

# Spectrum-Based Fault Localization in Model Transformations

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

Model transformations play a cornerstone role in Model-Driven Engineering 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 paper 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 *Kulczynski2*, *Mountford*, *Ochiai* and *Zoltar*, lead the debugger to inspect a maximum of three rules in order 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. 2017. Spectrum-Based Fault Localization in Model Transformations. *ACM Trans. Softw. Eng. Methodol.* 1, 1, Article 1 (January 2017), 49 pages.

<https://doi.org/0000001.0000001>

## 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 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

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: Department of Computer Languages and Systems, Universidad de Sevilla, Spain.

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](mailto:permissions@acm.org).

© 2017 Association for Computing Machinery.

1049-331X/2017/1-ART1 \$15.00

<https://doi.org/0000001.0000001>

50 strong theoretical foundations and applicability to real-world problems [71]. Most MT languages  
51 are composed of model transformation rules<sup>1</sup>. Each MT rule deals with the construction of part of  
52 the target model. They match input elements from the source model and generate output elements  
53 that compose the target model.

54 The correctness of software built using MDE techniques typically relies on the correctness of the  
55 operations executed using MTs. For this reason, it is critical in MDE to maintain and test them as it  
56 is done with source code in classical software engineering. However, checking whether the output  
57 of a MT is correct is a manual and error-prone task, which suffers the *oracle problem*. The oracle  
58 problem refers to, given an input for a system, the challenge of distinguishing the corresponding  
59 desired, correct behavior from potentially incorrect behavior [13]. In order to alleviate this problem  
60 in the model transformation domain, the formal specification of MTs has been proposed by the  
61 definition and use of *contracts* [14, 18, 21, 105], i.e., assertions that the execution of the MTs must  
62 satisfy. These assertions can be specified on the models resulting from the MTs, the models serving  
63 as input for the MTs, or both, and they can be tested in a black-box manner. These assertions are  
64 typically defined using the Object Constraint Language (OCL) [109].

65 However, even when using the assertions as oracle to test if MTs are faulty, it is still challenging  
66 to debug them and locate what parts of the MTs are wrong. The situation gets more critical as  
67 the size and complexity of MTs grow, where manual debugging is no longer possible, such as in  
68 aviation, medical data processing [107], automotive industry [95] or embedded and cyber-physical  
69 systems [83]. Therefore, there is an increasing need to count on methods, mechanisms and tools  
70 for debugging them.

71 Some works propose debugging model transformations by bringing them to a different domain  
72 such as Maude [103], DSLTrans [81] or Colored Petri Nets [111], where some specific analysis can  
73 be applied. The problem with these approaches is that the user needs to be familiar with such  
74 domains, besides, their performance and scalability can be worse than that of the original model  
75 transformation [103]. There are a few works that propose the use of contracts in order to debug  
76 model transformations [18, 22, 23]. Among them, the work by Burgueño et al. [18] is the closest to  
77 ours. They address the debugging of ATL model transformations based on contracts with a static  
78 approach that aims to identify the *guilty* rule, i.e., the faulty rule. It statically extracts the types  
79 appearing in the contracts as well as those of the MT rules and decides which rules are more likely  
80 to contain a bug. This is a *static* approach, since the transformation is not executed. Despite that, it  
81 achieves relatively good results on several case studies [18]. However, the effectiveness of *dynamic*  
82 approaches is an open question. Answering this question is one of the goals of this work.

83 Spectrum-Based Fault Localization (SBFL) is a popular technique used in software debugging  
84 for the localization of bugs [3, 116]. It uses the results of test cases and their corresponding code  
85 coverage information to estimate the likelihood of each program component (e.g., statements) of  
86 being faulty. A program spectrum details the execution information of a program from a certain  
87 perspective, such as branch or statement coverage [50]. SBFL entails identifying the part of the  
88 program whose activity correlates most with the detection of errors.

89 This paper presents and evaluates in detail the first approach that applies spectrum-based fault  
90 localization in model transformations, extending our paper with the initial ideas [100]. SBFL being  
91 a dynamic approach, our approach takes advantage of the information recovered after MT runs,  
92 what may help improve the results over static approaches [18], and at the same time complement  
93 them. We follow the approaches in [14, 18, 21, 105] and use the previously described contracts  
94 (assertions) as oracle to determine the correctness of MTs. Given a MT, a set of assertions and a  
95 set of source models, our approach indicates the violated assertions and uses the information of

96  
97 <sup>1</sup>Throughout the paper, we may also refer to *model transformation (MT) rules* as *transformation rules* or merely *rules*

99 the MT coverage to rank the transformation rules according to their suspiciousness of containing  
100 a bug. Also, out of the many existing techniques proposed in the literature for computing the  
101 suspiciousness values [86, 116, 118], we select 18 of them and compare their effectiveness in the  
102 context of MTs.

103 There is a plethora of frameworks and languages to define MTs. Among them, The ATLAS trans-  
104 formation language (ATL) [61, 85] has come to prominence in the MDE community both in the  
105 academic and the industrial arenas, so the testing of ATL transformations is of prime importance.  
106 This success is due to ATL's flexibility, support of the main metamodeling standards, usability that  
107 relies on strong tool integration within the Eclipse world, and a supportive development commu-  
108 nity [81]. In order to implement our approach and achieve automation, we have built a prototype  
109 for debugging ATL model transformations. However, we may mention that the proposed approach  
110 is applicable to any model transformation language as long as it is able to store the execution of  
111 the transformation in traces. Therefore, the approach could be trivially applied to languages such  
112 as QVT [45], Maude [26], Kermeta [56], and many more, since in most transformation languages it  
113 is possible to define the generation of an extra target model that stores the traces (cf. Section 2.2.3).

114 We have thoroughly evaluated the approach using the implemented prototype. To do so, we have  
115 selected four different case studies that differ regarding the application domains, size of metamodels  
116 and transformations, and the number and types of features of ATL used. For instance, the number of  
117 rules ranges from 8 to 39, and the lines of code from 53 to 1055. We have defined 117 OCL assertions  
118 for the four case studies, many of them taken from [18], and have applied mutation testing by  
119 creating 158 mutants using the operators presented in [99], where each mutant is a faulty variation  
120 of the original model transformation. Experimental results reveal that the best techniques place  
121 the faulty transformation rule among the three most suspicious rules in around 74% of the cases.  
122 Looking into each of the four case studies, the best techniques allow the tester to locate the fault  
123 by inspecting only 1.59, 2.99, 2.4 and 4.8 rules in each of the case studies, which are composed of 9,  
124 19, 8 and 39 rules, respectively. Furthermore, we compared our approach with a state-of-the-art  
125 approach based on the static analysis of transformation rules and assertions, observing a clear  
126 superiority of the proposed SBFL-based approach. The conclusions from our experiments serve as  
127 a proof of concept of the effectiveness of SBFL techniques to aid in the process of debugging model  
128 transformations.

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

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

141

142

143

## 143 2 BACKGROUND

144

145

146

147

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

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

## 151 2.1 Metamodeling

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

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

163 Every model must conform to a metamodel. Indeed, a metamodel defines the structure and  
164 constraints for a family of models [76]. Like everything in MDE, a metamodel is itself a model, and  
165 it is written in the language defined by its meta-metamodel. It specifies the concepts of a language,  
166 the relationships between these concepts, the structural rules that restrict the possible elements in  
167 the valid models and those combinations between elements with respect to the domain semantic  
168 rules.

169 A metamodel dictates what kind of models can be defined within a specific domain, i.e., it defines  
170 the abstract syntax of a DSL. The concrete syntax of DSLs can be defined in several ways, normally  
171 either graphically or textually. In order to provide a DSL with semantics and behavior, its defining  
172 metamodel may not be enough. Therefore, apart from its concrete and abstract syntaxes, also its  
173 semantics may need to be defined. For instance, model transformations can be used in order to  
174 give semantics to a DSL by translating it to a different domain where further analysis, simulations,  
175 and so on can be performed [104]. This mechanism enables the definition of flexible and reusable  
176 DSLs, where several kinds of analysis can be defined [32, 77].

## 177 2.2 Model Transformations

178 Model transformations play a cornerstone role in Model-Driven Engineering (MDE) since they  
179 provide the essential mechanisms for manipulating and transforming models [16, 97]. They allow  
180 querying, synthesizing and transforming models into other models or into code, so they are essential  
181 for building systems in MDE. A model transformation is a program executed by a transformation  
182 engine that takes one or more input models and produces one or more output models, as illustrated  
183 by the model transformation pattern [28] in Figure 1<sup>2</sup>. Model transformations are developed on  
184 the metamodel level, so they are reusable for all valid model instances. Most MT languages are  
185 composed of model transformation rules, where each rule deals with the construction of part of  
186 the target model. They match input elements from the source model and generate output elements  
187 that compose the target model.

188 There is a plethora of frameworks and languages to define MTs, such as Henshin [10], AGG [98],  
189 Maude [26], ATOM<sup>3</sup> [30], e-Motions [87], VIATRA [27], MOMoT [35–37], QVT [45], Kermeta [56],  
190 JTL [24], and ATL [62]. In most of these frameworks and languages, model transformations are

191 <sup>2</sup>In the paper, we use the terms *input/output* models/metamodels and *source/target* models/metamodels indistinctly.

197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245

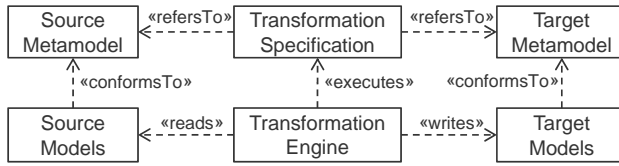


Fig. 1. Model transformation pattern (from [28])

composed of transformation rules. Among them, we focus in this paper on the ATL language due to its importance in both the academic and the industrial arenas.

**2.2.1 *ATL as 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, 85].

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<sup>3</sup>. 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 paper as running example. It transforms a BibTeXXML model to a DocBook composed document. BibTeXXML<sup>4</sup> is an XML-based format for the BibTeX bibliographic tool. DocBook [108] 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 (Fig. 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

<sup>3</sup>[https://wiki.eclipse.org/ATL/User\\_Guide\\_-\\_The\\_ATL\\_Language](https://wiki.eclipse.org/ATL/User_Guide_-_The_ATL_Language)

<sup>4</sup><http://bibtexml.sourceforge.net/>

246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294

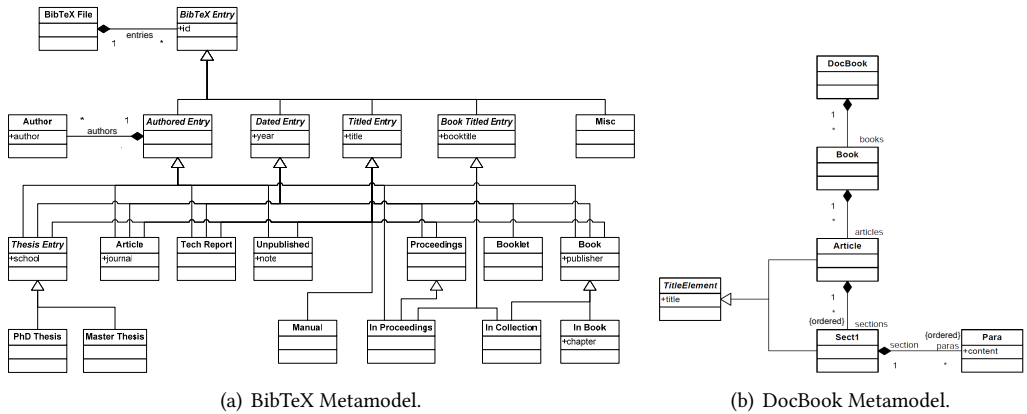


Fig. 2. Metamodels of the BibTeX2DocBook transformation (from [54])

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*. 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 (Fig. 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 9 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.

Listing 1. *BibTeX2DocBook* MT.

```

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)

```



```
295 12     endif)->sortedBy(e | e.title);
296 13
297 14 helper def: articleSet : Sequence(BibTeX!Article) =
298 15 BibTeX!Article.allInstances()->iterate(e; ret : Sequence(BibTeX!Article) = Sequence {} |
299 16   if ret->collect(e | e.journal)->includes(e.journal) then ret else ret->including(e)
300 17   endif)->sortedBy(e | e.journal);
301 18
302 19 helper context BibTeX!BibTeXEntry def: buildEntryPara() : String =
303 20 '[' + self.id + ']'
304 21 + '_' + self.oclcType().name
305 22 + (if self.oclcIsKindOf(BibTeX!TitledEntry) then '_' + self.title else '' endif)
306 23 + (if self.oclcIsKindOf(BibTeX!AuthoredEntry)
307 24   then self.authors->iterate(e; str : String = '' | str + '_' + e.author) else '' endif)
308 25 + (if self.oclcIsKindOf(BibTeX!DatedEntry) then '_' + self.year else '' endif)
309 26 + (if self.oclcIsKindOf(BibTeX!BookTitledEntry) then '_' + self.booktitle else '' endif)
310 27 + (if self.oclcIsKindOf(BibTeX!ThesisEntry) then '_' + self.school else '' endif)
311 28 + (if self.oclcIsKindOf(BibTeX!Article) then '_' + self.journal else '' endif)
312 29 + (if self.oclcIsKindOf(BibTeX!Unpublished) then '_' + self.note else '' endif)
313 30 + (if self.oclcIsKindOf(BibTeX!Book) then '_' + self.publisher else '' endif)
314 31 + (if self.oclcIsKindOf(BibTeX!InBook) then '_' + self.chapter.toString() else '' endif);
315 32
316 33
317 34 rule Main {                               -- tr1
318 35   from
319 36     bib : BibTeX!BibTeXFile
320 37   to
321 38     doc : DocBook!DocBook (books <- boo),
322 39     boo : DocBook!Book (articles <- art),
323 40     articles : DocBook!Article (title <- 'BibTeXXML_to_DocBook',
324 41       sections <- Sequence{se1, se2, se3, se4}),
325 42     se1 : DocBook!Sect1 (title <- 'References_List',
326 43       paras <- BibTeX!BibTeXEntry.allInstances()->sortedBy(e | e.id)),
327 44     se2 : DocBook!Sect1 (title <- 'Authors_List',
328 45       paras <- thisModule.authorSet),
329 46     se3 : DocBook!Sect1 (title <- 'Titles_List',
330 47       paras <- thisModule.titledEntrySet->collect(e | thisModule.resolveTemp(e, '
331 48         title_para'))),
332 49     se4 : DocBook!Sect1 (title <- 'Journals_List',
333 50       paras <- thisModule.articleSet->collect(e | thisModule.resolveTemp(e, '
334 51         journal_para'))
335 52   }
336 53 rule Author {                               -- tr2
337 54   from
338 55     a : BibTeX!Author (thisModule.authorSet->includes(a))
339 56   to
340 57     p1 : DocBook!Para (content <- a.author)
341 58   }
342 59 rule UntitledEntry {                       -- tr3
343 60   from
344 61     e : BibTeX!BibTeXEntry (not e.oclcIsKindOf(BibTeX!TitledEntry))
345 62   to
346 63     p : DocBook!Para (content <- e.buildEntryPara())
347 64   }
348 65
349 66 rule TitledEntry_Title_NoArticle {        -- tr4
350 67   from
351 68     e : BibTeX!TitledEntry (thisModule.titledEntrySet->includes(e) and
352 69       not e.oclcIsKindOf(BibTeX!Article))
353 70   to
354 71     entry_para : DocBook!Para (content <- e.buildEntryPara()),
355 72     title_para : DocBook!Para (content <- e.title)
356 73   }
357 74
358 75 rule TitledEntry_NoTitle_NoArticle {     -- tr5
359 76   from
360 77     e : BibTeX!TitledEntry (not thisModule.titledEntrySet->includes(e) and
361 78       not e.oclcIsKindOf(BibTeX!Article))
362 79   to
```

```
344 80     entry_para : DocBook!Para (content <- e.buildEntryPara())
345 81 }
346 82
347 83 rule Article_Title_Journal {           -- tr6
348 84   from
349 85     e : BibTeX!Article (thisModule.titledEntrySet->includes(e) and
350 86                       thisModule.articleSet->includes(e))
351 87   to
352 88     entry_para : DocBook!Para (content <- e.buildEntryPara()),
353 89     title_para : DocBook!Para (content <- e.title),
354 90     journal_para : DocBook!Para (content <- e.journal)
355 91 }
356 92
357 93 rule Article_NoTitle_Journal {         -- tr7
358 94   from
359 95     e : BibTeX!Article (not thisModule.titledEntrySet->includes(e) and
360 96                       thisModule.articleSet->includes(e))
361 97   to
362 98     entry_para : DocBook!Para (content <- e.buildEntryPara()),
363 99     journal_para : DocBook!Para (content <- e.journal)
364 100 }
365 101
366 102 rule Article_Title_NoJournal {        -- tr8
367 103   from
368 104     e : BibTeX!Article (thisModule.titledEntrySet->includes(e) and
369 105                       not thisModule.articleSet->includes(e))
370 106   to
371 107     entry_para : DocBook!Para (content <- e.buildEntryPara()),
372 108     title_para : DocBook!Para (content <- e.title)
373 109 }
374 110
375 111 rule Article_NoTitle_NoJournal {      -- tr9
376 112   from
377 113     e : BibTeX!Article (not thisModule.titledEntrySet->includes(e) and
378 114                       not thisModule.articleSet->includes(e))
379 115   to
380 116     entry_para : DocBook!Para (content <- e.buildEntryPara())
381 117 }
```

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* in order 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].

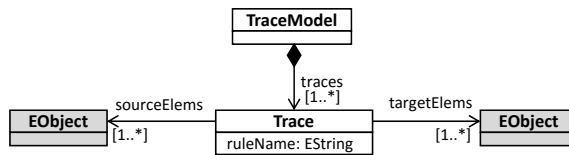
**2.2.3 ATL Internal Traces Mechanism.** The ATL engine works in two steps. First, all elements are created. Second, their features are initialized. This second phase implies to resolve the corresponding references. For instance, in the transformation shown in Listing 1, line 45 initializes the *paras* 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.



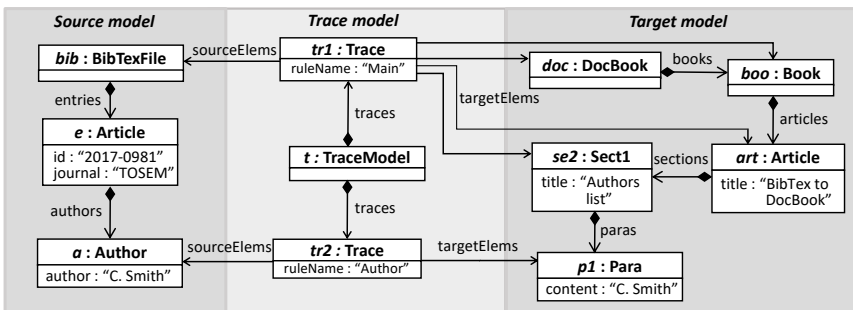
393 In order to resolve these references, ATL uses an internal tracing mechanism. Thereby, every time  
 394 a rule is executed, it creates a new trace and stores it in the internal trace model. A trace model can  
 395 be automatically obtained from a transformation execution, e.g., by using Jouault's *TraceAdder* [60],  
 396 and is composed of a set of traces, one for each rule execution. In our approach, we obtain trace  
 397 models that conform to the metamodel displayed in Figure 3(a). A trace captures the name of the  
 398 applied rule and the elements instantiating classes of the source metamodel (*sourceElems* reference)  
 399 that are used to create new elements that instantiate classes in the target metamodel (*targetElems*  
 400 reference). The elements pointed by such references are represented with *EObject* because, when the  
 401 metamodel is instantiated in a specific trace model, they can be any element of the source and target  
 402 models, respectively. The execution of both imperative *-(unique) lazy* and *called-* and declarative  
 403 *-matched-* rules are stored in the traces. This means that we have three models (the source model,  
 404 the target model and the trace model) linked by several so-called inter-model references. Therefore,  
 405 by navigating the traces, the ATL engine obtains information of which target element(s) have been  
 406 created from which source element(s) and by which rule.

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

415 The interesting aspect in this figure is the *paras* reference between *se2* and *p1* in the target  
 416 model, created using the traces. The process how ATL resolves such association is the following.  
 417 As mentioned before, after creating all target elements in the first phase, it resolves the references  
 418



(a) Trace Metamodel.



(b) Trace Model Sample.

Fig. 3. Traces in Model Transformations

442 in the second phase. In our example, it means resolving the binding in line 45, where helper  
443 *authorSet* returns all *Authors* in the source model. Therefore, such binding is expressing that the  
444 *paras* reference from element *se2* in the target model should point all elements of type *Author* in the  
445 source model, *a* in our case. Of course, target elements cannot point source elements, so the ATL  
446 engine searches in the traces in order to recover the target elements created from *Author* elements.  
447 In our example, it recovers element *p1*, created from *Author a*, by inspecting trace *tr2*.

448 When ATL resolves the references, it recovers the first target element created by the correspond-  
449 ing rule, i.e., the first one that is specified in the rule (right after the *to* part of the rule). For instance,  
450 element *entry\_para* (line 71 in Listing 1) would be the one recovered when resolving a binding  
451 reference with rule *TitledEntry\_Title\_NoArticle*. In order to specify a different element to recover  
452 when resolving the references, ATL provides the *resolveTemp* function. For instance, the *resolveTemp*  
453 function used in line 47 makes the engine retrieve the target element identified with the string  
454 “*title\_para*”, i.e., the one created in line 72.

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

458 Please note that the availability of such a simple trace model is useful in many different model  
459 transformations languages, what also increases the applicability of our approach beyond ATL. For  
460 instance, Falleri et al. [33] propose a simple trace metamodel for Kermeta model transformations,  
461 and Anastasakis et al. [8] simply require the link between source and target models in an Alloy  
462 transformation. Rose et al. [88] mention the Fujaba and the MOLA (graphical transformation  
463 language developed at the University of Latvia) traceability associations, similar to the one we  
464 use, and Troya and Vallecillo [103] apply the same trace metamodel in order to represent model  
465 transformations in Maude. Due to the importance of trace models, Jiménez et al. [73] propose  
466 a toolkit that allows not only the definition of model transformations but also supports trace  
467 generation.

