW'15)*. 1–10. DOI : <http://dx.doi.org/10.1109/ICSTW.2015.7107455>
- [99] Javier Troya, Sergio Segura, José Antonio Parejo, and Antonio Ruiz-Cortes. 2017. An approach for debugging model transformations applying spectrum-based fault localization. In *XXII Jornadas de Ingeniería del Software y Bases de Datos (JISBD'17)*. https://biblioteca.sistedes.es/submissions/uploaded-files/JISBD_2017_paper_32.pdf.
- [100] Javier Troya, Sergio Segura, José Antonio Parejo, and Antonio Ruiz-Cortés. 2017. Spectrum-Based Fault Localization in Model Transformations. Retrieved from <https://gestionproyectos.us.es/projects/itim/wiki>.
- [101] Javier Troya, Sergio Segura, and Antonio Ruiz-Cortés. 2018. Automated inference of likely metamorphic relations for model transformations. *J. Syst. Softw.* 136 (2018), 188–208. DOI : <http://dx.doi.org/10.1016/j.jss.2017.05.043>
- [102] Javier Troya and Antonio Vallecillo. 2011. A rewriting logic semantics for ATL. *J. Obj. Technol.* 10, 5 (2011), 1–29. DOI : <http://dx.doi.org/10.5381/jot.2011.10.1.a5>
- [103] Javier Troya and Antonio Vallecillo. 2014. Specification and simulation of queuing network models using domain-specific languages. *Comput. Stand. Intef.* 36, 5 (2014), 863–879. DOI : <http://dx.doi.org/10.1016/j.csi.2014.01.002>
- [104] Antonio Vallecillo and Martin Gogolla. 2012. Typing model transformations using tracts. In *Proceedings of the 5th International Conference on Theory and Practice of Model Transformations (ICMT'12)*. Springer, 56–71. DOI : http://dx.doi.org/10.1007/978-3-642-30476-7_4
- [105] András Vargha and Harold D. Delaney. 2000. A critique and improvement of the CL common language effect size statistics of McGraw and Wong. *J. Educ. Behav. Stat.* 25, 2 (2000), 101–132.
- [106] D. Wagelaar. 2014. Using ATL/EMFTVM for import/export of medical data. In *Proceedings of the 2nd Software Development Automation Conference*.
- [107] Leonard Muellner Norman Walsh. 1999. *DocBook: The Definitive Guide*. O'Reilly & Associates. <http://www.docbook.org/tgd/>
- [108] Jos Warmer and Anneke Kleppe. 2003. *The Object Constraint Language: Getting Your Models Ready for MDA*. Addison-Wesley.
- [109] Manuel Wimmer and Loli Burgueño. 2013. Testing M2T/T2M transformations. In *Proceedings of the International Conference on Model Driven Engineering Languages and Systems (MODELS'13)*. Springer, 203–219.
- [110] Manuel Wimmer, Gerti Kappel, Johannes Schönböck, Angelika Kusel, Werner Retschitzegger, and Wieland Schwinger. 2009. A Petri Net based debugging environment for QVT relations. In *Proceedings of the IEEE/ACM International Conference on Automated Software Engineering (ASE'09)*. IEEE, 3–14.
- [111] Manuel Wimmer, Salvador Martínez, Frédéric Jouault, and Jordi Cabot. 2012. A catalogue of refactorings for model-to-model transformations. *J. Obj. Technol.* 11, 2 (Aug. 2012), 2:1–40. DOI : <http://dx.doi.org/10.5381/jot.2012.11.2.a2>
- [112] Claes Wohlin, Per Runeson, Martin Höst, Magnus C. Ohlsson, and Björn Regnell. 2012. *Experimentation in Software Engineering*. Springer.
- [113] W. E. Wong, V. Debroy, R. Gao, and Y. Li. 2014. The DStar method for effective software fault localization. *IEEE Trans. Reliabil.* 63, 1 (Mar. 2014), 290–308. DOI : <http://dx.doi.org/10.1109/TR.2013.2285319>
- [114] W. E. Wong, V. Debroy, Y. Li, and R. Gao. 2012. Software fault localization using DStar (D*). In *Proceedings of the IEEE 6th International Conference on Software Security and Reliability*. 21–30. DOI : <http://dx.doi.org/10.1109/SERIE.2012.12>
- [115] W. Eric Wong, Ruizhi Gao, Yihao Li, Rui Abreu, and Franz Wotawa. 2016. A survey on software fault localization. *IEEE Trans. Softw. Eng.* 42, 8 (2016), 707–740. DOI : <http://dx.doi.org/10.1109/TSE.2016.2521368>

- [116] Xiaoyuan Xie. 2012. *On the Analysis of Spectrum-based Fault Localization*. Ph.D. Dissertation. Faculty of Information and Communication Technologies, Swinburne University of Technology, Australia.
- [117] Xiaoyuan Xie, Tsong Yueh Chen, Fei-Ching Kuo, and Baowen Xu. 2013. A theoretical analysis of the risk evaluation formulas for spectrum-based fault localization. *ACM Trans. Softw. Eng. Methodol.* 22, 4, Article 31 (Oct. 2013), 40 pages. DOI : <http://dx.doi.org/10.1145/2522920.2522924>
- [118] X. Xue and A. S. Namin. 2013. How significant is the effect of fault interactions on coverage-based fault localizations? In *Proceedings of the 2013 ACM/IEEE International Symposium on Empirical Software Engineering and Measurement*. 113–122. DOI : <http://dx.doi.org/10.1109/ESEM.2013.22>
- [119] Shin Yoo. 2012. Evolving human competitive spectra-based fault localisation techniques. In *Search Based Software Engineering*, Gordon Fraser and Jerffeson Teixeira de Souza (Eds.). Springer, Berlin, 244–258.
- [120] Yanbing Yu, James A. Jones, and Mary Jean Harrold. 2008. An empirical study of the effects of test-suite reduction on fault localization. In *Proceedings of the 30th International Conference on Software Engineering (ICSE'08)*. ACM, New York, NY, 201–210. DOI : <http://dx.doi.org/10.1145/1368088.1368116>
- [121] Xiangyu Zhang, Sriraman Tallam, Neelam Gupta, and Rajiv Gupta. 2007. Towards locating execution omission errors. In *Proceedings of the 28th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI'07)*. ACM, 415–424. DOI : <http://dx.doi.org/10.1145/1250734.1250782>
- [122] Zhenyu Zhang, W. K. Chan, T. H. Tse, Peifeng Hu, and Xinming Wang. 2009. Is non-parametric hypothesis testing model robust for statistical fault localization? *Inf. Softw. Technol.* 51, 11 (2009), 1573–1585. DOI : <http://dx.doi.org/10.1016/j.infsof.2009.06.013>
- [123] Álvaro Jiménez, Juan M. Vara, Verónica A. Bollati, and Esperanza Marcos. 2015. MeTAGeM-trace: Improving trace generation in model transformation by leveraging the role of transformation models. *Sci. Comput. Program.* 98 (2015), 3–27. DOI : <http://dx.doi.org/10.1016/j.scico.2014.09.003>

Received October 2017; revised May 2018; accepted July 2018