### 469 2.3 Spectrum-Based Fault Localization

470 *Spectrum-Based Fault Localization (SBFL)* uses the results of test cases and their corresponding code  
471 coverage information to estimate the likelihood of each program component (e.g., statements) of  
472 being faulty. A program spectrum details the execution information of a program from a certain  
473 perspective, such as branch or statement coverage [50]. Table 1 depicts an example showing how  
474 the technique is applied to a sample program [116]. This program receives a natural number, *a*.  
475 If it is bigger than 1, the program must print the result of adding 1 to such number as well as its  
476 double. Otherwise, it must print the number minus 1 as well as the number itself.

477 Having a look at the table, it horizontally shows the code statements of the program, i.e., its  
478 components. Note that a bug is seeded in statement *s<sub>7</sub>*, so that it does not multiply the number by  
479 2. Also note that SBFL considers all lines as statements, so, for instance, the line containing only  
480 the character that closes a branch, ‘}’, conforms statement *s<sub>11</sub>*. However, statement *s<sub>5</sub>* includes a  
481 condition as well as the opening of a branch with character ‘{’. Therefore, the way of writing a  
482 program may have an impact in the results returned by SBFL techniques. Vertically, the table shows  
483 three test cases of the program. For each test case (i.e., *tc<sub>1</sub>*, *tc<sub>2</sub>*, and *tc<sub>3</sub>*), a cell is marked with “●” if  
484 the program statement of the row has been exercised by the test case of the column, creating what  
485 is known as *coverage matrix* [4]. Additionally, the final row depicts the outcome of each test case,  
486 either “Successful” or “Failed”, conforming the so-called *error vector* [4]. Based on this information,  
487 it is possible to identify which components were involved in a failure (and which ones were not),  
488 narrowing the search for the faulty component that made the execution fail.

490

Table 1. An example showing the suspiciousness value computed using Tarantula (taken from [116])

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; // <b>BUG: 2i → i</b>			•	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				

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 [118, 121], which is the percentage of statements in a program that has to be examined until the first faulty statement is reached, i.e.,

$$EXAM_{Score} = \frac{\text{Number of statements examined}}{\text{Total number of statements}}$$

It is noteworthy that suspiciousness techniques may often provide the same value for different statements, being these tied for the same position in the ranking, e.g., statements  $s_6$  and  $s_7$  in Table 1. In order 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 [116]. 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.

Table 2. Tarantula [59] suspiciousness values for the simplified *BibTeX2DocBook* MT when  $OCL_2$  fails

T. Rule	$tc_{02}$	$tc_{12}$	$tc_{22}$	$tc_{32}$	$tc_{42}$	$tc_{52}$	$tc_{62}$	$tc_{72}$	$tc_{82}$	$tc_{92}$	$N_{CF}$	$N_{UF}$	$N_{CS}$	$N_{US}$	$N_C$	$N_U$	SuspRank	
$tr_1$	•	•	•	•	•	•	•	•	•	•	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	F								

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 in order 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, it has the EXAM score  $\frac{13}{13} = 1$ . Therefore, the set of values for the EXAM score, from best to worse, is  $\{\frac{1}{num\_statements}, \frac{2}{num\_statements}, \dots, \frac{num\_statements}{num\_statements}\}$ .

### 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 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 by 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

589 is not satisfied for a MT when there is, at least, one test case where it is violated. This means, for  
590 instance, that for  $ocl_1$  to be satisfied, it must be satisfied in  $\{tc_{11}, tc_{21}, \dots, tc_{n1}\}$ .

591 We may recall that this paper focuses on debugging and not testing. Thus, we do not impose any  
592 constraint on how the source models are generated, either manually or automatically, neither on  
593 the type of OCL assertions used. However, please note that the effectiveness of SBFL is directly  
594 related to the design of the test cases. For instance, having only one test case is useless, since no  
595 coverage matrix can be created. Besides, the source models should have a good coverage of the MT  
596 and, at the same time, be diverse. This way, different source models will likely fire different parts  
597 of the model transformation, and together they will exercise all rules and will produce a useful  
598 spectra.

599 Table 2 depicts a sample coverage matrix constructed using our approach. Horizontally, the  
600 table shows the transformation rules in which we aim to locate bugs. In particular, we list the 9  
601 transformation rules  $\langle tr_1, tr_2, \dots, tr_9 \rangle$  of the *BibTeX2DocBook* example [54], where a bug has been  
602 seeded in  $tr_2$ . Vertically, the table shows 10 test cases aiming to check the correctness of constraint  
603  $OCL_2$  in Listing 2,  $\langle tc_{02}, tc_{12}, \dots, tc_{92} \rangle$ . A cell is marked with “•” if the transformation rule of the  
604 row has been exercised by the test case of the column. The information about the rules triggered by  
605 a given test case can be collected by inspecting the trace model, e.g., using Jouault’s *TraceAdder* [60]  
606 (cf. Section 2.2.3). The final row depicts the error vector with the outcome of each test case, either  
607 successful (“S”) or failed (“F”). In the example, all test cases fail except  $tc_{12}$ , i.e.,  $OCL_2$  is violated.

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

614  
615  
616  
617 Listing 2. Sample OCL assertions for the *BibTeX2DocBook* MT.

```
618 1 --OCL1. The main Article must be properly created and named  
619 2 TrgBook.allInstances()->forAll(b|b.articles->exists(a|a.title='BibTeXML_to_DocBook'))  
620 3 --OCL2. For each author, there must be a paragraph with the name of the author within a  
621 4 SrcAuthor.allInstances()->forAll(a|TrgPara.allInstances()->exists(p|p.content=a.author and  
622 5 --OCL3. The titles of all publications must appear in the content of a paragraph of a  
623 6 SrcTitledEntry.allInstances()->forAll(te|TrgSect1.allInstances()->exists(s|s.paras->exists  
624 7 --OCL4. There must be the same number of BibTeXFile and DocBook instances  
625 8 SrcBibTeXFile.allInstances()->size()=TrgDocBook.allInstances()->size()
```

626

627

628

### 629 3.2 Calculating Suspiciousness

630 The following notation will be used throughout the paper and is used in our implementation to  
631 compute the suspiciousness of transformation rules from the information collected in the coverage  
632 matrix and the error vector. This notation is directly translated from the context of SBFL in software  
633 programs [116] by using transformation rules as the components (e.g., instead of statements),  
634 namely:  
635

636

637

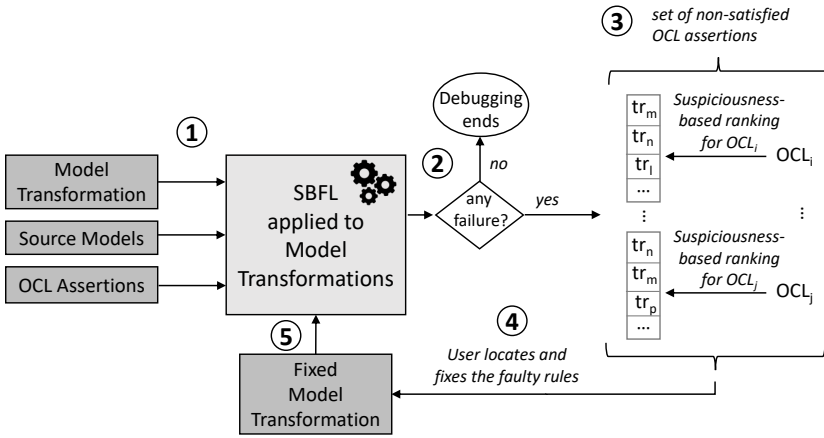


Fig. 4. Debugging of a MT applying our SBFL approach

- 655  $N_{CF}$  number of failed test cases that cover a rule
- 656  $N_{UF}$  number of failed test cases that do not cover a rule
- 657  $N_{CS}$  number of successful test cases that cover a rule
- 658  $N_{US}$  number of successful test cases that do not cover a rule
- 659  $N_C$  total number of test cases that cover a rule
- 660  $N_U$  total number of test cases that do not cover a rule
- 661  $N_S$  total number of successful test cases
- 662  $N_F$  total number of failed test cases

663 Table 2 shows the values of  $N_{CF}$ ,  $N_{UF}$ ,  $N_{CS}$ ,  $N_{US}$ ,  $N_C$  and  $N_U$  for each transformation rule. The  
 664 values of  $N_S$  and  $N_F$  are the same for all the rows, 9 and 1 respectively, and are omitted. Based on  
 665 this information, the suspiciousness of each transformation rule using *Tarantula* is depicted in the  
 666 column "Susp", followed by the position of each rule in the suspiciousness-based ranking. In the  
 667 example, the faulty rule  $tr_2$  is ranked first, tied with  $tr_7$ . Assuming that the faulty rule was inspected  
 668 in the second place (worst-case scenario), the EXAM score would be calculated as  $\frac{2}{9} = 0.222\%$ , i.e.,  
 669 22.2% of the transformation rules need to be examined in order to locate the bug.

### 671 3.3 Methodology

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

- 675 (1) The inputs have to be provided, namely the *Model Transformation* under test as well as the  
 676 sets of *Source Models* and *OCL Assertions*.
- 677 (2) The approach executes and indicates whether there is *any failure*, ending the process if there  
 678 is none.
- 679 (3) If there is a failure, it indicates the *set of non-satisfied OCL assertions*, i.e., those that are  
 680 violated for at least one test case. As explained in Section 3.1, it constructs a coverage matrix  
 681 and an error vector for each non-satisfied assertion and returns the *suspiciousness-based*  
 682 *rankings* in each case.
- 683 (4) The user picks the ranking of one of the OCL assertions in order to *locate and fix the faulty*  
 684 *rule* that made the assertion fail. As described in Section 4.2.5, we study the effectiveness of 18  
 685



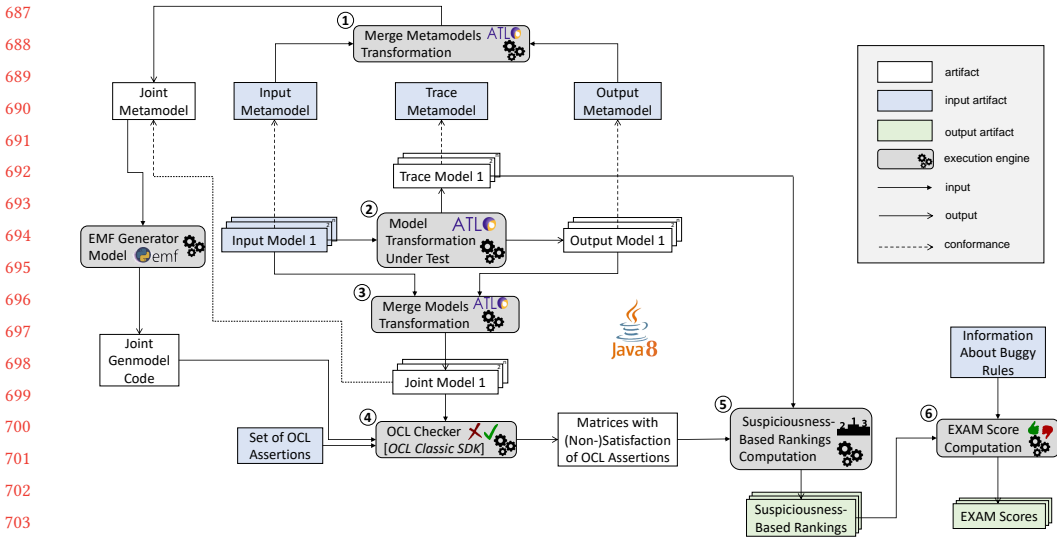


Fig. 5. Implementation and Automation of our Approach

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*, less 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:

- (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*<sup>5</sup>, 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

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

736 defined on both metamodels, namely the input and output metamodels. For this reason, we  
737 need to merge both metamodels into one, and same thing for the <input, output> model pairs.  
738 Therefore, the first step in our approach takes the *Input Metamodel* and *Output Metamodel* as  
739 input and, with the *Merge Metamodels Transformation*, it creates a *Joint Metamodel*. Due to  
740 the possibility of having classes with the same name in the input and output metamodels,  
741 this transformation puts the prefixes “Src” and “Trg” in all classes of the input and output  
742 metamodels, respectively.

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

745 (2) The next step in our approach is to run the *ATL Transformation Under Test* with all the *Input*  
746 *Models* in order to generate the corresponding *Output Models*. Our Java program orchestrates  
747 all these model transformation executions.

748 (3) For the same reason as explained in the first step, we need to merge the input and output  
749 models into the so-called *Joint Models*. These models must conform to the *Joint Metamodel*  
750 obtained in the first step. The *Merge Models Transformation* generates all the *Joint Models* for  
751 all the <*InputModel<sub>i</sub>*, *OutputModel<sub>i</sub>*> pairs.

752 (4) The next step is to check the *Set of OCL Assertions*. This must be done for all the *Joint Models*  
753 constructed after the executions of the model transformation. As explained in the first step,  
754 we need for this the Java code obtained from the *Joint Metamodel*. This step produces as  
755 output information about the satisfiability of the OCL assertions, captured in the figure  
756 as *Matrices with (Non-)Satisfaction of OCL Assertions*. This is different from the coverage  
757 matrices explained before, since the purpose now is to identify those OCL assertions that fail  
758 for at least one test case, so that coverage matrices and error vectors will be then computed  
759 for such assertions. This matrix, used internally by the program, has the OCL assertions  
760 as rows and the joint models as columns. Cell <*i*, *j*> is assigned 1 if the *i*-th OCL assertion  
761 is not satisfied when executing the model transformation with the *j*-th input model, and 0  
762 otherwise. Therefore, an OCL assertion has failed when there is at least a 1 in its row.

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

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

776 With the information of the coverage matrices and error vectors, we are able to automatically  
777 compute the 8 values described in Section 3.2 for computing the suspiciousness, namely  $N_{CF}$ ,  
778  $N_{UF}$ ,  $N_{CS}$ ,  $N_{US}$ ,  $N_C$ ,  $N_U$ ,  $N_S$ ,  $N_F$ . Finally, with these values and the formulae for calculating  
779 the suspiciousness with the 18 techniques considered in this study (cf. Section 4.2.5 and  
780 Table 6), we obtain the *Suspiciousness-Based Rankings*. These rankings, together with the  
781 coverage matrices, error vectors and values are returned as comma-separated values (CSV)

782

783

784

785 files by our tool. In particular, it returns a detailed CSV file for each non-satisfied OCL  
786 assertion.

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

## 795 4 EVALUATION

### 796 4.1 Research Questions

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

- 798 • **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  
799 for applying spectrum-based techniques for locating faults in model transformations, we  
800 want to answer whether this is feasible. This means, we want to check whether it is possible  
801 to automatically obtain for a model transformation a suspiciousness-based ranking, according  
802 to SBFL techniques, that indicates which rules should be inspected first in case of failure.
- 803 • **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  
804 been proposed in the literature in different fields, we want to determine how they behave,  
805 comparing among each other, in the context of model transformations. This means we want  
806 to study which techniques provide the best suspiciousness-based rankings and which ones  
807 provide the worst rankings.
- 808 • **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  
809 it is possible to state that applying spectrum-based techniques can accurately help the  
810 developer in the debugging of model transformations. This will be answered affirmatively if  
811 the techniques that are more effective, according to the answer to the previous RQ, provide  
812 accurate suspiciousness-based rankings.
- 813 • **RQ4 - Dynamic vs Static.** *How does our approach behave in comparison with a static approach?*  
814 Being our approach dynamic, we want to compare its performance with a notable approach  
815 for locating bugs in model transformations applying a static approach [18].

### 816 4.2 Experimental Setup

817  
818 4.2.1 *Case Studies.* We have used four case studies in order to evaluate our approach and  
819 developed solution. Two of them have been taken from the open-source ATL Zoo repository [12]  
820 and the two others from research projects and tools. They all differ regarding the application  
821 domains, size of metamodels and the number and types of ATL features used. Table 3 summarizes  
822 some information regarding the transformations. For instance, the size of the metamodels vary  
823 from 4 to 33 classes in the input metamodels and from 8 to 77 classes in the output metamodels.  
824 Regarding the size of the transformations, the number of rules range from 8 to 39 (in the table, *M*  
825 stands for *matched rules*, *(U)L* for *(unique) lazy rules* and *C* for *called rules*), and the lines of code  
826 (LoC) from 53 to 1055. This means that the smaller transformation is approximately 20 times smaller,  
827 in terms of LoC, than the biggest one. Furthermore, the transformations contain from no helper to  
828  
829  
830  
831  
832  
833

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	×	✓	✓	✓	✓

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 in order 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 [81], 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 [101] 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 [112], and we have considered the version proposed in [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  $R_8, R_7 < R_6; R_6, R_5 < R_4; R_4, R_3, R_2 < R_1$ . The presence of inheritance may worsen the results of SBFL techniques. Imagine we have, for instance,  $R_3 < R_2 < R_1$  in a model transformation and rule  $R_3$  is executed. In the trace, it is stored not only the execution of  $R_3$ , but also the execution of  $R_2$  and  $R_1$ , 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, the suspiciousness rankings will not make any difference between the three rules, having the three of them the same suspiciousness value.
- **BibTeX2DocBook.** This case study is the one used as running example in our paper. It is shown in Listing 1, and a complete description is available on [54].
- **CPL2SPL.** This transformation, described in [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* [87]. It converts models conforming to the Ecore metamodel into models that conform to the Maude [26] metamodel, in order to apply some formal reasoning on them afterwards.

4.2.2 *Test Suite.* Since this is a dynamic approach, we need input models in order to trigger the model transformations. For evaluating our work, we have developed a light-weight 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.

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	1400	18	90
BibTeX2DocBook	100	27 / 16	2700	40	269
CPL2SPL	100	34 / 15	3400	50	150
Ecore2Maude	100	42 / 3	4200	50	155
Total	400	117 / 44	11700	158	664

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<sup>6</sup>. 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 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 1400, 2700, 3400 and 4200 test cases in each of the transformations, having a total of 11700 test cases.

**4.2.3 Mutants.** In order 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 [99] and have applied them in the different case studies. The aim of these operators is the same as the ones presented by Mottu et al. [78], 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. [78] present several operators related to model *navigation*, such as *ROCC: relation to another class change*, which in [99] 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, in order 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. In order to be able to have such restricting mutants and as many other approaches do [9, 46, 78], we have generated them manually using the operators proposed in [99].

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

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



932 *binding deletion* is used for deleting the only binding of rule 2. Our set of OCL assertions in the  
933 different case studies is complete enough as to kill all 158 mutants, i.e., all mutants make at least  
934 one OCL assertion fail, indicating there is an error in the MT. Table 4 displays the number of  
935 mutants that have been created for each case study (column 5), while Table 5 presents the mutation  
936 operators [99] that have been used for creating the mutants, and the number of mutants in which  
937 the operators are applied. Please note that fewer mutants have been created for the *UML2ER* case  
938 study. This is due to the fact that this model transformation is the smallest one, and it is actually  
939 almost four times smaller than the second smaller one in terms of lines of code (cf. Table 3).

940 We may mention that more than one mutation operators can be used for constructing one mutant.  
941 For instance, we can combine *out-pattern element addition* with *binding addition* in order to create  
942 a mutant that adds one more target element and updates one of its features. We have also created  
943 mutants with more than one faulty rule. The reason for including higher-order mutants [57, 82] is  
944 the definition of realistic mutants, i.e., mutants that produce valid models and reproduce typical  
945 mistakes caused by developers. In fact, as presented in the survey on software fault localization by  
946 Wong et al. [116], having programs with a single bug (i.e., each faulty program has exactly one  
947 bug) is not the case for real-life software, which in general contains multiple bugs. Results of a  
948 study [49] based on an analysis of fault and failure data from two large, real-world projects show  
949 that individual failures are often triggered by multiple bugs spread throughout the system. Another  
950 study [69] also reports a similar finding. The very same reality occurs in model transformations,  
951 where it is not common to have isolated faults located in only one rule. Indeed, since some rules  
952 have implicit relations among them (cf. Section 2.2.3), it is very common to have errors spread in  
953 several rules.

954  
955 Listing 3. Excerpt of a mutant of *BibTeX2DocBook* MT.

```
956 1 rule Author {                               -- tr2  
957 2   from  
958 3   a : BibTeX!Author (thisModule.authorSet->includes(a))  
959 4   to  
960 5   p1 : DocBook!Para () --binding deletion  
961 6 }
```

962 **4.2.4 Set of Non-Satisfied OCL Assertions.** As described, we have produced 158 mutants that  
963 correspond to buggy versions of the model transformations in the different case studies. Each one  
964 of them may violate one or more of the OCL assertions defined for the corresponding case study  
965 (of course, more than one mutant may violate the same assertion). In total, the 158 mutants make  
966 664 OCL assertions fail, as displayed in the last column of Table 4, so the results of our evaluation  
967 are extracted from the 664 suspiciousness-based rankings obtained, one for each violated assertion.  
968 These rankings are the results of suspiciousness values calculated with 664 coverage matrices and  
969 with the corresponding 664 error vectors. These coverage matrices have different sizes depending  
970 on the case study. All of them have 100 columns, since we are using 100 input models, and the  
971 number of rows is determined by the number of rules in the model transformation.

972 **4.2.5 Techniques for Suspiciousness Computation.** We are interested in studying how different  
973 techniques<sup>7</sup> for computing the suspiciousness of program components behave in the context of  
974 model transformations. To this end, we have surveyed papers that apply spectrum-based fault  
975 localization techniques in different contexts and have selected the 18 techniques that, together with  
976 their corresponding formulae, are displayed in Table 6. *Tarantula* [59] is one of the best-known fault  
977 localization techniques. It follows the intuition that statements that are executed primarily by more

978  
979 <sup>7</sup>Throughout the evaluation, we use the terms *techniques* and *metrics* indistinctly



981 failed test cases are highly likely to be faulty. Additionally, statements that are executed primarily  
982 by more successful test cases are less likely to be faulty. The *Ochiai* similarity coefficient is known  
983 from the biology domain and it has been proved to outperform several other coefficients used in  
984 fault localization and data clustering [3]. This can be attributed to the *Ochiai* coefficient being  
985 more sensitive to activity of potential fault locations in failed runs than to activity in passed runs.  
986 *Ochiai2* is an extension of such technique [11, 79, 116]. *Kulczynski2*, taken from the field of artificial  
987 intelligence, and *Cohen* have showed promising results in preliminary experiments [79, 118]. *Russel-*  
988 *Rao* has shown different results in previous experiments [86, 117, 118], while *Simple Matching* has  
989 been used in clustering [79]. *Reogers & Tanimoto* presented a high similarity with *Simple Matching*  
990 when ranking in the study performed in [79]. The framework called *Barinel* [70] combines spectrum-  
991 based fault localization and model-based debugging to localize single and multiple bugs in programs.  
992 *Zoltar* [55] is also a tool set for fault localization. *Arithmetic Mean*, *Phi* (Geometric Mean), *Op2*  
993 and *Pierce* have been considered in some comparative studies with other metrics [79, 116, 118].  
994 *Mountford* behaves as the second best technique, among 17 of them, for a specific program in a  
995 study performed in [115], where *Baroni-Urbani & Buser* is also studied. As for  $D^*$ , its numerator,  
996  $(N_{CF})^*$ , depends on the value of ‘\*’ selected. This technique resulted the best technique in the study  
997 performed in [114], where ‘\*’ was assigned a value of 2. We have followed the same approach, so  
998 we have  $(N_{CF})^2$  in the numerator of the formula.

999 Note that the computation of these formulae may result in having zero in a denominator. Different  
1000 approaches mention how to deal with such cases [80, 119, 120]. Following the guidelines of these  
1001 works, if a denominator is zero and the numerator is also zero, our computation returns zero.  
1002 However, if the numerator is not 0, then it returns 1 [120].

1003  
1004 **4.2.6 Evaluation Metric.** In order to compare the effectiveness of the different SBFL techniques,  
1005 we apply the EXAM score described in Section 2.3. In the context of this work, this score indicates  
1006 the percentage of transformation rules that need to be examined until the faulty rule is reached. Its  
1007 value is in the range  $[1/(\text{num rules}), 1]$ , and the higher its value, the worse.

1008  
1009  
1010 Table 5. Mutation operators used and number of mutants where they are applied

Mutation Operator (from [99])	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

Table 6. Techniques applied for suspiciousness computation

Technique	Formula
Arithmetic Mean [118]	$\frac{2(N_{CF} \times N_{US} - N_{UF} \times N_{CS})}{(N_{CF} + N_{CS}) \times (N_{US} + N_{UF}) + (N_{CF} + N_{UF}) \times (N_{CS} + N_{US})}$
Barinel [2]	$1 - \frac{N_{CS}}{N_{CS} + N_{CF}}$
Baroni-Urbani & Buser [115]	$\frac{\sqrt{N_{CF} \times N_{US} + N_{CF}}}{\sqrt{N_{CF} \times N_{US} + N_{CF} + N_{CS} + N_{UF}}}$
Braun-Banquet [116]	$\frac{N_{CF}}{\max(N_{CF} + N_{CS}, N_{CF} + N_{UF})}$
Cohen [79]	$\frac{2 \times (N_{CF} \times N_{US} - N_{UF} \times N_{CS})}{(N_{CF} + N_{CS}) \times (N_{US} + N_{CS}) + (N_{CF} + N_{UF}) \times (N_{UF} + N_{US})}$
D* [114]	$\frac{(N_{CF})^*}{N_{CS} + N_F + N_{CF}}$
Kulczynski2 [79]	$\frac{1}{2} \times \left( \frac{N_{CF}}{N_{CF} + N_{UF}} + \frac{N_{CF}}{N_{CF} + N_{CS}} \right)$
Mountford [115]	$\frac{N_{CF}}{0.5 \times ((N_{CF} \times N_{CS}) + (N_{CF} \times N_{UF})) + (N_{CS} \times N_{UF})}$
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 [79]	$N_{CF} - \frac{N_{CS}}{N_S + 1}$
Phi [75]	$\frac{N_{CF} \times N_{US} - N_{UF} \times N_{CS}}{\sqrt{(N_{CF} + N_{CS}) \times (N_{CF} + N_{UF}) \times (N_{CS} + N_{US}) \times (N_{UF} + N_{US})}}$
Pierce [116]	$\frac{(N_{CF} \times N_{UF}) + (N_{UF} \times N_{CS})}{(N_{CF} \times N_{UF}) + (2 \times N_{UF} \times N_{US}) + (N_{CS} \times N_{US})}$
Rogers & Tanimoto [74]	$\frac{N_{CF} + N_{US}}{N_{CF} + N_{US} + 2(N_{UF} + N_{CS})}$
Russel-Rao [86]	$\frac{N_{CF}}{N_{CF} + N_{UF} + N_{CS} + N_{US}}$
Simple Matching [116]	$\frac{N_{CF} + N_{US}}{N_{CF} + N_{CS} + N_{US} + N_{UF}}$
Tarantula [59]	$\frac{\frac{N_{CF}}{N_F}}{\frac{N_{CF}}{N_F} + \frac{N_{CS}}{N_S}}$
Zoltar [55]	$\frac{N_{CF}}{N_{CF} + N_{UF} + N_{CS} + \frac{10000 \times N_{UF} \times N_{CS}}{N_{CF}}}$

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, 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, 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, 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

1079 the middle place of a tie. Therefore, if it is in a tie with  $(n - 1)$  other rules, it will be inspected in  
1080 the  $(n/2)$ th position.

1081 4.2.7 *Execution Environment.* All the runs have been executed in a PC running the 64-bits OS  
1082 Windows 10 Pro with processor Intel Core i7-4770 @ 3.40GHz and 16 GB of RAM. We have used  
1083 Eclipse Modeling Tools version Mars Release 2 (4.5.2), and we had to install the plugins ATL (we  
1084 have used version 3.6.0) and ATL/EMFTVM (version 3.8.0). Finally, Java 8 is required.

### 1086 4.3 Experimental Results

1087 Table 7 shows the descriptive statistics of the EXAM score for each suspiciousness computation  
1088 technique when applied to each case study on the three evaluation scenarios (average-, best- and  
1089 worst-case scenarios)—ignore for now the rows with the numbers for *Burgueño'15*, as those numbers  
1090 are commented in Section 4.5.3. We may recall that the EXAM score, in the range  $(0, 1]$ , indicates the  
1091 percentage of transformation rules that need to be inspected in order to locate the faulty rule. This  
1092 score is never 0, since the inspection of the faulty rule counts. For this reason, since the MTs under  
1093 test for each case study contain a different number of rules (9 in *Bibtex2DocBook*, 19 in *CPL2SPL*,  
1094 39 in *Ecore2Maude*, and 8 in *URML2ER*), the best possible values (the case where the faulty rule is  
1095 ranked first in the suspiciousness rank) for the score are:  $\frac{1}{9} = 0.\bar{1}$ ,  $\frac{1}{19} = 0.052631$ ,  $\frac{1}{39} = 0.025641$ ,  
1096 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  
1097 rule is ranked last). The table also shows, in the last two columns, the average mean and standard  
1098 deviation values considering all case studies.

1099 Having a look at the average EXAM scores in the average-case scenario, we observe there are  
1100 8 techniques where less than 25% of the rules need to be inspected in order to locate the faulty  
1101 rule, i.e., their EXAM score is below 0.25. These are, ordered from lower to higher percentage,  
1102 *Mountford*, *Kulczynski2*, *Ochiai*, *Zoltar*, *Phi*, *Arithmetic Mean*, *Braun-Banquet*, and *Op2*. If we have a  
1103 look at these 8 techniques in the best-case scenario, we see that *Phi* and *Arithmetic Mean* have the  
1104 lowest, therefore best, numbers. However, their numbers are the worst among these techniques  
1105 in the worst-case scenario, implying that these techniques produced quite a large number of ties.  
1106 Observing the 8 techniques in the worst-case scenario, we see that *Mountford* and *Kulczynski2* are  
1107 able to locate the faulty rule by inspecting less than 23% of the rules, so they seem to be the best  
1108 techniques. In particular, *Kulczynski2* is able to locate the faulty rule first in the rank in 45% of the  
1109 cases in the worst-case scenario, and it ranks the guilty rule in the top 3 —i.e., only up to 3 rules  
1110 need to be inspected in order to locate the fault— in 70% of the cases. The EXAM scores for these  
1111 two techniques in the worst-case scenario are similar to the ones in the best- and average-case  
1112 scenarios, concluding that there are not many ties. These two techniques are closely followed by  
1113 *Ochiai* and *Zoltar*, techniques that do not produce many ties either and that are able to locate the  
1114 faulty rule by inspecting less than 25% of the rules in the worst-case scenario. We can conclude that  
1115 the four techniques with best results are, in this order, *Kulczynski2*, *Mountford*, *Ochiai* and *Zoltar*.  
1116 However, this ordering is not strict, since they behave slightly differently among them depending  
1117 on the case study. In particular, in order to locate the fault, these techniques lead the debugger to  
1118 inspect between 1.59 and 1.84 (out of 9) rules in *BibTex2DocBook*, 2.98 and 3.5 (out of 19) rules in  
1119 *CPL2SPL*, between 4.78 and 7.68 (out of 39) rules in *Ecore2Maude*, and between 2.65 and 2.69 (out of  
1120 8) rules in *UML2ER* in the average-case scenario. The average standard deviation in all scenarios is  
1121 around 0.2 for these four techniques, meaning that the results they provide are quite stable.

1122 Going back to the average-case scenario and looking for techniques that give bad results, there  
1123 are 5 techniques that need to inspect more than 30% of the rules in order to locate the faulty one,  
1124 namely *Barinel*, *Russel Rao*, *Tarantula*, *Dstar* and *Pierce*. Interestingly, the worst technique, so-called  
1125 *Pierce* [116], needs to inspect more than 63% of the rules. This means that it performs even worse  
1126

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	
1128	AC	Arithmetic Mean	.111	.284	.240	.105	.166	.171	.077	.175	.181	.188	.313	.229	.234	.205
1129		Barinel	.333	.391	.216	.184	.245	.146	.192	.237	.172	.406	.363	.176	.309	.178
1130		Braun-Banquet	.222	.277	.182	.105	.192	.199	.090	.160	.179	.188	.337	.319	.242	.219
1131		B-U & Buser	.333	.365	.236	.105	.172	.189	.077	.134	.142	.188	.336	.319	.252	.221
1132		Cohen	.333	.343	.229	.105	.169	.170	.077	.176	.182	.188	.313	.229	.250	.202
1133		Dstar	.222	.326	.310	.263	.293	.213	.231	.423	.344	.469	.537	.298	.395	.292
1134		Kulczynski2	.111	.177	.142	.105	.185	.203	.077	.133	.139	.188	.331	.313	.207	.199
1135		Mountford	.111	.204	.151	.053	.157	.194	.077	.123	.111	.188	.337	.316	.205	.193
1136		Ochiai	.111	.188	.147	.105	.185	.193	.077	.134	.143	.188	.333	.318	.210	.200
1137		Ochiai2	.444	.443	.270	.105	.180	.191	.077	.175	.182	.188	.313	.229	.278	.218
1138		Op2	.111	.182	.149	.105	.226	.217	.154	.245	.228	.188	.331	.313	.246	.227
1139		Phi	.111	.268	.237	.105	.166	.172	.077	.172	.179	.188	.313	.229	.230	.204
1140		Pierce	.833	.682	.285	.737	.636	.283	.667	.596	.203	.688	.611	.301	.631	.268
1141		Rogers & Tani.	.556	.454	.277	.053	.206	.235	.077	.132	.140	.188	.302	.289	.273	.235
1142		Russel Rao	.222	.255	.121	.105	.240	.222	.333	.367	.182	.375	.438	.262	.325	.197
1143	Simple Matching	.556	.454	.277	.053	.206	.235	.077	.132	.140	.188	.302	.289	.273	.235	
1144	Tarantula	.333	.398	.221	.092	.164	.191	.167	.211	.172	.438	.499	.259	.318	.211	
1145	Zoltar	.111	.177	.142	.105	.182	.198	.154	.197	.185	.188	.331	.313	.222	.209	
1146	Burgueño'15	.388	.436	.245	.105	.239	.224	.167	.317	.312	.375	.476	.297	.367	.269	
1147	BC	Arithmetic Mean	.111	.260	.233	.105	.161	.165	.026	.073	.112	.125	.196	.173	.173	.171
1148		Barinel	.333	.342	.235	.158	.229	.141	.051	.095	.112	.125	.168	.139	.208	.157
1149		Braun-Banquet	.222	.277	.182	.105	.180	.178	.026	.107	.175	.125	.308	.325	.218	.215
1150		B-U & Buser	.333	.365	.236	.105	.163	.171	.026	.081	.130	.125	.307	.326	.229	.216
1151		Cohen	.333	.320	.228	.105	.164	.163	.026	.075	.116	.125	.196	.173	.189	.170
1152		Dstar	.222	.325	.309	.263	.284	.202	.205	.372	.335	.438	.494	.310	.369	.289
1153		Kulczynski2	.111	.177	.142	.105	.176	.185	.026	.080	.132	.125	.301	.320	.184	.195
1154		Mountford	.111	.203	.151	.053	.148	.175	.026	.069	.097	.125	.304	.325	.181	.187
1155		Ochiai	.111	.188	.147	.105	.176	.176	.026	.081	.135	.125	.304	.325	.187	.196
1156		Ochiai2	.444	.416	.272	.105	.171	.174	.026	.072	.106	.125	.196	.173	.214	.181
1157		Op2	.111	.182	.149	.105	.221	.210	.026	.193	.241	.125	.301	.320	.225	.230
1158		Phi	.111	.245	.228	.105	.161	.165	.026	.070	.108	.125	.196	.173	.168	.169
1159		Pierce	.667	.587	.260	.658	.605	.262	.359	.410	.169	.375	.461	.308	.516	.250
1160		Rogers & Tani.	.556	.450	.277	.053	.195	.229	.026	.080	.131	.125	.274	.292	.250	.232
1161		Russel Rao	.111	.141	.122	.053	.196	.205	.026	.171	.247	.125	.261	.319	.192	.223
1162	Simple Matching	.556	.450	.277	.053	.195	.229	.026	.080	.131	.125	.274	.292	.250	.232	
1163	Tarantula	.333	.349	.241	.053	.146	.176	.026	.068	.112	.125	.304	.330	.217	.215	
1164	Zoltar	.111	.177	.142	.105	.173	.180	.026	.144	.192	.125	.301	.320	.199	.208	
1165	Burgueño'15	.333	.342	.219	.0526	.106	.086	.154	.270	.279	.312	.458	.296	.253	.172	
1166	WC	Arithmetic Mean	.111	.307	.283	.105	.171	.179	.128	.276	.317	.250	.429	.359	.296	.285
1167		Barinel	.444	.441	.261	.211	.260	.155	.256	.380	.303	.625	.557	.299	.409	.255
1168		Braun-Banquet	.222	.277	.182	.105	.205	.224	.154	.213	.190	.250	.365	.315	.265	.228
1169		B-U & Buser	.333	.365	.236	.105	.181	.211	.128	.187	.161	.250	.365	.315	.275	.231
1170		Cohen	.333	.367	.269	.105	.174	.177	.128	.278	.318	.250	.429	.359	.312	.281
1171		Dstar	.222	.326	.311	.263	.302	.229	.282	.474	.357	.500	.579	.290	.420	.297
1172		Kulczynski2	.111	.177	.142	.105	.194	.224	.128	.186	.155	.250	.360	.308	.229	.208
1173		Mountford	.111	.204	.151	.053	.167	.217	.128	.178	.134	.250	.369	.310	.230	.203
1174		Ochiai	.111	.188	.147	.105	.194	.215	.128	.188	.159	.250	.363	.314	.233	.209
1175		Ochiai2	.444	.470	.303	.105	.189	.213	.128	.279	.322	.250	.429	.359	.342	.299
1176		Op2	.111	.182	.149	.105	.231	.224	.231	.298	.219	.250	.360	.308	.268	.225
1177		Phi	.111	.292	.282	.105	.171	.179	.128	.273	.317	.250	.429	.359	.291	.284
1178		Pierce	1,000	.777	.335	.737	.667	.311	1,000	.781	.297	1,000	.761	.350	.746	.323
1179		Rogers & Tani.	.556	.458	.277	.053	.217	.244	.128	.184	.156	.250	.331	.290	.297	.242
1180		Russel Rao	.333	.369	.156	.105	.284	.249	.641	.563	.161	.625	.615	.260	.458	.206
1181	Simple Matching	.556	.458	.277	.053	.217	.244	.128	.184	.156	.250	.331	.290	.297	.242	
1182	Tarantula	.444	.448	.264	.105	.182	.214	.231	.354	.303	.750	.693	.269	.419	.262	
1183	Zoltar	.111	.177	.142	.105	.191	.220	.231	.250	.185	.250	.360	.308	.244	.214	
1184	Burgueño'15	.444	.529	.317	.131	.372	.413	.179	.365	.371	0.5	.494	.303	.440	.351	

1177 than random testing, and this is true in all case studies. Regarding the other four, *Dstar* needs to  
1178 inspect almost 37% of the rules even in the best-case scenario, what is not a good result either. If  
1179 we go to the worst-case scenario, all these techniques need to inspect more than 40% of the rules,  
1180 so we can conclude that they do not behave good and therefore we do not recommend to use them  
1181 when applying SBFL in the MT domain.

1182 The distributions of the results of each technique are graphically depicted in the box-plots of  
1183 Figure 6, so they are useful in order to analyze each case study separately and see if the conclusions  
1184 drawn so far are confirmed. The figure contains one box-plot per scenario (average-case labelled as  
1185 AC, best-case labelled as BC, and worst-case labelled as WC) and case study, where the Y and X  
1186 axis indicate the EXAM score and technique, respectively. These box-plots gather the results of  
1187 the EXAM scores obtained with all mutants and depict them in vertical boxes—ignore for now the  
1188 boxes for *Burgueño'15*, since they are commented in Section 4.5.3. The dots outside the boxes are  
1189 known as outliers.

1190 Having a look at the average-case scenarios in Figure 6, we can appreciate how techniques are  
1191 categorized in two groups. On the one hand, techniques that perform well are represented by small  
1192 boxes located at the bottom of each box-plot. We refer to these techniques as *good-performers*. On  
1193 the other hand, the boxes of techniques with bad performance are stretched and typically located  
1194 around the middle of the plot. We will name this group of techniques *bad-performers*. It is worth  
1195 mentioning that among the group of good-performers, the most reliable ones are those with smaller  
1196 vertical line segments above the box. This means that in the cases where faulty rules are difficult to  
1197 locate, they provide lower EXAM scores than other good-performers.

1198 For instance, in the *UML2ER* case study for the average-case scenario, the set of most-reliable  
1199 good-performers comprises of *Kulczynski2*, *Mountford*, *Zoltar* and *Ochiai*. In fact, the boxes of these  
1200 four techniques are quite stable and similar in all scenarios of all case studies. Some other techniques,  
1201 such as *Op2*, seem to provide similar performance, since for instance its boxes in the *UML2ER*  
1202 case study are similar to those of these four techniques. However, the boxes are clearly worse in  
1203 the *Ecore2Maude* and *CPL2SPL* case studies. At the same time, *Mountford* shows slightly better  
1204 performance than the other three good-performers in some box-plots, such as in the *CPL2SPL* case  
1205 study. Regarding the 5 techniques mentioned before that give bad results in the table of descriptive  
1206 statistics, they fit in the profile of bad-performers. We can observe that their boxes are not uniform  
1207 when comparing box-plots, having some boxes even located in the top of the charts.

1208 Finally, it is worth noting that the box-plots in the *UML2ER* case study present the highest  
1209 disparity among the three scenarios and that most techniques seem to behave worse in this case  
1210 study than in the other three, showing larger boxes. This suggests that it is more challenging for the  
1211 techniques to properly rank the faulty rule in this case study than in the other three case studies.  
1212 This may be due to the high use of rules inheritance in this case study, what might complicate the  
1213 location of the fault as explained in Section 4.2.1. Please also note, as commented in Section 4.2.3,  
1214 that fewer mutants have been created for this case study compared with the other three. This could  
1215 also have an impact in the results.

1216 In order to study the data from a different perspective, namely the average values, we have  
1217 constructed Figure 7. The figure presents a matrix where the suspiciousness computation techniques  
1218 are represented by rows, and the mutants of the different case studies are represented by columns.  
1219 Each cell is therefore colored according to its  $EXAM_{m,t}^{Average}$ , where  $m$  represents the mutant (X  
1220 axis) and  $t$  the technique (Y axis).

1221 As we can see in the color key, cells with lower values are lightly colored, while cells with higher  
1222 values are darkly colored. The lighter the shade of cell  $\langle i, j \rangle$ , the better has performed technique  $j$   
1223 in mutant  $i$  on average. Observing the four techniques with good performance mentioned before,  
1224

1225

1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241  
1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274

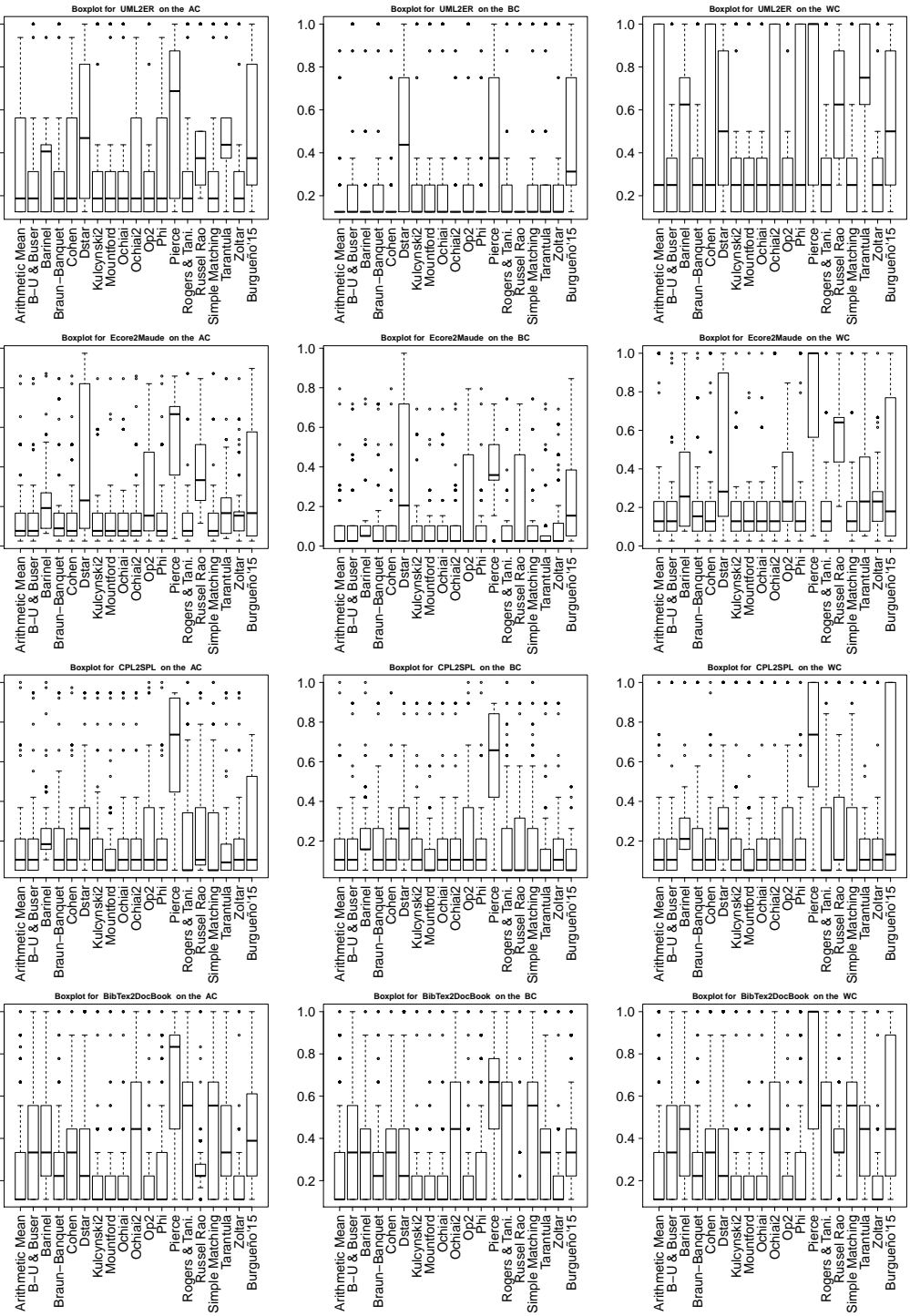


Fig. 6. Box-plot of the EXAM score of each technique per scenario and case study



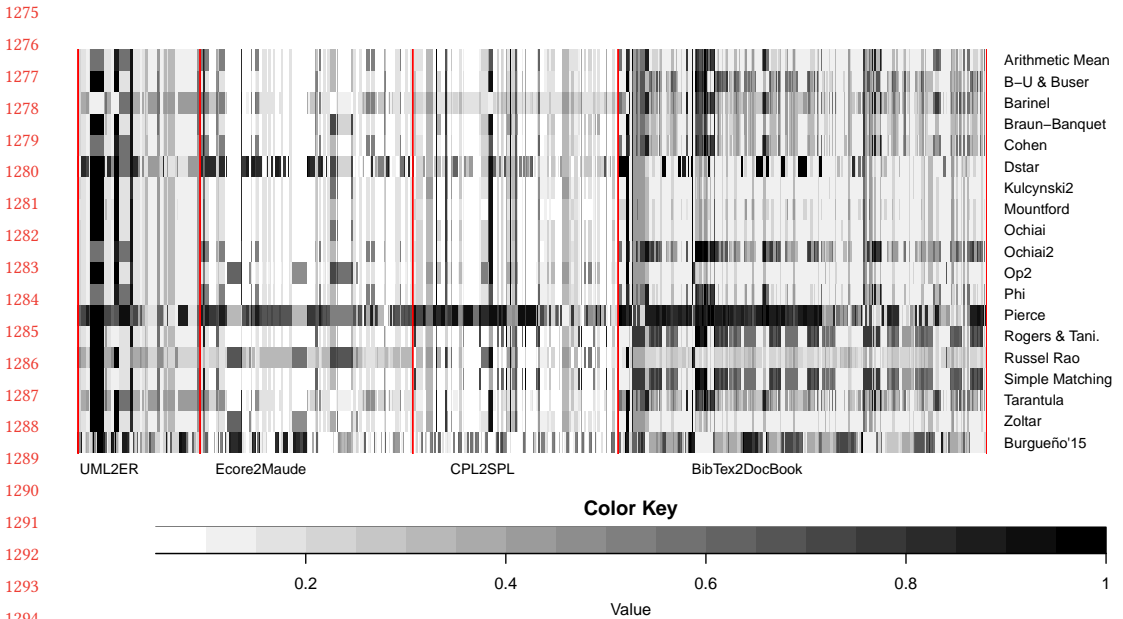


Fig. 7. Value per technique and case study

namely *Kulczynski2*, *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 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 75 seconds (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 paper 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,

1324 while the alternative hypothesis ( $H_1$ ) states that at least for one pair of techniques such difference is  
1325 statistically significant. Statistical tests provide the probability (named *p-value*) of getting the actual  
1326 observed results based on the assumption that the null hypothesis is true. P-values range in  $[0, 1]$ ,  
1327 for which researchers have established by convention that p-values under 0.05 represent so-called  
1328 statistically significant values, and are sufficient to reject the null hypothesis. The results of the  
1329 study do not follow a normal distribution, as confirmed by Shapiro-Wilk normality tests, thus the  
1330 Friedman test was used for the analysis [40]. The resulting p-values were  $< 1^{-10}$  for the results of the  
1331 four case studies, leading us to reject  $H_0$  for all of them. In order to find the specific techniques with  
1332 statistically significant differences, pairwise comparisons were performed using Conover-Iman's  
1333 Test [52]. More specifically, we compared all the possible pairs of techniques, out of 18 techniques  
1334 under study, i.e.,  $\binom{18}{2} = \frac{18!}{2!(18-2)!} = 153$  pairwise comparison per case study. Additionally, we applied  
1335 a correction of the p-values using the Holm post-hoc procedure [53], as recommended in [31]. The  
1336 results of the corrected p-values for the pairwise comparisons of all techniques are available on the  
1337 project's website [101]. In summary, the percentage of pairwise comparisons revealing statistically  
1338 significant differences was 96% in *Bibtex2DocBook*, 82% in *CPL2SPL*, 78% in *Ecore2Maude*, and 49%  
1339 in *UML2ER*. Again, these data highlight that the latter case study is the one giving worse and more  
1340 unstable results, which is consistent with the conclusion drawn in the analysis of the results in  
1341 Section 4.3.

1342  
1343 4.4.2 *Effect-size estimation.* In order to further investigate the differences between the different  
1344 suspiciousness computation techniques, Vargha and Delaney's  $\widehat{A}_{12}$  statistic [106] was used to  
1345 evaluate the effect size, i.e., determine which technique performs better and to what extent. Table 8  
1346 shows the effect size statistic for every pair of techniques. Each cell shows the  $\widehat{A}_{12}$  value obtained  
1347 when comparing the suspiciousness computation technique in the column against the technique in  
1348 the row. Vargha and Delaney [106] suggested thresholds for interpreting the effect size: 0.5 means  
1349 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)  
1350 or very large (0.72-1) difference in favour of the technique in the column; values below 0.5 indicate  
1351 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  
1352 of the technique in the row. Cells indicating medium, large, and very large differences in favor of  
1353 the column are shaded in light grey, grey, and dark grey, respectively. The values in boldface are  
1354 those where hypothesis test revealed statistical differences (p-value  $< 0.05$ ). As expected, there is  
1355 not a clear winner technique for all the case studies. However, the results confirm the superiority  
1356 of *Mountford*, *Kulczynski2*, *Ochiai* and *Zoltar*, showing from medium to large differences in 35, 30, 29  
1357 and 28 (out of 72) pairwise comparisons. Analogously, the results support the bad performance of  
1358 *Pierce* –outperformed by all of other techniques–, *Barinel* and *Tarantula*.

## 1360 4.5 Comparison Study

1361 In order to answer RQ4 and to study whether our approach performs well in the location of faults  
1362 in model transformations, we want to compare its effectiveness with a state-of-the-art approach  
1363 based on the static analysis of transformation rules and assertions that obtained good results [18].  
1364 We believe the comparison of our approach with this one is fair and adequate for several reasons.  
1365 First, the model transformation language used as proof of concept in both approaches is ATL.  
1366 Second, OCL assertions are used in both approaches as oracle, i.e., to determine whether a model  
1367 transformation is correct or not. Third, both approaches determine an order in which the rules  
1368 must be examined in order to locate the faulty rules. Fourth, we are using in the evaluation of  
1369 our approach the same four case studies proposed in [18]. Fifth, we are able to use the mutants  
1370 developed for evaluating our approach in order to evaluate the approach in [18], and we are also  
1371

Table 8. Effect size estimations

1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403  
1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421

	Arith. Mean	Braun-Banquet	Barinel	B-U & Buser	Cohen	Dstar	Kulczynski2	Mountford	Ochiai	Ochiai2	Op2	Phi	Pierce	Rogers & Tani.	Russel Rao	Simple Match	Tarantula	Zoltar	Burgueño'15	
BibTex2DocBook	Arith. Mean	.670	.330	.465	.387	.405	.475	.627	.568	.605	.327	.620	.526	.158	.336	.403	.336	.324	.626	.303
	Barinel	-	-	.659	.538	.572	.651	.813	.775	.799	.445	.807	.692	.230	.431	.703	.431	.493	.813	.449
	Braun-Banquet	.535	.341	-	.396	.425	.520	.679	.624	.659	.318	.672	.561	.150	.322	.489	.322	.335	.679	.308
	B-U & Buser	.613	.462	.604	-	.529	.591	.745	.702	.729	.419	.739	.634	.206	.405	.603	.405	.456	.745	.416
	Cohen	.595	.428	.575	.471	-	.574	.730	.685	.714	.389	.724	.618	.197	.387	.580	.387	.421	.730	.381
	Dstar	.525	.349	.480	.409	.426	-	.645	.585	.623	.358	.638	.550	.230	.373	.418	.373	.344	.644	.335
	Kulczynski2	.373	.187	.321	.255	.270	.355	-	.428	.473	.207	.493	.401	.084	.226	.220	.226	.182	.499	.178
	Mountford	.432	.225	.376	.298	.315	.415	.572	-	.545	.241	.564	.462	.102	.260	.298	.260	.219	.570	.212
	Ochiai	.395	.201	.341	.271	.286	.377	.527	.455	-	.220	.519	.423	.091	.239	.249	.239	.196	.525	.191
	Ochiai2	.673	.555	.682	.581	.611	.642	.793	.759	.780	-	.788	.690	.272	.476	.688	.476	.548	.793	.511
	Op2	.380	.193	.328	.261	.276	.362	.507	.436	.481	.212	-	.408	.087	.232	.229	.232	.188	.506	.184
	Phi	.474	.308	.439	.366	.382	.450	.599	.538	.577	.310	.592	-	.148	.318	.365	.318	.303	.598	.285
	Pierce	.842	.770	.850	.794	.803	.770	.916	.898	.909	.728	.913	.852	-	.749	.856	.749	.763	.916	.745
	Rogers & Tani.	.664	.569	.678	.595	.613	.627	.774	.740	.761	.524	.768	.682	.251	-	.651	.500	.562	.774	.516
	Russel Rao	.597	.297	.511	.397	.420	.582	.780	.702	.751	.312	.771	.635	.144	.349	-	.349	.290	.779	.275
	Simple Match	.664	.569	.678	.595	.613	.627	.774	.740	.761	.524	.768	.682	.251	.500	.651	-	.562	.774	.516
	Tarantula	.676	.507	.665	.544	.579	.656	.818	.781	.804	.452	.812	.697	.237	.438	.710	.438	-	.818	.457
Zoltar	.374	.187	.321	.255	.270	.356	.501	.430	.475	.207	.494	.402	.084	.226	.220	.226	.182	.499	.177	
Burgueño'15	.697	.550	.692	.584	.619	.664	.822	.788	.809	.489	.816	.715	.254	.483	.725	.483	.543	.822	-	
CPL2SPL	Arith. Mean	-	.274	.469	.500	.489	.285	.488	.551	.466	.477	.425	.502	.091	.505	.385	.505	.521	.488	.443
	Barinel	.726	-	.674	.738	.715	.435	.717	.777	.683	.694	.623	.738	.145	.687	.607	.687	.772	.720	.629
	Braun-Banquet	.531	.326	-	.526	.523	.337	.509	.570	.508	.518	.448	.526	.112	.525	.410	.525	.542	.511	.454
	B-U & Buser	.500	.262	.474	-	.491	.291	.490	.553	.473	.484	.431	.501	.099	.508	.392	.508	.524	.491	.446
	Cohen	.511	.285	.477	.509	-	.291	.497	.560	.476	.488	.432	.512	.092	.513	.394	.513	.531	.497	.451
	Dstar	.715	.565	.663	.709	.709	-	.684	.739	.687	.699	.608	.712	.184	.672	.586	.672	.722	.688	.604
	Kulczynski2	.512	.283	.491	.510	.503	.316	-	.562	.487	.496	.443	.511	.105	.519	.404	.519	.531	.504	.449
	Mountford	.449	.223	.430	.447	.440	.261	.438	-	.422	.432	.385	.447	.089	.465	.340	.465	.466	.438	.401
	Ochiai	.534	.317	.492	.527	.524	.313	.513	.578	-	.513	.447	.533	.107	.531	.411	.531	.546	.515	.461
	Ochiai2	.523	.306	.482	.516	.512	.301	.504	.568	.487	-	.437	.522	.104	.521	.399	.521	.536	.504	.454
	Op2	.575	.377	.552	.569	.568	.392	.557	.615	.553	.563	-	.572	.134	.568	.465	.568	.585	.562	.487
	Phi	.498	.262	.474	.499	.488	.288	.489	.553	.467	.478	.428	-	.091	.507	.390	.507	.522	.489	.446
	Pierce	.909	.855	.888	.901	.908	.816	.895	.911	.893	.896	.866	.909	-	.876	.857	.876	.906	.898	.852
	Rogers & Tani.	.495	.313	.475	.492	.487	.328	.481	.535	.469	.479	.432	.493	.124	-	.387	.500	.505	.482	.443
	Russel Rao	.615	.393	.590	.608	.606	.414	.596	.660	.589	.601	.535	.610	.143	.613	-	.613	.626	.599	.529
	Simple Match	.495	.313	.475	.492	.487	.328	.481	.535	.469	.479	.432	.493	.124	.500	.387	-	.505	.482	.443
	Tarantula	.479	.228	.458	.476	.469	.278	.469	.534	.454	.464	.415	.478	.094	.495	.374	.495	-	.469	.430
Zoltar	.512	.280	.489	.509	.503	.312	.496	.562	.485	.496	.438	.511	.102	.518	.401	.518	.531	-	.449	
Burgueño'15	.557	.371	.546	.554	.549	.396	.551	.599	.539	.546	.513	.554	.148	.557	.471	.557	.570	.551	-	
ECORE2MAUDE	Arith. Mean	-	.286	.501	.542	.500	.289	.539	.539	.539	.500	.392	.501	.071	.543	.184	.543	.397	.439	.406
	Barinel	.714	-	.720	.770	.713	.402	.768	.775	.766	.714	.583	.716	.096	.771	.253	.771	.600	.650	.511
	Braun-Banquet	.499	.280	-	.545	.499	.292	.545	.540	.545	.499	.404	.500	.077	.546	.137	.546	.370	.438	.403
	B-U & Buser	.458	.230	.455	-	.458	.270	.500	.495	.500	.458	.357	.459	.047	.502	.101	.502	.325	.390	.376
	Cohen	.500	.287	.501	.542	-	.289	.539	.539	.539	.500	.392	.501	.072	.543	.184	.543	.398	.439	.407
	Dstar	.711	.598	.708	.730	.711	-	.733	.730	.733	.710	.636	.713	.382	.730	.479	.730	.631	.667	.590
	Kulczynski2	.461	.232	.455	.500	.461	.267	-	.493	.499	.461	.357	.462	.046	.501	.099	.501	.327	.390	.375
	Mountford	.461	.225	.460	.505	.461	.270	.507	-	.507	.461	.355	.463	.032	.508	.084	.508	.323	.389	.371
	Ochiai	.461	.234	.455	.500	.461	.267	.501	.493	-	.461	.358	.462	.046	.501	.100	.501	.327	.391	.376
	Ochiai2	.500	.286	.501	.542	.500	.290	.539	.539	.539	-	.392	.501	.074	.543	.185	.543	.397	.439	.406
	Op2	.608	.417	.596	.643	.608	.364	.643	.645	.642	.608	-	.608	.131	.644	.261	.644	.497	.543	.482
	Phi	.499	.284	.500	.541	.499	.287	.538	.537	.538	.499	.392	-	.067	.542	.173	.542	.394	.438	.405
	Pierce	.929	.904	.923	.953	.928	.618	.954	.968	.954	.926	.869	.933	-	.956	.800	.956	.914	.915	.747
	Rogers & Tani.	.457	.229	.454	.498	.457	.270	.499	.492	.499	.457	.356	.458	.044	-	.095	.500	.322	.387	.374
	Russel Rao	.816	.747	.863	.899	.816	.521	.901	.916	.900	.815	.739	.827	.200	.905	-	.905	.753	.792	.633
	Simple Match	.457	.229	.454	.498	.457	.270	.499	.492	.499	.457	.356	.458	.044	.500	.095	-	.322	.387	.374
	Tarantula	.603	.400	.630	.675	.602	.369	.673	.677	.673	.603	.503	.606	.086	.678	.247	.678	-	.558	.480
Zoltar	.561	.350	.562	.610	.561	.333	.610	.611	.609	.561	.457	.562	.085	.613	.208	.613	.442	-	.445	
Burgueño'15	.594	.489	.597	.624	.593	.410	.625	.629	.624	.594	.518	.595	.253	.626	.367	.626	.520	.555	-	
UML2ER	Arith. Mean	-	.416	.509	.510	.500	.271	.514	.499	.513	.500	.198	.542	.361	.611	.542	.288	.514	.321	
	Barinel	.584	-	.624	.624	.584	.319	.629	.623	.629	.584	.629	.584	.244	.653	.465	.653	.364	.629	.420
	Braun-Banquet	.491	.376	-	.502	.491	.275	.503	.487	.502	.491	.503	.491	.261	.532	.259	.532	.261	.503	.322
	B-U & Buser	.490	.376	.498	-	.490	.274	.501	.485	.500	.490	.501	.490	.259	.530	.257	.530	.260	.501	.321
	Cohen	.500	.416	.509	.510	-	.271	.514	.499	.513	.500	.514	.500	.198	.542	.311	.542	.288	.514	.321
	Dstar	.729	.681	.725	.726	.729	-	.733	.725	.731	.729	.733	.729	.43						

1422 able to compute the EXAM values (cf. Section 4.5.2) for such approach, so we can fairly compare  
1423 both approaches. Finally, in the set of OCL assertions that we have created for each case study in  
1424 this work, we have included all the OCL assertions the authors in [18] proposed for evaluating their  
1425 approach<sup>8</sup> (cf. third column of Table 4). We have used those assertions, as well as some more that  
1426 we have defined, for evaluating our approach. Since the tools developed for the static approach [18]  
1427 are available (cf. Section 4.5.1), we have been able to run the approach with them. We have used  
1428 the complete set of OCL assertions for comparing both approaches. In order to demonstrate that  
1429 the results are not biased due to the new defined OCL assertions, we have also made a comparison  
1430 considering only the OCL assertions defined in [18]. This comparison is presented in the appendix  
1431 of this paper, where the results presented and conclusions drawn are very similar.

1432 In the following, we first summarize how the approach by Burgueño et al. works and computes  
1433 the rankings. Then, we explain how we are able to obtain EXAM values for such approach. Finally,  
1434 we present and discuss the results of the comparison.

1435  
1436 *4.5.1 Static Fault Localization in Model Transformations.* The paper by Burgueño et al. [18]  
1437 proposes a static approach for the localization of faults in model transformations. As in our approach,  
1438 ATL is the language used as proof of concept and the assertions that the transformations must  
1439 satisfy are also defined in OCL. Therefore, it follows the same methodology as proposed in this paper  
1440 (cf. Section 3.3). Also, like our approach, theirs is backed up by a tool. However, for determining  
1441 if any OCL assertion fails (step 2 in the methodology), their approach relies on an external tool,  
1442 namely *TractsTool* [19, 110]. This means that the user also needs to get familiarized with this other  
1443 tool.

1444 When executed, this static approach computes, for all OCL assertions, the order in which rules  
1445 should be inspected in order to locate the bug. To do so, it computes, for each pair <assertion,  
1446 rule>, the probability that the assertion failure comes from the rule making use of the common  
1447 denominator that both have, namely the structural elements belonging to the metamodels. The  
1448 approach builds on the following steps:

- 1449 (1) *Footprint Extraction.* The *structural elements*, referred to as *footprints*, of both model transfor-  
1450 mation and assertions are extracted.
- 1451 (2) *Footprint Matching.* The footprints extracted are compared for each rule and assertion.
- 1452 (3) *Matching Tables Calculation.* The percentage of types overlapping, so-called *alignment*, for  
1453 each transformation rule and assertion is calculated. This information is used to produce the  
1454 matching tables.
- 1455 (4) *Matching Tables Interpretation.* The resulting tables are analyzed for identifying the order in  
1456 which rules should be inspected in case any OCL assertion fails.

1457 In order to apply this approach, three tools need to be executed, two of which are proposed  
1458 and implemented in [18]. First, as mentioned before, the *TractsTool* is executed to check which  
1459 OCL assertions fail. Then, the *ATL Transformation Types Extractor* is executed to generate a model  
1460 with the footprints of the ATL transformation. Finally, the *Matching Tables Calculator* uses, among  
1461 others, such model as input and generates the matching tables, indicating also the order in which  
1462 rules should be inspected in case of failure<sup>9</sup>.

1463 Three matching tables are generated by this approach. They are matrices that have the OCL  
1464 assertions as rows and transformation rules as columns. Two of them need to be inspected in  
1465

1466  
1467 <sup>8</sup>The OCL assertions used in [18] are available on [http://atenea.lcc.uma.es/index.php/Main\\_Page/Resources/MTB](http://atenea.lcc.uma.es/index.php/Main_Page/Resources/MTB)

1468 <sup>9</sup>The *ATL Transformation Types Extractor* and *Matching Tables Calculator* tools are available on [http://atenea.lcc.uma.es/index.php/Main\\_Page/Resources/MTB/MTB](http://atenea.lcc.uma.es/index.php/Main_Page/Resources/MTB/MTB)

1471 order to determine the order in which rules should be inspected in case of failure –for a detailed  
1472 explanation, the interested reader is referred to [18].

1473  
1474 4.5.2 *EXAM Values Computation.* In order to obtain the suspiciousness-based rankings for the  
1475 approach in [18], we have obtained the matching tables of all 158 mutants, for which we have  
1476 made use of the available tools mentioned before, namely *ATL Transformation Types Extractor*  
1477 and *Matching Tables Calculator*. We have developed a program that, for each case study, takes the  
1478 matching tables as input together with a CSV file that contains information of the buggy rules  
1479 in each mutant and the OCL assertions that fail in such mutant. With those inputs, this program  
1480 computes the order in which the rules should be inspected for each of the OCL assertions that fail,  
1481 which is the same concept as the suspiciousness-based ranking proposed in spectrum-based fault  
1482 localization. Therefore, with this rules ordering, and since we have as input information of the  
1483 buggy rule of each mutant, we are able to easily compute the EXAM score. As output, our program  
1484 generates a CSV file indicating, for each mutant and each OCL assertion that fails, the EXAM score  
1485 in the best-, worst- and average-case scenarios (cf. Section 4.2.6).

1486 All the artifacts used for the comparison, namely the 158 mutants and 117 OCL assertions,  
1487 together with all the matching tables generated for all case studies are available on our project's  
1488 website [101].

1489  
1490 4.5.3 *Static Approach vs Dynamic Approach: Results.* As described before, Table 7 provides the  
1491 descriptive statistics of the EXAM score, where *Mountford*, *Kulczynski2*, *Ochiai* and *Zoltar* show  
1492 the best numbers. Regarding the static technique proposed by Burgueño et al. [18] (*Burgueño'15*  
1493 in the table), it performs worse than these techniques. In the average-case scenario, the static  
1494 approach needs to inspect around 37% of the rules in order to locate the fault, which is much  
1495 more than the 20% that needs to be inspected by the best techniques. In particular, for each case  
1496 study in the average-case scenario, the static technique needs to inspect 2.3 (out of 9) more rules in  
1497 *BibTex2DocBook* (25.9% of the MT), 1.387 (out of 19) more rules in *CPL2SPL* (7.3% of the MT), 7.18  
1498 (out of 39) more rules in *Ecore2Maude* (18.4% of the MT), and 1.3 (out of 8) more rules in *UML2ER*  
1499 (16.2% of the MT) compared with the best techniques in each case. Regarding the number of ties,  
1500 there is not a uniform behavior. For instance, in *BibTex2DocBook* and *CPL2SPL* there are clearly  
1501 more ties in the static technique compared to the best dynamic techniques, since the difference in  
1502 the EXAM score in the best- and worst-case scenarios is bigger. As for *Ecore2Maude* and *UML2ER*,  
1503 the number of ties seems to be similar among both approaches.

1504 Looking at the worst dynamic techniques, the static approach seems to behave better than some  
1505 of them. Having a look at the average mean (penultimate column), it behaves much better than  
1506 *Pierce* in the average-case scenario, since the latter technique needs to inspect more than 63% of  
1507 the rules in order to locate the fault. It also performs better than *Dstar* in this scenario, since this  
1508 technique needs to inspect more than 39% of the rules. Finally, the static technique by Burgueño et  
1509 al. performs worse than *Russel Rao* in the average-case scenario, but a bit better in the worst-case  
1510 scenario. Therefore, for now we can conclude that the static technique may behave better than 2  
1511 dynamic techniques and clearly behaves worse than other 15 techniques, but let us delve deeper  
1512 into the results.

1513 We further analyze the results by looking at each case study in the box-plots of Figure 6. In  
1514 general, we notice that the results of the static approach are typically similar among the three  
1515 scenarios, although the boxes are larger than those of most dynamic techniques, indicating a worse  
1516 performance. We can appreciate that the static approach behaves normally better than *Pierce*,  
1517 confirming our previous finding. As for *Dstar*, its boxes are in many plots larger than the ones  
1518 of the static approach. However, in other plots its boxes are smaller, so we cannot confirm the  
1519



1520 superiority of the static technique with regards to *Dstar*. For instance, in the *BibTex2DocBook* and  
1521 *CPL2SPL* case studies, the shape of the box-plots for *Dstar* seems to be normally better. Finally,  
1522 regarding *Russel Rao*, its boxes have in most cases better shapes than those of the static approach.

1523 The effect-size estimations of the statistical analysis for the static approach are displayed in  
1524 Table 8. To begin with, we can see in the *BibTex2DocBook* case study that the four best SBFL  
1525 techniques are clearly better than the static approach by Burgueño et al. [18], since the values in  
1526 the row of the static approach are above 0.78 for the corresponding cells, indicating a very-large  
1527 difference in favor of the technique in the column. Also, the technique that seemed to be similar to  
1528 the static approach, namely *Dstar*, is proved to be better in this case study. In general, the color of  
1529 the row shows that most techniques behave better than the static one.

1530 In fact, looking at the four case studies, the numbers in the cells of the rows of the static approach  
1531 and the columns with the best SBFL techniques –*Kulczynski2*, *Mountford*, *Zoltar* and *Ochiai*– are  
1532 always above 0.55, leaving no doubt that the static approach behaves worse. Besides, all these cells  
1533 reveal statistical differences (p-value <0.05, displayed in boldface in the table).

1534 The superiority of the static approach regarding *Pierce* is confirmed in all case studies. However,  
1535 it can not be concluded that it is better than any other of the techniques, since the rows of the static  
1536 technique do not present a value <0.5 in more than two case studies for any of the other techniques.  
1537 Finally, we see that in the *UML2ER* case study the static approach behaves generally much worse  
1538 than most techniques. An explanation can be that the static approach, based on types matching,  
1539 does not behave well in the presence of rule inheritance.

1540 In summary, we can confirm that all SBFL techniques have a better performance when locating  
1541 the faulty rule than the static technique, except for *Pierce*, where the static technique behaves clearly  
1542 better. Besides, the static approach normally presents more ties than the best dynamic techniques.

1543 Regarding performance in terms of runtime, static approaches are typically faster since they do  
1544 not need to execute the program under test. This is the same in our case, where the static approach  
1545 is faster. In any case, it requires to perform footprints extractions –both in OCL assertions and  
1546 ATL transformation rules– and footprints matching, that also requires some resources. Altogether,  
1547 the static approach takes from less than 1 second (in *UML2ER*) to 42 seconds (in *Ecore2Maude*) per  
1548 mutant, less than required by our dynamic approach (from 4 to 75 seconds, cf. Section 4.3).

## 1549 4.6 Discussion

1551 The results of the exhaustive experiments described in the previous sections allow us to answer  
1552 the research questions formulated in Section 4.1.

1553  
1554 4.6.1 *RQ1 - Feasibility*. The first research question, related to the feasibility of the approach, “*RQ1:*  
1555 *Is it possible to automate the process of locating faults in model transformations applying spectrum-*  
1556 *based techniques?*”, can be answered affirmatively. Indeed, we have automated the process of locating  
1557 the faulty rules in model transformations by means of a Java program<sup>10</sup> that orchestrates ATL model  
1558 transformations and uses the information stored in the traces to compute the suspiciousness-based  
1559 rankings based on the program spectra. This automation is explained in Section 3.4. All the artifacts  
1560 used as input and generated as output are available on our project’s website [101].

1561 Furthermore, even though our program has been implemented for ATL model transformations,  
1562 we are confident that it can be adapted for any transformation language that is able to store in a  
1563 trace model the result of the execution. In fact, the trace model is nothing but an output model.  
1564 Therefore, any model transformation language that is able to produce more than one output model  
1565 can generate a trace model as output.

1566  
1567 <sup>10</sup>Available on Github: [https://github.com/javitroya/SBFL\\_MT](https://github.com/javitroya/SBFL_MT)



1569 4.6.2 *RQ2 - Effectiveness.* The second research question has to do with the comparison of  
1570 the techniques evaluated with our automated approach: “*RQ2: How effective are state-of-the-art*  
1571 *techniques for suspiciousness computation in the localization of faulty rules in model transformations?*”  
1572 This question has to do with how well the different techniques are able to position the faulty rule  
1573 in the suspiciousness-based ranking. According to the results presented in Sections 4.3 and 4.4,  
1574 we can conclude that the top 4 most-effective techniques are *Kulczynski2*, *Mountford*, *Ochiai* and  
1575 *Zoltar*, the first two presenting slightly better overall results. At the other end, we have *Pierce* as  
1576 the least-effective technique. The top 3 of non-effective techniques is completed by *Barinel* and  
1577 *Tarantula*.

1578 4.6.3 *RQ3 - Accuracy.* The third research question is related to the accuracy of the approach:  
1579 “*RQ3: Is our approach able to accurately locate faulty rules in model transformations?*” The answer to  
1580 this question is related to the previous one, since depending on the effectiveness of the techniques  
1581 we will conclude whether the approach is accurate or not. In particular, we need to look at the  
1582 most effective techniques. Evaluation results revealed that the best techniques place the faulty  
1583 transformation rule among the three most suspicious rules in around 74% of the cases. Looking into  
1584 each of the four case studies, the best techniques allow the tester to locate the fault by inspecting,  
1585 on average, only 1.59 rules (out of 9) in *BibTex2DocBook*, 2.99 rules (out of 19) in *CPL2SPL*, 4.8 rules  
1586 (out of 39) in *Ecore2Maude* and 2.4 rules (out of 8) in *UML2ER*. According to these numbers, we  
1587 can conclude that the application of spectrum-based fault localization is accurate in the context of  
1588 model transformations if techniques such as *Mountford*, *Kulczynski2*, *Zoltar* and *Ochiai* are applied,  
1589 so we shall recommend to apply this approach to debug model transformations. These conclusions  
1590 are supported by the evaluation of more case studies available on our project’s website [101].  
1591

1592 4.6.4 *RQ4 - Dynamic vs Static.* Our last research question has to do with the comparison of our  
1593 dynamic approach with a notable static approach [18]: “*RQ4: How does our approach behave in com-*  
1594 *parison with a static approach?*” In summary, we can conclude that most dynamic techniques based  
1595 on spectrum computation are better than the static approach for the localization of faults in model  
1596 transformations. This was expected, since dynamic techniques execute the model transformation  
1597 –from which they extract a lot of information–, and the static approach does not. However, the static  
1598 approach is still clearly better than one dynamic technique, namely *Pierce*. Furthermore, it also  
1599 behaves better than other techniques, such as *Dstar*, *Tarantula*, *Simple Matching*, *Rogers & Tanimoto*  
1600 and *Barinel* in some case studies. It is also noteworthy that both approaches are complementary  
1601 and so it should be possible to define heuristics for the selection of the best technique on each  
1602 application scenario, or even combine them. For example, it is better to apply the static approach in  
1603 environments with low resources or when the transformations are very expensive to execute [65],  
1604 for instance in the case of transforming very large models [15, 25], and when it is not possible to  
1605 get model instances of the source metamodel at the time of developing the model transformation.  
1606

## 1607 4.7 Threats to Validity

1608 According to Wohlin et al. [113], there are four basic types of validity threats that can affect the  
1609 validity of our study. We cover each of these in the following paragraphs.  
1610

1611 4.7.1 *Conclusion Validity Threats.* Threats to the conclusion validity are concerned with the  
1612 issues that affect the ability to draw correct conclusions from the data obtained from the experiments.  
1613 In order to mitigate these threats, we have applied statical analysis to confirm the conclusions drawn  
1614 from the means and figures, and we have used the specific statistical tests and effect size measures  
1615 recommended by the guidelines on empirical methodology. Furthermore, all the assumptions  
1616

1618 required for the application of the tests were checked, and the raw data and scripts for replication  
1619 are available in the companion materials of this paper [101].  
1620

1621 *4.7.2 Construct Validity Threat.* It is concerned with the relationship between theory and what  
1622 is observed. A possible construct validity threat (known as the mono-method bias) is related to the  
1623 use of one single metric, the so-called EXAM score, to evaluate the performance of the approach  
1624 and the suspiciousness computation techniques compared. Other metrics have been proposed [116],  
1625 such as the T-score [68], P-score [123] and N-score [43]. However, EXAM score is an accepted  
1626 metric for measuring the quality of spectrum-based fault localization techniques, and has been used  
1627 in a variety of works.[116] Moreover, we have decided to obtain the EXAM score as it is directly  
1628 applicable in the context of model transformations. By considering the transformation rules as  
1629 units of examination, the EXAM score is easy to understand, since it is directly proportional to the  
1630 amount of rules to be examined, rather than to an indirect measurement in terms of the amount of  
1631 code that does not need to be examined, as proposed by other scores.

1632 Another possible construct validity threat is the mono-operation bias, which is related to the use  
1633 of a single treatment or technique that could bias our conclusions. Since we compare the approach  
1634 with a static alternative [18] and have used up to 18 suspiciousness-computation techniques and 4  
1635 use cases in our experiments, we consider that this threat is neutralized.  
1636

1637 *4.7.3 Internal Validity Threats.* These threats are related to those factors that might affect the  
1638 results of our evaluation. First of all, we may remark that this is a debugging approach, not a  
1639 testing approach. Therefore, the objective of this work is not to generate high-quality test models,  
1640 something addressed in related papers [6, 39, 44, 48, 96], but to localize the faults that triggered  
1641 test failures. In fact, a key point in favor of our approach is that it can be used in conjunction  
1642 with any method for test model generation, either random or guided. For evaluating our work, we  
1643 have developed a light-weight random model generator that, given any metamodel, produces a  
1644 user-defined number of random model instances, as explained in Section 4.2.2. With this generator  
1645 we have obtained a set composed of 100 source models in the test suite of each case study, so a  
1646 total of 400 models have been generated. These models have achieved full coverage – all rules  
1647 and lines of code have been exercised – in all case studies. However, using more complex input  
1648 model generation approaches [6, 48] may be required in those cases where random generation is  
1649 not enough to achieve a sufficient coverage.

1650 A second threat is that we have used in total 117 OCL assertions in the first study and 44 in  
1651 the dynamic-vs-static comparison study (cf. Table 4). We have tried to minimize this threat by  
1652 constructing a set of OCL assertions that cover much of the specifications of the transformations.  
1653 Besides, for the comparison study to be fair, we have taken the OCL assertions proposed in [18].  
1654 Third, we have tried to create a large set of mutants, composed of 158 of them, and we have aimed  
1655 at maximizing the variation of semantic faults and mutation operators used. Having used more or  
1656 fewer mutants could have had an impact in the results. For instance, we recall that fewer mutants  
1657 have been created for the *UML2ER* case study than for any of the other case studies, as commented  
1658 in Section 4.2.3. Having used different mutation operators could have also had an impact in the  
1659 results. For instance, some approaches propose mutation operators that yield run-time errors, such  
1660 as the work by Sánchez-Cuadrado et al. [92], which proposes a powerful approach that relies on  
1661 static analysis and type inference to locate, among others, run-time errors. However, please bear in  
1662 mind that the purpose of the approach presented in this paper is to localize semantic faults, i.e., it  
1663 needs the model transformation to finish and produce output models, so that their satisfaction can  
1664 be checked against the set of OCL assertions available. That is why we have used a subset [99] of the  
1665 operators defined in [92] and that mimic semantic faults likely to be made by programmers [78], as  
1666

1667 explained in Section 4.2.3. In any case, our approach is complementary to those aiming at spotting  
1668 bugs that produce run-time errors [92].

1669 As a final threat to internal validity, we may mention a weakness of SBFL, and generally of all  
1670 fault localization techniques [116], which is the incapability of locating bugs resulting from missing  
1671 code [122]. Same thing happens with our approach, it is likely to produce bad results if there  
1672 are missing rules. For this reason, and as commented above, our approach can be complemented  
1673 with Sánchez-Cuadrado et al.'s [92] approach, which identifies rules absence with a static analysis.  
1674 Indeed, the target elements created in a transformation rule typically reference or are referenced  
1675 by target elements created in other transformation rules, so static analysis is a good technique for  
1676 identifying the absence of rules that should create referencing or referenced target elements. For  
1677 instance, in the model transformation shown in Listing 1, the target elements created by rule *Main*  
1678 reference the target elements created by all the other rules. Likewise, the target elements created  
1679 by all the other rules are referenced by those created by rule *Main*. Therefore, the absence of any  
1680 of these rules can be detected with a static analysis tailored at examining that there will be no  
1681 dangling references among the target elements created. Finally, this threat can also be mitigated  
1682 with the definition of proper OCL assertions. For instance, in the transformation of Listing 1, the  
1683 specification should dictate that there must be an element of class *DocBook* created for each element  
1684 of class *BibTexFile*, so that the number of instances of both classes must be the same after executing  
1685 the model transformation. This can be expressed with assertion *OCL4* in Listing 2. Therefore, even  
1686 if we do not count on approaches like the one by Sánchez-Cuadrado et al. [92], the non-satisfaction  
1687 of assertions such as *OCL4* can help the developer realize a rule is missing.

1688  
1689 **4.7.4 External Validity.** These threats have to do with the extent to which it is possible to  
1690 generalize the findings of the experiments. The first threat is that the results of our experiments  
1691 have been obtained with four case studies, which externally threatens the generalizability of our  
1692 results. To mitigate this threat, we have tried to select a set of model transformations that considers  
1693 all ATL constructs and where the model transformations differ in their domains, size of metamodels  
1694 and transformation, and variability of features used within the transformations, as reflected in  
1695 Table 3. Furthermore, we have selected the same case studies as those used in the related paper  
1696 compared to our approach in Section 4.5, published in 2015 in the *IEEE Transactions on Software*  
1697 *Engineering* journal [18]. Second, we have analyzed a set of 18 techniques for the computation of  
1698 the suspiciousness-based rankings. Although it is a large set, result of doing a thorough literature  
1699 review, we might have left aside some techniques that could give better results than the ones  
1700 obtained with the best techniques of our study. Also, we have implemented our approach for ATL  
1701 due to its importance both in industria and academia, so it would be interesting to test it with other  
1702 transformation languages. However, we do believe our approach would produce similar results for  
1703 any model transformation language based in rules as long as the result of their executions can be  
1704 stored in traces (cf. Section 2.2.3), that allows to construct the coverage matrix and error vector  
1705 and, therefore, apply SBFL techniques.

1706 There are two other threats related to the external validity of the results that have to do with the  
1707 program spectra creation in our implementation. In particular, we have used in our prototype the  
1708 ATLAS transformation language and have considered the rules, of any type, as unit of examination  
1709 and therefore as the components to be considered for constructing the spectra. Should we also have  
1710 considered helpers in the spectra, the results of techniques effectiveness could have been different.  
1711 This decision has been made considering related works that also check (ATL) transformations  
1712 correctness against OCL assertions. While some approaches only check whether an assertion is  
1713 violated or not by a model transformation [7, 19, 42, 81, 110], others propose to locate the fault  
1714 when an assertion is not satisfied [7, 22, 23], but none of them inspect the helpers—they remain at

1715

1716 the rule level. Crucial for our decision has been the static approach for locating faults proposed  
1717 by Burgueño et al. [18], which does not consider helpers either and only locate faults in ATL  
1718 rules. Should we have considered them, the thorough comparison with this approach presented in  
1719 Section 4.5 would have been unfair. After having proved the effectiveness of SBFL techniques in  
1720 the model transformation domain according to the extensive evaluation presented in this paper, a  
1721 natural evolution of this work is to perform a thorough study considering helpers to check if these  
1722 techniques remain effective. In any case, if our current approach determines that a rule is faulty,  
1723 and it is calling a helper, then the user of the approach would inspect the rule and, if (s)he sees no  
1724 fault, (s)he would proceed by inspecting the helper, so this threat is reduced.

1725 Finally, the other threat is that the components considered in our approach might be too coarse-  
1726 grained: our approach works at rule level. This means that the user would need, for example, to put  
1727 more effort in locating a bug in a big rule than when doing it in a small rule. However, the complexity  
1728 of transformation rules and model transformations is inherent to the bridges they try to build  
1729 among different semantic domains, and different types of model transformations can be written  
1730 depending on the problem to be solved [28, 64]. For instance, the creator of ATL recommends  
1731 to use declarative code as much as possible<sup>11</sup>. Besides, some approaches exist for modularizing  
1732 model transformations, so that they become as easy-to-understand and reusable as possible [34, 90].  
1733 Like with the threat before, another reason that led us to work at rule level in this approach is  
1734 that related works that aim to locate bugs in model transformations against OCL satisfaction also  
1735 propose approaches at rule level [22, 23], and specially the work with which we do an extensive  
1736 comparison [18].

## 1737 5 RELATED WORK

1738 Due to the lack of oracles and formal semantics in model transformation languages, some approaches  
1739 propose to translate the transformation specifications to other domains where formal treatment  
1740 is possible. For instance, Troya and Vallecillo propose to translate ATL to the rewriting logic  
1741 framework Maude [103], where some formal analysis can be performed, although the translation is  
1742 not fully automated. Anastasakis et al. propose to translate QVT model transformations to Alloy in  
1743 order to verify if some properties hold for the transformation, and there are also approaches for  
1744 verifying contracts for ATL transformations based on the Coq proof assistant [20, 84]. Oakes et al.  
1745 propose to translate the declarative part of ATL to the visual graph-based model transformation  
1746 engine DSLTrans [81]. Visual contracts similar to our OCL assertions but less expressive can be  
1747 then tested for satisfaction in DSLTrans. Similar visual contracts, using a visual language with  
1748 formal semantics called PaMoMo, are used by Guerra et al. [47], but in this case their approach  
1749 compiles such contracts into QVT and their satisfiability is checked with the PACO-Checker tool. A  
1750 big difference of our approach with these is that we do not need to leave the model transformation  
1751 development environment in order to check for the correctness of the MTs, so our approach stays  
1752 within the Eclipse Modeling Framework dealing with Ecore metamodels and XMI models and the  
1753 user does not need to be familiar with any other domain-specific language such as Maude, DSLTrans,  
1754 Alloy or Coq. Furthermore, our approach helps locate the faulty rules, that is not addressed in these  
1755 approaches.

1756 As in our approach, Cheng et al. [22] propose to verify if ATL transformations satisfy OCL  
1757 assertions. However, in order to prove the correctness of the ATL transformation, they encode both  
1758 the OCL assertions and the ATL transformation specification into the Boogie language [89]. Boogie  
1759 is a procedure-oriented language that is based on Hoare-logic. Then, their developed VeriATL  
1760 verification system indicates whether the ATL specification satisfies the specified OCL assertions or  
1761

1762 <sup>11</sup><http://www.idi.ntnu.no/emner/ttd4250/Slides/M2M-atl-intria1117.pdf>  
1763

1765 not. However, this approach does not report useful feedback to help the transformation developers  
1766 fix the fault, which is the main objective of our approach. Cheng and Tisi [23] then build on this  
1767 approach and tool (VeriATL) with the goal of localizing the fault by applying natural deduction  
1768 and program slicing. However, instead of offering the developer with a rules ranking according to  
1769 their chance to contain a bug, their approach determines scenarios, which are slices of the model  
1770 transformation under test, where a certain OCL assertion is not satisfied together with the proof  
1771 tree. This is achieved by deriving sub-goals from the OCL assertions. Since this approach aims at  
1772 locating a fault from a different perspective than ours, they can complement each other.

1773 Burgueño et al. [18] propose an approach with a similar purpose as ours, but their approach is  
1774 static. They also count on ATL model transformations and OCL assertions that must be checked  
1775 for correctness and try to locate the faulty rule without translating the OCL assertions nor ATL  
1776 transformations to any formal language. Their approach proposes to locate the faulty rules based  
1777 on matching functions that automatically establish alignments among the metamodels footprints  
1778 appearing in the transformation rules and those present in the OCL assertions. A comparison of  
1779 this approach and the one we present in this paper has been done in Section 4.5, where we have  
1780 seen that most techniques for the spectrum-based localization of model transformations give better  
1781 results than this static approach. Furthermore, this approach does not check if an OCL assertion is  
1782 satisfied, but it resorts to the Tracts tool [19]. Contrarily, our approach does not need any input  
1783 from external tools. A good aspect of the static approach is that it does not need any input model,  
1784 since actually the transformation is not executed, and it requires shorter runtimes. This aspect  
1785 makes this approach very useful in several situations. For example, it is better to apply the static  
1786 approach in environments with low resources or when the transformations are very expensive to  
1787 execute [65], for instance in the case of transforming very large models [15, 25], and when it is  
1788 not possible to get model instances of the source metamodel at the time of developing the model  
1789 transformation. Both approaches are therefore tangential.

1790 There are other approaches that propose static analysis for debugging model transformations.  
1791 Sánchez-Cuadrado et al. [92] combine static analysis and constraint solving in order to discover  
1792 errors in ATL model transformations such as navigation errors (like invalid collection operations  
1793 and operators), disconformities between the types used in the transformation and those declared  
1794 in its source/target metamodels, integrity constraints regarding the semantics of ATL, problems  
1795 related to dependencies between transformation rules and, in summary, any error that the current  
1796 syntactic checker of ATL is not able to identify. They even provide possible suitable quick fixes  
1797 based on speculative analysis [91]. These approaches have meant an important milestone in the  
1798 evolution of ATL. Our approach is orthogonal to these and, consequently, can serve to complement  
1799 them. Finally, Sánchez-Cuadrado et al. [93] have built, on top of their so-called *anATL*yer tool  
1800 described in the previous cited papers, an approach for checking contracts specified in the target  
1801 language. Their approach translates these target contracts into source contracts by using the model  
1802 transformation, so that they can predict, without the need to execute the model transformation,  
1803 whether any specific input model will yield (in)correct target models. Since they use the model  
1804 transformation for generating source contracts, it has to be correct. Therefore, different from our  
1805 approach, this is not targeted to debugging model transformations, but to statically check target  
1806 constraints in a light-weight manner.

1807 The approach we present in this work is perfectly in line with the approach we presented in [102]  
1808 with the aim of locating bugs in three application scenarios of model transformations, namely  
1809 regression testing, incremental transformations and migrations among transformation languages.  
1810 Thus, the approach in [102] proposes to automatically derive OCL assertions from a given ATL  
1811 model transformation, which are satisfied by the transformation. The approach applies a technique  
1812  
1813



1814 known as metamorphic testing [94]. By applying metamorphic testing, and after identifying a  
1815 set of patterns that normally takes place in the trace information stored after the execution of a  
1816 model transformation, it is able to automatically derive so-called likely metamorphic relations,  
1817 which can be seen as precisely the OCL assertions used in the current work. In this way, (i) in  
1818 regression testing, (ii) when an original transformation is migrated to a different transformation  
1819 language or (iii) an incremental transformation is developed with the same behavior of the original  
1820 transformation, the approach presented in this paper can be used in order to check whether the  
1821 OCL assertions obtained for the original transformation by the approach in [102] are satisfied in  
1822 the latter evolved or modified transformations. Metamorphic testing has also been applied by He et  
1823 al. [51], in this case for bidirectional model transformation testing.

1824 We recall that this paper focuses on debugging and not testing. Thus, we do not impose any  
1825 constraint on how the source models are generated, either manually or automatically. In any case,  
1826 some proposals for generating models have been proposed in the literature, where some of them  
1827 require input by the tester. For instance, the model generator in [17] requires the tester to provide  
1828 metamodel fragments as input, or the one in [96] requires input from the MMCC external tool [38]  
1829 to provide model fragments. Other approaches propose the generation of models in different formats  
1830 such as the Human Usable Textual Notation [41], so they need to be transformed prior to their  
1831 use as input for model transformation languages integrated in the Eclipse Modeling Framework  
1832 such as ATL. Some other more sophisticated model generators try to derive a set of input models  
1833 from model transformations [44], what is not desired in our case because we may be debugging  
1834 erroneous transformations, and from OCL constraints [6, 48]. Most of these approaches can be  
1835 used for generating test models for our approach. However, as explained in Section 4.2.2, we have  
1836 used a light-weight model generator that, given a metamodel, it returns a set of random models  
1837 conforming to such metamodel, where the models present certain variability among them with  
1838 respect to the classes of the metamodel. Since none of the case studies contain complex graph  
1839 constraints as preconditions that are difficult to cover with random graph generation, the models  
1840 we have generated have obtained full coverage in all case studies, i.e., they have exercised all rules.  
1841 However, in other cases, obtaining models with full coverage may require the use of more complex  
1842 and computationally-expensive methods and tools [5, 6, 44, 48, 96].

1843

1844

## 1845 6 CONCLUSION

1846 In this paper we have presented the first approach for debugging model transformations following  
1847 a spectrum-based fault localization (SBFL) approach. We have implemented and automated it for  
1848 the ATLAS transformation language due to its importance in both industry and academia. However,  
1849 we are confident that the approach can be extensible to any model transformation language as  
1850 long as it can store the output of its execution in a trace model. The implemented automation has  
1851 allowed us to perform a thorough evaluation.

1852 Taking as input the model transformation under test and a set of source models and OCL  
1853 assertions that serve as oracle, our approach determines which assertions are not satisfied and,  
1854 for each of them, it ranks the transformation rules according to their suspiciousness of being  
1855 the faulty rule causing the failure. We have compared the effectiveness of 18 state-of-the-art  
1856 techniques proposed in the literature for the suspiciousness computation of program components  
1857 (e.g., statements) in the context of model transformations. The evaluation has been carried out  
1858 using four case studies that differ regarding the application domains, size of metamodels and the  
1859 number and types of ATL features used. Our experiments conclude that the best techniques place  
1860 the faulty transformation rule among the three most suspicious rules in around 74% of the cases.

1861

1862



1863 These conclusions are supported by more case studies, other than the four presented in this paper,  
1864 whose evaluation is available on our project's website [101].

1865 We have also evaluated our approach by comparing it with a static approach that presented  
1866 notable results [18]. The conclusion is that applying dynamic techniques based on spectra compu-  
1867 tation allows to identify the faulty rule more quickly. However, the runtime of the static technique  
1868 is shorter, and it does not need any input model, since the model transformation is not executed.  
1869 Therefore, both approaches are tangential and can complement each other.

1870 Summarizing, we have proved the effectiveness in the context of model transformations of  
1871 SBFL, a technique never applied before for localizing faults in this domain. We have proved it  
1872 is feasible to automate such technique in this domain, offering novel ways of debugging model  
1873 transformations. Despite we have obtained good effectiveness results, further experiments can be  
1874 performed as future work. For instance, we can consider helpers in the program spectra, and even  
1875 each line of code could be considered as a component. In both cases, the trace model used has to  
1876 be extended. Also, in order to break ties in the suspiciousness rankings, we could use the rules  
1877 execution frequency, as some works have proposed for procedural programs [1, 67].  
1878

1879

## 1879 VERIFIABILITY

1880 For the sake of verifiability, our prototype as well as all artifacts of the experiments are available  
1881 on our project's website [101]. For each case study, it is available the transformation and its  
1882 metamodels, the OCL assertions defined, the transformation mutants together with information  
1883 of the mutation operators applied and the OCL assertions that fail with each mutant, as well as  
1884 the CSV files with the results generated by our program for all mutants and all OCL assertions.  
1885 For the comparison study, it is available for each case study the subset of mutants used together  
1886 with the matching tables generated with the approach in [18] for each mutant, and the subset of  
1887 OCL assertions obtained from [18]. Several files with statistical results and raw data and scripts  
1888 for replication are also available. Finally, the implemented prototype is available on Github: [https://github.com/javitroya/SBFL\\_MT](https://github.com/javitroya/SBFL_MT).  
1889  
1890

1891

## 1891 REFERENCES

- 1892
- 1893 [1] ABREU, R., GONZALEZ-SANCHEZ, A., AND VAN GEMUND, A. J. C. Exploiting count spectra for bayesian fault localization.  
1894 In *Proc. of the 6th International Conference on Predictive Models in Software Engineering* (2010), PROMISE '10, ACM,  
1895 pp. 12:1–12:10.
  - 1896 [2] ABREU, R., ZOETWEIJ, P., AND GEMUND, A. J. C. v. Spectrum-Based Multiple Fault Localization. In *Proc. of IEEE/ACM*  
1897 *International Conference on Automated Software Engineering* (Washington, DC, USA, 2009), ASE '09, IEEE Computer  
1898 Society, pp. 88–99.
  - 1899 [3] ABREU, R., ZOETWEIJ, P., GOLSTEIJN, R., AND VAN GEMUND, A. J. A practical evaluation of spectrum-based fault  
1900 localization. *Journal of Systems and Software* 82, 11 (2009), 1780 – 1792.
  - 1901 [4] ABREU, R., ZOETWEIJ, P., AND VAN GEMUND, A. J. C. On the Accuracy of Spectrum-based Fault Localization. In  
1902 *Testing: Academic and Industrial Conference Practice and Research Techniques - MUTATION (TAICPART-MUTATION*  
1903 *2007)* (2007), pp. 89–98.
  - 1904 [5] ALI, S., IQBAL, M. Z., AND ARCURI, A. Improved heuristics for solving ocl constraints using search algorithms. In  
1905 *Proceedings of the 2014 Annual Conference on Genetic and Evolutionary Computation* (New York, NY, USA, 2014),  
1906 GECCO '14, ACM, pp. 1231–1238.
  - 1907 [6] ALI, S., IQBAL, M. Z., ARCURI, A., AND BRIAND, L. C. Generating test data from ocl constraints with search techniques.  
1908 *IEEE Transactions on Software Engineering* 39, 10 (2013), 1376–1402.
  - 1909 [7] ALMENDROS-JIMÉNEZ, J. M., AND BECERRA-TERÓN, A. Automatic generation of ecore models for testing ATL transfor-  
1910 mations. In *6th International Conference on Model and Data Engineering (MEDI)* (2016), vol. 9893 of LNCS, Springer,  
1911 pp. 16–30.
  - [8] ANASTASAKIS, K., BORDBAR, B., GEORG, G., AND RAY, I. On challenges of model transformation from uml to alloy.  
*Software & Systems Modeling* 9, 1 (2008).

- 1912 [9] ARANEGA, V., MOTTU, J.-M., ETIEN, A., DEGUEULE, T., BAUDRY, B., AND DEKEYSER, J.-L. Towards an automation of  
1913 the mutation analysis dedicated to model transformation. *Software Testing, Verification and Reliability* 25, 5-7 (2015),  
1914 653–683.
- 1915 [10] ARENDT, T., BIERMANN, E., JURACK, S., KRAUSE, C., AND TAENTZER, G. Henshin: Advanced Concepts and Tools for  
1916 In-Place EMF Model Transformations. In *Proc. of MODELS* (2010), vol. 6394, pp. 121–135.
- 1917 [11] ASSIRI, F. Y., AND BIEMAN, J. M. Fault localization for automated program repair: effectiveness, performance, repair  
1918 correctness. *Software Quality Journal* 25, 1 (2017), 171–199.
- 1919 [12] ATL. ATL Zoo. <http://www.eclipse.org/atl/atlTransformations>, 2006.
- 1920 [13] BARR, E. T., HARMAN, M., MCMINN, P., SHAHBAZ, M., AND YOO, S. The Oracle Problem in Software Testing: A Survey.  
1921 *IEEE Transactions on Software Engineering* 41, 5 (May 2015), 507–525.
- 1922 [14] BAUDRY, B., DINH-TRONG, T., MOTTU, J.-M., SIMMONDS, D., FRANCE, R., GHOSH, S., FLEUREY, F., AND LE TRAON, Y. Model  
1923 Transformation Testing Challenges. In *Proc. of the ECMDA workshop on Integration of Model Driven Development and*  
1924 *Model Driven Testing* (2006).
- 1925 [15] BENELALLAM, A., GÓMEZ, A., TISI, M., AND CABOT, J. Distributed model-to-model transformation with atl on mapreduce.  
1926 In *Proceedings of the 2015 ACM SIGPLAN International Conference on Software Language Engineering* (New York, NY,  
1927 USA, 2015), SLE 2015, ACM, pp. 37–48.
- 1928 [16] BRAMBILLA, M., CABOT, J., AND WIMMER, M. *Model-Driven Software Engineering in Practice*. Morgan&Claypool, 2012.
- 1929 [17] BROTTIER, E., FLEUREY, F., STEEL, J., BAUDRY, B., AND TRAON, Y. L. Metamodel-based test generation for model  
1930 transformations: an algorithm and a tool. In *Proc. of ISSRE 2006* (2006), pp. 85–94.
- 1931 [18] BURGUEÑO, L., TROYA, J., WIMMER, M., AND VALLECILLO, A. Static Fault Localization in Model Transformations. *IEEE*  
1932 *Tansactions on Software Engineering* 41, 5 (May 2015), 490–506.
- 1933 [19] BURGUEÑO, L., WIMMER, M., TROYA, J., AND VALLECILLO, A. Tractstool: Testing MTs based on contracts. In *Invited*  
1934 *Talks, Demonstration Session, Poster Session, and ACM Student Research Competition (MODELS 2013)* (Oct. 2013), CEUR.  
1935 <http://ceur-ws.org/Vol-1115/poster5.pdf>.
- 1936 [20] CALEGARI, D., LUNA, C., SZASZ, N., AND TASISTRO, A. A Type-Theoretic Framework for Certified Model Transformations.  
1937 In *Proc. of SBMF* (2010), pp. 112–127.
- 1938 [21] CARIOU, E., MARVIE, R., SEINTURIER, L., AND DUCHIEN, L. OCL for the Specification of Model Transformation  
1939 Contracts. In *Proc. of the OCL and Model Driven Engineering Workshop* (2004).
- 1940 [22] CHENG, Z., MONAHAN, R., AND POWER, J. F. *A Sound Execution Semantics for ATL via Translation Validation*. Springer  
1941 International Publishing, Cham, 2015, pp. 133–148.
- 1942 [23] CHENG, Z., AND TISI, M. *A Deductive Approach for Fault Localization in ATL Model Transformations*. Springer Berlin  
1943 Heidelberg, Berlin, Heidelberg, 2017, pp. 300–317.
- 1944 [24] CICHETTI, A., DI RUSCIO, D., ERAMO, R., AND PIERANTONIO, A. JTL: A Bidirectional and Change Propagating  
1945 Transformation Language. In *Software Language Engineering*, vol. 6563 of LNCS. Springer, 2011, pp. 183–202.
- 1946 [25] CLASEN, C., DIDONET DEL FABRO, M., AND TISI, M. Transforming Very Large Models in the Cloud: a Research  
1947 Roadmap. In *First International Workshop on Model-Driven Engineering on and for the Cloud* (Copenhagen, Denmark,  
1948 July 2012), Springer.
- 1949 [26] CLAVEL, M., DURÁN, F., EKER, S., LINCOLN, P., MARTÍ-OLIET, N., MESEGUER, J., AND TALCOTT, C. *All About Maude – A*  
1950 *High-Performance Logical Framework*, vol. 4350 of LNCS. Springer, 2007.
- 1951 [27] CSERTÁN, G., HUSZERL, G., MAJZIK, I., PAP, Z., PATARICZA, A., AND VARRÓ, D. VIATRA - visual automated transformations  
1952 for formal verification and validation of UML models. In *Proc. of the 17th International Conference on Automated*  
1953 *Software Engineering (ASE'02)* (2002), IEEE/ACM, pp. 267–270.
- 1954 [28] CZARNECKI, K., AND HELSEN, S. Feature-based survey of model transformation approaches. *IBM Systems Journal* 45, 3  
1955 (2006), 621–646.
- 1956 [29] DA SILVA, A. R. Model-driven engineering: A survey supported by the unified conceptual model. *Computer Languages,*  
1957 *Systems & Structures* 43 (2015), 139 – 155.
- 1958 [30] DE LARA, J., AND VANGHELUWE, H. AToM3: A Tool for Multi-formalism and Meta-modelling. In *Proc. of the 5th*  
1959 *International Conference on Fundamental Approaches to Software Engineering (FASE'02)*, vol. 2306 of LNCS. Springer,  
1960 2002, pp. 174–188.
- [31] DERRAC, J., GARCÍA, S., MOLINA, D., AND HERRERA, F. A practical tutorial on the use of nonparametric statistical  
tests as a methodology for comparing evolutionary and swarm intelligence algorithms. *Swarm and Evolutionary*  
*Computation* 1, 1 (2011), 3 – 18.
- [32] DURÁN, F., ZSCHALER, S., AND TROYA, J. On the Reusable Specification of Non-functional Properties in DSLs. In *5th*  
*International Conference on Software Language Engineering (SLE 2012). Revised Selected Papers* (2013), LNCS, Springer,  
pp. 332–351.
- [33] FALLERI, J.-R., HUCHARD, M., AND NEBUT, C. Towards a Traceability Framework for Model Transformations in

- 1961 Kermeta. In *ECMDA-TW'06: ECMDA Traceability Workshop* (2006), pp. 31–40.
- 1962 [34] FLECK, M., TROYA, J., KESSENTINI, M., WIMMER, M., AND ALKHAZI, B. Model transformation modularization as a  
1963 many-objective optimization problem. *IEEE Transactions on Software Engineering* 43, 11 (2017), 1009–1032.
- 1964 [35] FLECK, M., TROYA, J., AND WIMMER, M. Marrying Search-based Optimization and Model Transformation Technology.  
1965 In *Proc. of NasBASE* (2015), pp. 1–16.
- 1966 [36] FLECK, M., TROYA, J., AND WIMMER, M. Search-Based Model Transformations. *Journal of Software: Evolution and  
1967 Process* 28, 12 (2016), 1081–1117.
- 1968 [37] FLECK, M., TROYA, J., AND WIMMER, M. Search-Based Model Transformations with MOMoT. In *Proc. of Theory and  
1969 Practice of Model Transformations (ICMT)* (2016), Springer, pp. 79–87.
- 1970 [38] FLEUREY, F., BAUDRY, B., MULLER, P., AND TRAON, Y. L. Qualifying input test data for model transformations. *Software  
1971 and System Modeling* 8, 2 (2009), 185–203.
- 1972 [39] FLEUREY, F., BAUDRY, B., MULLER, P.-A., AND TRAON, Y. L. Qualifying input test data for model transformations.  
1973 *Software & Systems Modeling* 8, 2 (2009), 185–203.
- 1974 [40] FRIEDMAN, M. A comparison of alternative tests of significance for the problem of m rankings. *The Annals of  
1975 Mathematical Statistics* 11, 1 (1940), 86–92.
- 1976 [41] GINER, P., AND PELECHANO, V. Test-driven development of model transformations. In *Proc. of MODELS'09*, vol. 5795  
1977 of LNCS. Springer, 2009, pp. 748–752.
- 1978 [42] GOGOLLA, M., AND VALLECILLO, A. *Tractable Model Transformation Testing*. Springer Berlin Heidelberg, 2011, pp. 221–  
1979 235.
- 1980 [43] GONG, C., ZHENG, Z., LI, W., AND HAO, P. Effects of Class Imbalance in Test Suites: An Empirical Study of Spectrum-  
1981 Based Fault Localization. In *IEEE 36th Annual Computer Software and Applications Conference Workshops* (2012),  
1982 pp. 470–475.
- 1983 [44] GONZÁLEZ, C. A., AND CABOT, J. *ATLTest: A White-Box Test Generation Approach for ATL Transformations*. Springer,  
1984 2012, pp. 449–464.
- 1985 [45] GREENYER, J., AND KINDLER, E. Comparing relational model transformation technologies: implementing Query/View/  
1986 Transformation with Triple Graph Grammars. *Software and System Modeling* 9, 1 (2010), 21–46.
- 1987 [46] GUERRA, E. Specification-driven test generation for model transformations. In *Theory and Practice of Model Transfor-  
1988 mations* (Berlin, Heidelberg, 2012), Z. Hu and J. de Lara, Eds., Springer Berlin Heidelberg, pp. 40–55.
- 1989 [47] GUERRA, E., DE LARA, J., WIMMER, M., KAPPEL, G., KUSEL, A., RETSCHITZEGGER, W., SCHÖNBÖCK, J., AND SCHWINGER,  
1990 W. Automated verification of model transformations based on visual contracts. *Automated Software Engineering* 20, 1  
1991 (2013), 5–46.
- 1992 [48] GUERRA, E., AND SOEKEN, M. Specification-driven model transformation testing. *Software & Systems Modeling* 14, 2  
1993 (2015), 623–644.
- 1994 [49] HAMILL, M., AND GOSEVA-POPSTOJANOVA, K. Common trends in software fault and failure data. *IEEE Transactions on  
1995 Software Engineering* 35, 4 (2009), 484–496.
- 1996 [50] HARROLD, M. J., ROTHERMEL, G., SAYRE, K., WU, R., AND YI, L. An empirical investigation of the relationship between  
1997 spectra differences and regression faults. *Software Testing, Verification and Reliability* 10, 3 (2000), 171–194.
- 1998 [51] HE, X., CHEN, X., CAI, S., ZHANG, Y., AND HUANG, G. Testing Bidirectional Model Transformation Using Metamorphic  
1999 Testing. *Information and Software Technology* (2018).
- 2000 [52] HOLLANDER, M., WOLFE, D. A., AND CHICKEN, E. *Nonparametric statistical methods*. John Wiley & Sons, 2013.
- 2001 [53] HOLM, S. A simple sequentially rejective multiple test procedure. *Scand. J. Statist.* 6, 2 (1979), 65–70.
- 2002 [54] INRIA. ATL Transformation Example: BibTeXML to DocBook, 2005. [https://www.eclipse.org/atl/atlTransformations/  
2003 BibTeXML2DocBook/ExampleBibTeXML2DocBook\[v00.01\].pdf](https://www.eclipse.org/atl/atlTransformations/BibTeXML2DocBook/ExampleBibTeXML2DocBook[v00.01].pdf).
- 2004 [55] JANSSEN, T., ABREU, R., AND VAN GEMUND, A. J. Zoltar: a spectrum-based fault localization tool. In *Proc. of the 2009  
2005 ESEC/FSE workshop on Software integration and evolution @ runtime (SINTER 2009)* (New York, NY, USA, 2009), ACM,  
2006 pp. 23–30.
- 2007 [56] JÉZÉQUEL, J.-M., BARAIS, O., AND FLEUREY, F. Model Driven Language Engineering with Kermeta. In *Generative and  
2008 Transformational Techniques in Software Engineering III*, vol. 6491 of LNCS. Springer, 2011, pp. 201–221.
- 2009 [57] JIA, Y., AND HARMAN, M. Higher order mutation testing. *Information and Software Technology* 51, 10 (2009), 1379 –  
1393. Source Code Analysis and Manipulation, SCAM 2008.
- [58] JIA, Y., AND HARMAN, M. An Analysis and Survey of the Development of Mutation Testing. *IEEE Transactions on  
Software Engineering* 37, 5 (2011), 649–678.
- [59] JONES, J. A., AND HARROLD, M. J. Empirical Evaluation of the Tarantula Automatic Fault-localization Technique. In  
*Proc. of the 20th IEEE/ACM International Conference on Automated Software Engineering* (New York, NY, USA, 2005),  
ASE '05, ACM, pp. 273–282.
- [60] JOUAULT, F. Loosely Coupled Traceability for ATL. In *Workshop Proc. of ECMDA* (2005).

- 2010 [61] JOUAULT, F., ALLILAIRE, F., BÉZIVIN, J., AND KURTEV, I. ATL: A Model Transformation Tool. *Sci. Comput. Program.* 72,  
2011 1-2 (2008), 31–39.
- 2012 [62] JOUAULT, F., ALLILAIRE, F., BÉZIVIN, J., KURTEV, I., AND VALDURIEZ, P. ATL: A QVT-like Transformation Language. In  
2013 *Companion to the 21st ACM SIGPLAN Symposium on Object-oriented Programming Systems, Languages, and Applications*  
2014 (2006), OOPSLA '06, ACM, pp. 719–720.
- 2015 [63] JOUAULT, F., BÉZIVIN, J., CONSEL, C., KURTEV, I., AND LATRY, F. Building DSLs with AMMA/ATL, a Case Study on SPL  
2016 and CPL Telephony Languages. In *ECOOP Workshop on Domain-Specific Program Development* (Nantes, France, July  
2017 2006).
- 2018 [64] KIRSIE, B. G. Guideline and Evaluation of Model Transformation Engineering Approaches. Master's thesis, KTH  
Industrial Engineering and Management, Sweden, 2010.
- 2019 [65] KOLOVOS, D. S., ROSE, L. M., MATRAGKAS, N., PAIGE, R. F., GUERRA, E., CUADRADO, J. S., DE LARA, J., RÁTH, I., VARRÓ,  
2020 D., TISI, M., AND CABOT, J. A research roadmap towards achieving scalability in model driven engineering. In  
2021 *Proceedings of the Workshop on Scalability in Model Driven Engineering* (New York, NY, USA, 2013), BigMDE '13, ACM,  
2022 pp. 2:1–2:10.
- 2023 [66] KÜHNE, T. Matters of (meta-) modeling. *Software & Systems Modeling* 5, 4 (2006), 369–385.
- 2024 [67] LEE, H. J., NAISH, L., AND RAMAMOCHANARAO, K. Effective software bug localization using spectral frequency weighting  
2025 function. In *IEEE 34th Annual Computer Software and Applications Conference* (2010), pp. 218–227.
- 2026 [68] LIU, C., FEI, L., YAN, X., HAN, J., AND MIDKIFF, S. P. Statistical debugging: A hypothesis testing-based approach. *IEEE*  
*Transactions on Software Engineering* 32, 10 (Oct. 2006), 831–848.
- 2027 [69] LUCIA, THUNG, F., LO, D., AND JIANG, L. Are faults localizable? In *2012 9th IEEE Working Conference on Mining Software*  
*Repositories (MSR)* (2012), pp. 74–77.
- 2028 [70] LUCIA, L., LO, D., JIANG, L., THUNG, F., AND BUDI, A. Extended comprehensive study of association measures for fault  
2029 localization. *Journal of Software: Evolution and Process* 26, 2 (2014), 172–219.
- 2030 [71] LÚCIO, L., AMRANI, M., DINGEL, J., LAMBERS, L., SALAY, R., SELIM, G., SYRIANI, E., AND WIMMER, M. Model Transfor-  
2031 mation Intents and Their Properties. *Software and System Modeling* (2014), 1–35.
- 2032 [72] LUDEWIG, J. Models in software engineering – an introduction. *Software and Systems Modeling* 2, 1 (2003), 5–14.
- 2033 [73] ÁLVARO JIMÉNEZ, VARA, J. M., BOLLATI, V. A., AND MARCOS, E. Metagem-trace: Improving trace generation in model  
2034 transformation by leveraging the role of transformation models. *Science of Computer Programming* 98 (2015), 3 – 27.
- 2035 [74] MAO, X., LEI, Y., DAI, Z., QI, Y., AND WANG, C. Slice-based Statistical Fault Localization. *J. Syst. Softw.* 89 (Mar. 2014),  
2036 51–62.
- 2037 [75] MAXWELL, A. E., AND PILLINER, A. E. G. Deriving coefficients of reliability and agreement for ratings. *British Journal*  
*of Mathematical and Statistical Psychology* 21, 1 (1968), 105–116.
- 2038 [76] MELLOR, S. J., SCOTT, K., UHL, A., WEISE, D., AND SOLEY, R. M. *MDA distilled: principles of model-driven architecture*,  
2039 vol. 88. Addison-Wesley, 2004.
- 2040 [77] MORENO-DELGADO, A., DURÁN, F., ZSCHALER, S., AND TROYA, J. Modular DSLs for Flexible Analysis: An e-Motions  
2041 Reimplementation of Palladio. In *10th European Conference on Modelling Foundations and Applications (ECMFA 2014)*  
2042 (2014), LNCS, Springer, pp. 132–147.
- 2043 [78] MOTTU, J.-M., BAUDRY, B., AND LE TRAON, Y. *Mutation Analysis Testing for Model Transformations*. Springer Berlin  
2044 Heidelberg, Berlin, Heidelberg, 2006, pp. 376–390.
- 2045 [79] NAISH, L., LEE, H. J., AND RAMAMOCHANARAO, K. A Model for Spectra-based Software Diagnosis. *ACM Trans. Softw.*  
*Eng. Methodol.* 20, 3 (Aug. 2011), 11:1–11:32.
- 2046 [80] NAISH, L., NEELOFAR, AND RAMAMOCHANARAO, K. Multiple bug spectral fault localization using genetic programming.  
2047 In *2015 24th Australasian Software Engineering Conference* (2015), pp. 11–17.
- 2048 [81] OAKES, B. J., TROYA, J., LÚCIO, L., AND WIMMER, M. Full Contract Verification for ATL using Symbolic Execution.  
*Journal on Software & System Modeling* (2016), 1–35.
- 2049 [82] OMAR, E., GHOSH, S., AND WHITLEY, D. Subtle higher order mutants. *Information and Software Technology* 81 (2017), 3  
2050 – 18.
- 2051 [83] PERSSON, M., TÖRNGREN, M., QAMAR, A., WESTMAN, J., BIEHL, M., TRIPAKIS, S., VANGHELUWE, H., AND DENIL, J. A  
2052 characterization of integrated multi-view modeling in the context of embedded and cyber-physical systems. In *Proc.*  
*of the 11th ACM International Conference on Embedded Software* (Piscataway, NJ, USA, 2013), EMSOFT '13, IEEE Press,  
2053 pp. 10:1–10:10.
- 2054 [84] POERNOMO, I., AND TERRELL, J. Correct-by-Construction Model Transformations from Partially Ordered Specifications  
2055 in Coq. In *Proc. of ICFEM* (2010), pp. 56–73.
- 2056 [85] PROJECT, E. M. Atlas Transformation Language – ATL. <http://eclipse.org/atl>, 2015.
- 2057 [86] QI, Y., MAO, X., LEI, Y., AND WANG, C. Using Automated Program Repair for Evaluating the Effectiveness of Fault  
2058 Localization Techniques. In *Proc. of the 2013 International Symposium on Software Testing and Analysis* (New York,



- 2059 NY, USA, 2013). ISSTA 2013, ACM, pp. 191–201.
- 2060 [87] RIVERA, J. E., DURAN, F., AND VALLECILLO, A. A Graphical Approach for Modeling Time-dependent Behavior of  
2061 DSLs. In *Proc. of the IEEE Symposium on Visual Languages and Human-Centric Computing (VL/HCC'09)* (2009), IEEE,  
2062 pp. 51–55.
- 2063 [88] ROSE, L. M., HERRMANNSDOERFER, M., MAZANEK, S., VAN GORP, P., BUCHWALD, S., HORN, T., KALNINA, E., KOCH, A.,  
2064 LANO, K., SCHÄTZ, B., AND WIMMER, M. Graph and model transformation tools for model migration. *Software &  
2065 Systems Modeling* 13, 1 (2014), 323–359.
- 2066 [89] RUSTAN, K., AND LEINO, M. This is boogie 2. Tech. rep., 2008. Manuscript KRML 178.
- 2067 [90] SÁNCHEZ-CUADRADO, J., AND GARCÍA-MOLINA, J. Modularization of model transformations through a phasing  
2068 mechanism. *Software & Systems Modeling* 8, 3 (Jul 2009), 325–345.
- 2069 [91] SÁNCHEZ-CUADRADO, J., GUERRA, E., AND DE LARA, J. Quick fixing ATL transformations with speculative analysis.  
2070 *Software & Systems Modeling* (2016), 1–35.
- 2071 [92] SÁNCHEZ-CUADRADO, J., GUERRA, E., AND DE LARA, J. Static analysis of model transformations. *IEEE Transactions on  
2072 Software Engineering* 43, 9 (2017), 868–897.
- 2073 [93] SÁNCHEZ-CUADRADO, J., GUERRA, E., DE LARA, J., CLARISÓ, R., AND CABOT, J. Translating Target to Source Constraints in  
2074 Model-to-Model Transformations. In *ACM/IEEE 20th International Conference on Model Driven Engineering Languages  
2075 and Systems (MODELS)* (2017), pp. 12–22.
- 2076 [94] SEGURA, S., FRASER, G., SANCHEZ, A., AND RUIZ-CORTES, A. A survey on metamorphic testing. *IEEE Transactions on  
2077 Software Engineering* 42, 9 (2016), 805–824.
- 2078 [95] SELIM, G. M. K., WANG, S., CORDY, J. R., AND DINGEL, J. *Model Transformations for Migrating Legacy Models: An  
2079 Industrial Case Study*. Springer Berlin Heidelberg, Berlin, Heidelberg, 2012, pp. 90–101.
- 2080 [96] SEN, S., BAUDRY, B., AND MOTTU, J.-M. Automatic Model Generation Strategies for Model Transformation Testing. In  
2081 *Proc. of ICMT'09* (2009), vol. 5563 of LNCS, Springer, pp. 148–164.
- 2082 [97] SENDALL, S., AND KOZACZYNSKI, W. Model Transformation: The Heart and Soul of Model-Driven Software Development.  
2083 *IEEE Software* 20, 5 (2003), 42–45.
- 2084 [98] TAENTZER, G. AGG: A Graph Transformation Environment for Modeling and Validation of Software. In *Proc. of the  
2085 2nd Workshop on Applications of Graph Transformations with Industrial Relevance (AGTIVE'03)*, vol. 3062 of LNCS.  
2086 Springer, 2003, pp. 446–453.
- 2087 [99] TROYA, J., BERGMAYR, A., BURGUENO, L., AND WIMMER, M. Towards systematic mutations for and with ATL model  
2088 transformations. In *Proc. of the IEEE 8th Int. Conference on Software Testing, Verification and Validation Workshops  
2089 (ICSTW), 2015 IEEE Eighth International Conference on* (2015), pp. 1–10.
- 2090 [100] TROYA, J., SEGURA, S., PAREJO, J. A., AND RUIZ-CORTES, A. An Approach for Debugging Model Transformations  
2091 Applying Spectrum-Based Fault Localization. In *XXII Jornadas de Ingeniería del Software y Bases de Datos (JISBD)*  
2092 (2017).
- 2093 [101] TROYA, J., SEGURA, S., PAREJO, J. A., AND RUIZ-CORTÉS, A. Spectrum-Based Fault Localization in Model Transformations.  
2094 <https://gestionproyectos.us.es/projects/itim/wiki>, 2017.
- 2095 [102] TROYA, J., SEGURA, S., AND RUIZ-CORTÉS, A. Automated inference of likely metamorphic relations for model transfor-  
2096 mations. *Journal of Systems and Software* 136 (2018), 188 – 208.
- 2097 [103] TROYA, J., AND VALLECILLO, A. A Rewriting Logic Semantics for ATL. *Journal of Object Technology* 10 (2011), 5:1–29.
- 2098 [104] TROYA, J., AND VALLECILLO, A. Specification and simulation of queuing network models using domain-specific  
2099 languages. *Computer Standards & Interfaces* 36, 5 (2014), 863 – 879.
- 2100 [105] VALLECILLO, A., AND GOGOLLA, M. Typing model transformations using tracts. In *Proc. of 5th Int. Conf. on Theory and  
2101 Practice of Model Transformations (ICMT 2012)* (2012), Springer, pp. 56–71.
- 2102 [106] VARGHA, A., AND DELANEY, H. D. A critique and improvement of the cl common language effect size statistics of  
2103 mcgraw and wong. *Journal of Educational and Behavioral Statistics* 25, 2 (2000), 101–132.
- 2104 [107] WAGELAAR, D. Using ATL/EMFTVM for import/export of medical data. In *2nd Software Development Automation  
2105 Conference* (2014). <https://es.slideshare.net/DennisWagelaar/wagelaar-sda2014>.
- 2106 [108] WALSH, L. M. N. *DocBook: The Definitive Guide*. O'Reilly & Associates, 1999.
- 2107 [109] WARMER, J., AND KLEPPE, A. *The Object Constraint Language: Getting your models ready for MDA*. Addison Wesley,  
2003.
- [110] WIMMER, M., AND BURGUENO, L. Testing M2T/T2M Transformations. In *Proc. of Int. Conference on Model Driven  
Engineering Languages and Systems (MoDELS'13)* (2013), Springer, pp. 203–219.
- [111] WIMMER, M., KAPPEL, G., SCHÖNBÖCK, J., KUSEL, A., RETSCHITZEGGER, W., AND SCHWINGER, W. A Petri Net based  
debugging environment for QVT Relations. In *Proc. of ASE'09* (2009), IEEE, pp. 3–14.
- [112] WIMMER, M., MARTÍNEZ, S., JOUAULT, F., AND CABOT, J. A catalogue of refactorings for model-to-model transformations.  
*Journal of Object Technology* 11, 2 (Aug. 2012), 2:1–40.

- 2108 [113] WOHLIN, C., RUNESON, P., HÖST, M., OHLSSON, M. C., AND REGNELL, B. *Experimentation in Software Engineering*.  
2109 Springer, 2012.
- 2110 [114] WONG, W. E., DEBROY, V., GAO, R., AND LI, Y. The DStar Method for Effective Software Fault Localization. *IEEE*  
2111 *Transactions on Reliability* 63, 1 (March 2014), 290–308.
- 2112 [115] WONG, W. E., DEBROY, V., LI, Y., AND GAO, R. Software Fault Localization Using DStar (D\*). In *Proc. of IEEE Sixth*  
2113 *International Conference on Software Security and Reliability* (June 2012), pp. 21–30.
- 2114 [116] WONG, W. E., GAO, R., LI, Y., ABREU, R., AND WOTAWA, F. A Survey on Software Fault Localization. *IEEE Transactions*  
2115 *on Software Engineering* 42, 8 (2016), 707–740.
- 2116 [117] XIE, X. *On the Analysis of Spectrum-based Fault Localization*. PhD thesis, Faculty of Information and Communication  
2117 Technologies, Swinburne University of Technology, Australia, 2012.
- 2118 [118] XIE, X., CHEN, T. Y., KUO, F.-C., AND XU, B. A Theoretical Analysis of the Risk Evaluation Formulas for Spectrum-based  
2119 Fault Localization. *ACM Trans. Softw. Eng. Methodol.* 22, 4 (Oct. 2013), 31:1–31:40.
- 2120 [119] XUE, X., AND NAMIN, A. S. How significant is the effect of fault interactions on coverage-based fault localizations? In  
2121 *2013 ACM / IEEE International Symposium on Empirical Software Engineering and Measurement* (2013), pp. 113–122.
- 2122 [120] YOO, S. Evolving human competitive spectra-based fault localisation techniques. In *Search Based Software Engineering*  
2123 (Berlin, Heidelberg, 2012), G. Fraser and J. Teixeira de Souza, Eds., Springer, pp. 244–258.
- 2124 [121] YU, Y., JONES, J. A., AND HARROLD, M. J. An Empirical Study of the Effects of Test-suite Reduction on Fault Localization.  
2125 In *Proc. of the 30th International Conference on Software Engineering* (New York, NY, USA, 2008), ICSE '08, ACM,  
2126 pp. 201–210.
- 2127 [122] ZHANG, X., TALLAM, S., GUPTA, N., AND GUPTA, R. Towards Locating Execution Omission Errors. In *Proceedings of*  
2128 *the 28th ACM SIGPLAN Conference on Programming Language Design and Implementation* (2007), PLDI '07, ACM,  
2129 pp. 415–424.
- 2130 [123] ZHANG, Z., CHAN, W., TSE, T., HU, P., AND WANG, X. Is non-parametric hypothesis testing model robust for statistical  
2131 fault localization? *Information and Software Technology* 51, 11 (2009), 1573 – 1585. Third IEEE International Workshop  
2132 on Automation of Software Test (AST 2008) Eighth International Conference on Quality Software (QSIC 2008).

## 2130 A APPENDIX - STATIC-VS-DYNAMIC COMPARISON WITH REDUCED SET OF OCL 2131 ASSERTIONS

2132 The comparison study presented in Section 4.5 has compared our approach with the static approach  
2133 by Burgueño et al. [18]. In that comparison, we have used all OCL assertions: those taken from [18]  
2134 and several others defined for evaluating this work. This appendix is devoted to present the figures  
2135 and results for the comparison using only the OCL assertions defined in [18]. This way we show  
2136 that the new OCL assertions defined for evaluating our approach are not tailored to defeat the  
2137 approach by Burgueño et al.

2138 As it is shown in the second part of the third column in Table 9, 44 OCL assertions, out of the  
2139 total of 117 assertions created for the four case studies, have been taken from the static approach  
2140 we want to compare our approach with [18]. First of all, out of the 158 mutants we have created for  
2141 the four case studies, we select those that make any of the 44 OCL assertions fail. They are a total  
2142 of 104 mutants, so they are the ones to be considered in this comparison. The second part of the  
2143 fifth and third columns of Table 9 display the number of mutants and OCL assertions considered in  
2144 each case study for the comparison study, respectively. All the artifacts used for the comparison,  
2145 namely the 104 mutants and 44 OCL assertions, together with all the matching tables generated for  
2146 all case studies are available on our project's website [101].

2147 The approach by Burgueño et al. as well as the way we compute the EXAM values are explained  
2148 in Sections 4.5.1 and 4.5.2, respectively. The descriptive statistics of the EXAM score provided by  
2149 the techniques when applied to the 104 MT mutants are shown in Table 10.

2150 First of all, it is worth noting that the conclusions drawn from the experiments considering  
2151 all OCL assertions and mutants (cf. Sections 4.3 and 4.4) hold for this study with the 104 MT  
2152 mutants, i.e., *Mountford*, *Kulczynski2*, *Ochiai* and *Zoltar* have again the best numbers. Regarding  
2153 the static technique proposed by Burgueño et al. [18], it performs worse than these techniques. In  
2154 the average-case scenario, the static approach needs to inspect around 35% of the rules in order to  
2155



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 (out of 9) more rules in *BibTex2DocBook* (24.1% of the MT), 0.916 (out of 19) more rules in *CPL2SPL* (4.82% of the MT), 5 (out of 39) more rules in *Ecore2Maude* (12.8% of the MT), and 1.58 (out 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 in order 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 –*Kulczynski2*, *Mountford*, *Zoltar* and *Ochiai*– are

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	1400	18 / 16	90
BibTeX2DocBook	100	27 / 16	2700	40 / 40	269
CPL2SPL	100	34 / 15	3400	50 / 39	150
Ecore2Maude	100	42 / 3	4200	50 / 9	155
Total	400	117 / 44	11700	158 / 104	664

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
2206															
2207															
2208															
2209															
2210															
2211															
2212															
2213															
2214															
2215															
2216															
2217															
2218															
2219															
2220															
2221															
2222															
2223															
2224															
2225															
2226															
2227															
2228															
2229															
2230															
2231															
2232															
2233															
2234															
2235															
2236															
2237															
2238															
2239															
2240															
2241															
2242															
2243															
2244															
2245															
2246															
2247															
2248															
2249															
2250															
2251															
2252															
2253															
2254															

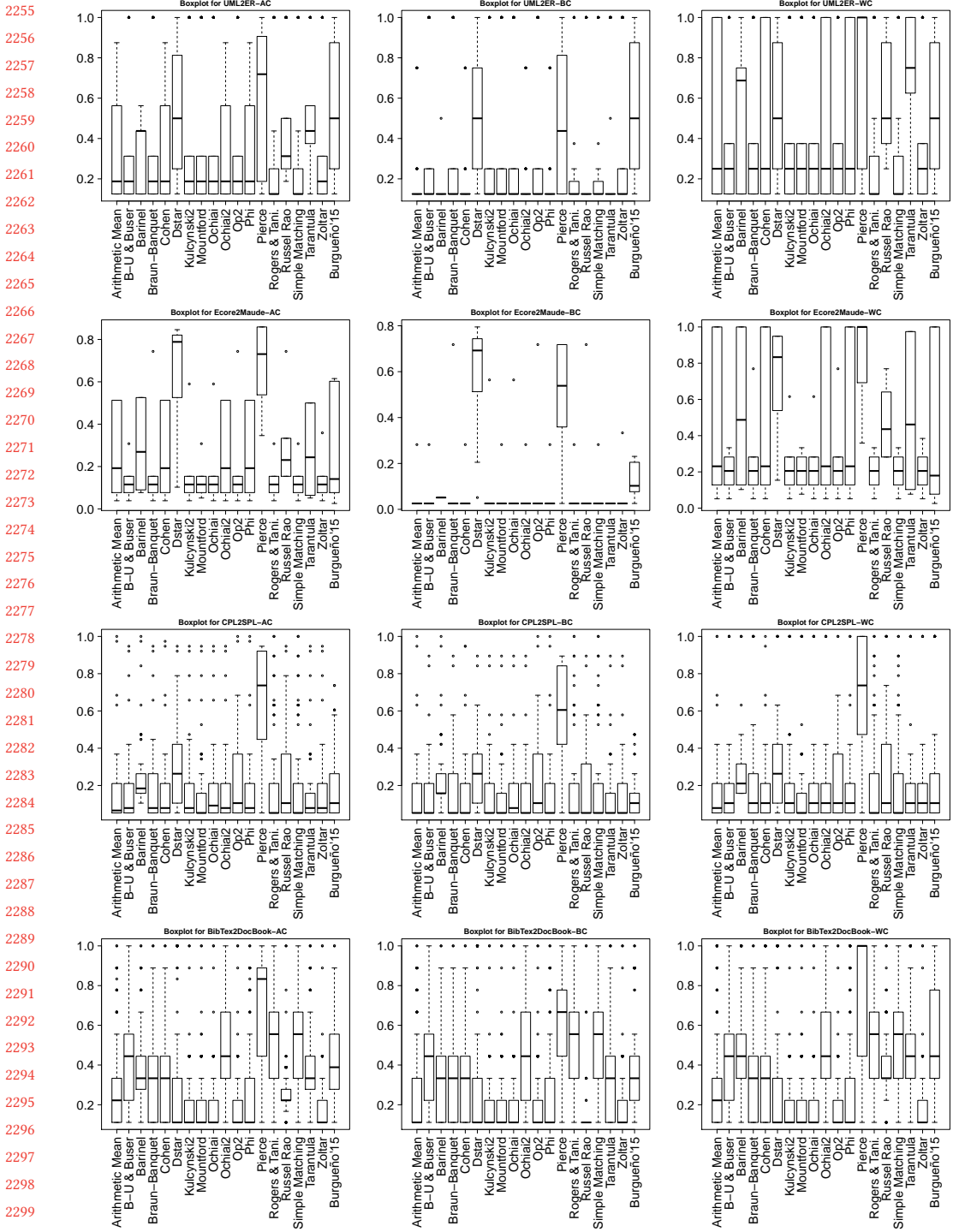


Fig. 8. Box-plot of the EXAM score of each technique per scenario and case study including [18] (Burgueno'15)

2304 always above 0.55, leaving no doubt that the static approach behaves worse. Besides, all these  
2305 cells reveal statistical differences (p-value <0.05, displayed in boldface in the table), except for the  
2306 *Ecore2Maude* case study. The latter is due to the fact that the results in *Ecore2Maude* have been  
2307 taken from only 9 mutants (cf. second part of fifth column in Table 9), which are the ones that make  
2308 the 3 OCL assertions considered in this case study fail, since only these assertions are defined in  
2309 the evaluation of the static approach by Burgueño et al. (cf. [18] –second part of third column in  
2310 Table 9). Indeed, in the comparison with the complete set of OCL assertions (cf. Section 4.5.3), the  
2311 cells of the *Ecore2Maude* also reveal statistical differences, since 42, instead of 3, OCL assertions are  
2312 considered. Please note that the conclusions of both comparisons is the same.

2313 The superiority of the static approach regarding *Pierce* is confirmed in the other three case  
2314 studies. However, it can not be concluded that it is better than any other of the techniques, since  
2315 the rows of the static technique do not present a value <0.5 in more than one case study for any of  
2316 the other techniques. Finally, we see that in the *UML2ER* case study the static approach behaves  
2317 generally much worse than most techniques. An explanation can be that the static approach, based  
2318 on types matching, does not behave well in the presence of rule inheritance.

2319 In summary, we can confirm that all SBFL techniques have a better performance when locating  
2320 the faulty rule than the static technique, except for *Pierce*, where the static technique behaves clearly  
2321 better. Besides, the static approach normally presents more ties than the best dynamic techniques.

2322 Received July 2017; revised March 2009; accepted June 2009  
2323

2324  
2325  
2326  
2327  
2328  
2329  
2330  
2331  
2332  
2333  
2334  
2335  
2336  
2337  
2338  
2339  
2340  
2341  
2342  
2343  
2344  
2345  
2346  
2347  
2348  
2349  
2350  
2351  
2352

Table 11. Effect size estimations for the comparison with [18]

	Cohen	Braun-Banquet	Simple Match	Kulczynski2	Barinel	Arith. Mean	Mountford	Zollar	Ochiai	Phi	Op2	Russel Rao	B-U & Buser	Pierce	Ochiai2	Rogers & Tani.	Dstar	Tarantula	Burgueno [18]	
BibTex2DocBook	Arithmetic Mean	0,000	<b>0,296</b>	<b>0,411</b>	<b>0,318</b>	<b>0,368</b>	<b>0,532</b>	<b>0,641</b>	<b>0,573</b>	<b>0,611</b>	<b>0,271</b>	<b>0,636</b>	<b>0,535</b>	<b>0,139</b>	<b>0,249</b>	<b>0,393</b>	<b>0,249</b>	<b>0,292</b>	<b>0,640</b>	<b>0,300</b>
	Barinel	<b>0,704</b>	0,000	<b>0,647</b>	<b>0,488</b>	<b>0,578</b>	<b>0,739</b>	<b>0,844</b>	<b>0,802</b>	<b>0,827</b>	<b>0,414</b>	<b>0,840</b>	<b>0,730</b>	<b>0,214</b>	<b>0,335</b>	<b>0,749</b>	<b>0,335</b>	<b>0,495</b>	<b>0,844</b>	<b>0,494</b>
	Braun-Banquet	<b>0,589</b>	<b>0,353</b>	0,000	<b>0,362</b>	<b>0,439</b>	<b>0,631</b>	<b>0,742</b>	<b>0,683</b>	<b>0,717</b>	<b>0,295</b>	<b>0,738</b>	<b>0,622</b>	<b>0,150</b>	<b>0,250</b>	<b>0,565</b>	<b>0,250</b>	<b>0,348</b>	<b>0,742</b>	<b>0,355</b>
	B-U & Buser	<b>0,682</b>	<b>0,512</b>	<b>0,638</b>	0,000	<b>0,579</b>	<b>0,709</b>	<b>0,811</b>	<b>0,769</b>	<b>0,793</b>	<b>0,434</b>	<b>0,807</b>	<b>0,707</b>	<b>0,212</b>	<b>0,354</b>	<b>0,694</b>	<b>0,354</b>	<b>0,507</b>	<b>0,811</b>	<b>0,498</b>
	Cohen	<b>0,632</b>	<b>0,422</b>	<b>0,561</b>	<b>0,421</b>	0,000	<b>0,666</b>	<b>0,771</b>	<b>0,722</b>	<b>0,751</b>	<b>0,352</b>	<b>0,768</b>	<b>0,660</b>	<b>0,183</b>	<b>0,297</b>	<b>0,640</b>	<b>0,297</b>	<b>0,418</b>	<b>0,772</b>	<b>0,412</b>
	Dstar	<b>0,468</b>	<b>0,261</b>	<b>0,369</b>	<b>0,291</b>	<b>0,334</b>	0,000	<b>0,608</b>	<b>0,535</b>	<b>0,575</b>	<b>0,252</b>	<b>0,603</b>	<b>0,504</b>	<b>0,163</b>	<b>0,244</b>	<b>0,330</b>	<b>0,244</b>	<b>0,257</b>	<b>0,606</b>	<b>0,276</b>
	Kulczynski2	<b>0,359</b>	<b>0,156</b>	<b>0,258</b>	<b>0,189</b>	<b>0,229</b>	<b>0,392</b>	0,000	<b>0,419</b>	<b>0,464</b>	<b>0,155</b>	<b>0,495</b>	<b>0,396</b>	<b>0,077</b>	<b>0,153</b>	<b>0,191</b>	<b>0,153</b>	<b>0,153</b>	<b>0,498</b>	<b>0,173</b>
	Mountford	<b>0,427</b>	<b>0,198</b>	<b>0,317</b>	<b>0,231</b>	<b>0,278</b>	<b>0,465</b>	<b>0,581</b>	0,000	<b>0,545</b>	<b>0,189</b>	<b>0,575</b>	<b>0,466</b>	<b>0,095</b>	<b>0,180</b>	<b>0,276</b>	<b>0,180</b>	<b>0,194</b>	<b>0,579</b>	<b>0,215</b>
	Ochiai	<b>0,389</b>	<b>0,173</b>	<b>0,283</b>	<b>0,207</b>	<b>0,249</b>	<b>0,425</b>	<b>0,536</b>	<b>0,455</b>	0,000	<b>0,168</b>	<b>0,581</b>	<b>0,428</b>	<b>0,085</b>	<b>0,165</b>	<b>0,228</b>	<b>0,165</b>	<b>0,169</b>	<b>0,534</b>	<b>0,190</b>
	Ochiai2	<b>0,729</b>	<b>0,586</b>	<b>0,705</b>	<b>0,566</b>	<b>0,648</b>	<b>0,748</b>	<b>0,845</b>	<b>0,811</b>	<b>0,832</b>	0,000	<b>0,842</b>	<b>0,750</b>	<b>0,262</b>	<b>0,410</b>	<b>0,763</b>	<b>0,410</b>	<b>0,581</b>	<b>0,846</b>	<b>0,576</b>
	Op2	<b>0,364</b>	<b>0,160</b>	<b>0,262</b>	<b>0,193</b>	<b>0,232</b>	<b>0,397</b>	<b>0,505</b>	<b>0,425</b>	<b>0,469</b>	<b>0,158</b>	0,000	<b>0,401</b>	<b>0,078</b>	<b>0,156</b>	<b>0,197</b>	<b>0,156</b>	<b>0,157</b>	<b>0,503</b>	<b>0,175</b>
	Phi	<b>0,465</b>	<b>0,270</b>	<b>0,378</b>	<b>0,293</b>	<b>0,340</b>	<b>0,496</b>	<b>0,604</b>	<b>0,534</b>	<b>0,572</b>	<b>0,250</b>	<b>0,599</b>	<b>0,000</b>	<b>0,129</b>	<b>0,231</b>	<b>0,342</b>	<b>0,231</b>	<b>0,266</b>	<b>0,602</b>	<b>0,277</b>
	Pierce	<b>0,861</b>	<b>0,786</b>	<b>0,850</b>	<b>0,788</b>	<b>0,817</b>	<b>0,837</b>	<b>0,923</b>	<b>0,905</b>	<b>0,915</b>	<b>0,738</b>	<b>0,922</b>	<b>0,871</b>	0,000	<b>0,731</b>	<b>0,869</b>	<b>0,731</b>	<b>0,781</b>	<b>0,923</b>	<b>0,779</b>
	Rogers & Tanimoto	<b>0,751</b>	<b>0,665</b>	<b>0,750</b>	<b>0,646</b>	<b>0,703</b>	<b>0,756</b>	<b>0,847</b>	<b>0,820</b>	<b>0,835</b>	<b>0,590</b>	<b>0,844</b>	<b>0,769</b>	<b>0,269</b>	<b>0,000</b>	<b>0,758</b>	<b>0,500</b>	<b>0,661</b>	<b>0,847</b>	<b>0,628</b>
	Russel Rao	<b>0,607</b>	<b>0,251</b>	<b>0,435</b>	<b>0,306</b>	<b>0,360</b>	<b>0,670</b>	<b>0,809</b>	<b>0,724</b>	<b>0,772</b>	<b>0,237</b>	<b>0,803</b>	<b>0,658</b>	<b>0,131</b>	<b>0,242</b>	0,000	<b>0,242</b>	<b>0,246</b>	<b>0,808</b>	<b>0,278</b>
	Simple Matching	<b>0,751</b>	<b>0,665</b>	<b>0,750</b>	<b>0,646</b>	<b>0,703</b>	<b>0,756</b>	<b>0,847</b>	<b>0,820</b>	<b>0,835</b>	<b>0,590</b>	<b>0,844</b>	<b>0,769</b>	<b>0,269</b>	<b>0,500</b>	<b>0,758</b>	<b>0,000</b>	<b>0,661</b>	<b>0,847</b>	<b>0,628</b>
	Tarantula	<b>0,708</b>	<b>0,505</b>	<b>0,652</b>	<b>0,493</b>	<b>0,582</b>	<b>0,743</b>	<b>0,847</b>	<b>0,806</b>	<b>0,831</b>	<b>0,419</b>	<b>0,843</b>	<b>0,734</b>	<b>0,219</b>	<b>0,339</b>	<b>0,754</b>	<b>0,339</b>	0,000	<b>0,847</b>	<b>0,499</b>
Zollar	<b>0,360</b>	<b>0,156</b>	<b>0,258</b>	<b>0,189</b>	<b>0,228</b>	<b>0,394</b>	<b>0,502</b>	<b>0,421</b>	<b>0,466</b>	<b>0,154</b>	<b>0,497</b>	<b>0,398</b>	<b>0,077</b>	<b>0,153</b>	<b>0,192</b>	<b>0,153</b>	<b>0,153</b>	<b>0,000</b>	<b>0,172</b>	
Burgueno [18]	<b>0,700</b>	<b>0,506</b>	<b>0,645</b>	<b>0,509</b>	<b>0,588</b>	<b>0,724</b>	<b>0,827</b>	<b>0,785</b>	<b>0,810</b>	<b>0,424</b>	<b>0,825</b>	<b>0,723</b>	<b>0,221</b>	<b>0,372</b>	<b>0,722</b>	<b>0,372</b>	<b>0,501</b>	<b>0,828</b>	<b>0,000</b>	
CPL2SPL	Arithmetic Mean	0,000	<b>0,274</b>	<b>0,472</b>	<b>0,503</b>	<b>0,497</b>	<b>0,274</b>	<b>0,493</b>	<b>0,547</b>	<b>0,474</b>	<b>0,483</b>	<b>0,436</b>	<b>0,504</b>	<b>0,103</b>	<b>0,507</b>	<b>0,391</b>	<b>0,507</b>	<b>0,520</b>	<b>0,493</b>	<b>0,444</b>
	Barinel	<b>0,726</b>	0,000	<b>0,686</b>	<b>0,750</b>	<b>0,728</b>	<b>0,436</b>	<b>0,731</b>	<b>0,786</b>	<b>0,700</b>	<b>0,645</b>	<b>0,745</b>	<b>0,169</b>	<b>0,706</b>	<b>0,632</b>	<b>0,706</b>	<b>0,799</b>	<b>0,737</b>	<b>0,509</b>	<b>0,671</b>
	Braun-Banquet	<b>0,528</b>	<b>0,314</b>	0,000	<b>0,526</b>	<b>0,524</b>	<b>0,322</b>	<b>0,514</b>	<b>0,565</b>	<b>0,513</b>	<b>0,521</b>	<b>0,463</b>	<b>0,527</b>	<b>0,113</b>	<b>0,523</b>	<b>0,422</b>	<b>0,523</b>	<b>0,533</b>	<b>0,516</b>	<b>0,466</b>
	B-U & Buser	<b>0,497</b>	<b>0,250</b>	<b>0,474</b>	0,000	<b>0,493</b>	<b>0,277</b>	<b>0,492</b>	<b>0,546</b>	<b>0,479</b>	<b>0,488</b>	<b>0,441</b>	<b>0,500</b>	<b>0,101</b>	<b>0,505</b>	<b>0,396</b>	<b>0,505</b>	<b>0,517</b>	<b>0,493</b>	<b>0,443</b>
	Cohen	<b>0,503</b>	<b>0,272</b>	<b>0,476</b>	<b>0,507</b>	0,000	<b>0,274</b>	<b>0,497</b>	<b>0,552</b>	<b>0,477</b>	<b>0,486</b>	<b>0,440</b>	<b>0,507</b>	<b>0,102</b>	<b>0,511</b>	<b>0,395</b>	<b>0,511</b>	<b>0,525</b>	<b>0,497</b>	<b>0,447</b>
	Dstar	<b>0,726</b>	<b>0,564</b>	<b>0,678</b>	<b>0,723</b>	<b>0,726</b>	0,000	<b>0,706</b>	<b>0,752</b>	<b>0,707</b>	<b>0,714</b>	<b>0,642</b>	<b>0,728</b>	<b>0,184</b>	<b>0,686</b>	<b>0,622</b>	<b>0,686</b>	<b>0,737</b>	<b>0,709</b>	<b>0,657</b>
	Kulczynski2	<b>0,507</b>	<b>0,269</b>	<b>0,486</b>	<b>0,508</b>	<b>0,503</b>	<b>0,294</b>	0,000	<b>0,553</b>	<b>0,488</b>	<b>0,496</b>	<b>0,449</b>	<b>0,509</b>	<b>0,106</b>	<b>0,510</b>	<b>0,404</b>	<b>0,510</b>	<b>0,522</b>	<b>0,501</b>	<b>0,449</b>
	Mountford	<b>0,453</b>	<b>0,214</b>	<b>0,435</b>	<b>0,454</b>	<b>0,448</b>	<b>0,248</b>	<b>0,447</b>	0,000	<b>0,434</b>	<b>0,442</b>	<b>0,401</b>	<b>0,453</b>	<b>0,092</b>	<b>0,466</b>	<b>0,350</b>	<b>0,466</b>	<b>0,463</b>	<b>0,447</b>	<b>0,400</b>
	Ochiai	<b>0,526</b>	<b>0,300</b>	<b>0,487</b>	<b>0,512</b>	<b>0,523</b>	<b>0,293</b>	<b>0,512</b>	<b>0,566</b>	0,000	<b>0,508</b>	<b>0,454</b>	<b>0,526</b>	<b>0,105</b>	<b>0,526</b>	<b>0,412</b>	<b>0,526</b>	<b>0,536</b>	<b>0,513</b>	<b>0,462</b>
	Ochiai2	<b>0,517</b>	<b>0,291</b>	<b>0,479</b>	<b>0,512</b>	<b>0,514</b>	<b>0,286</b>	<b>0,504</b>	<b>0,558</b>	<b>0,492</b>	0,000	<b>0,446</b>	<b>0,517</b>	<b>0,103</b>	<b>0,519</b>	<b>0,404</b>	<b>0,519</b>	<b>0,526</b>	<b>0,504</b>	<b>0,454</b>
	Op2	<b>0,564</b>	<b>0,355</b>	<b>0,537</b>	<b>0,559</b>	<b>0,560</b>	<b>0,358</b>	<b>0,551</b>	<b>0,599</b>	<b>0,546</b>	<b>0,554</b>	<b>0,000</b>	<b>0,563</b>	<b>0,129</b>	<b>0,548</b>	<b>0,459</b>	<b>0,548</b>	<b>0,565</b>	<b>0,552</b>	<b>0,496</b>
	Phi	<b>0,496</b>	<b>0,255</b>	<b>0,473</b>	<b>0,500</b>	<b>0,493</b>	<b>0,272</b>	<b>0,491</b>	<b>0,547</b>	<b>0,474</b>	<b>0,483</b>	<b>0,437</b>	<b>0,000</b>	<b>0,102</b>	<b>0,506</b>	<b>0,393</b>	<b>0,506</b>	<b>0,516</b>	<b>0,491</b>	<b>0,441</b>
	Pierce	<b>0,897</b>	<b>0,831</b>	<b>0,887</b>	<b>0,899</b>	<b>0,898</b>	<b>0,816</b>	<b>0,894</b>	<b>0,908</b>	<b>0,895</b>	<b>0,897</b>	<b>0,871</b>	<b>0,898</b>	0,000	<b>0,853</b>	<b>0,863</b>	<b>0,853</b>	<b>0,903</b>	<b>0,897</b>	<b>0,877</b>
	Rogers & Tanimoto	<b>0,493</b>	<b>0,294</b>	<b>0,477</b>	<b>0,495</b>	<b>0,489</b>	<b>0,314</b>	<b>0,490</b>	<b>0,534</b>	<b>0,474</b>	<b>0,481</b>	<b>0,452</b>	<b>0,494</b>	<b>0,147</b>	0,000	<b>0,401</b>	<b>0,500</b>	<b>0,503</b>	<b>0,490</b>	<b>0,453</b>
	Russel Rao	<b>0,609</b>	<b>0,368</b>	<b>0,578</b>	<b>0,604</b>	<b>0,605</b>	<b>0,378</b>	<b>0,596</b>	<b>0,650</b>	<b>0,588</b>	<b>0,596</b>	<b>0,541</b>	<b>0,607</b>	<b>0,137</b>	<b>0,599</b>	0,000	<b>0,599</b>	<b>0,613</b>	<b>0,598</b>	<b>0,537</b>
	Simple Matching	<b>0,493</b>	<b>0,294</b>	<b>0,477</b>	<b>0,495</b>	<b>0,489</b>	<b>0,314</b>	<b>0,490</b>	<b>0,534</b>	<b>0,474</b>	<b>0,481</b>	<b>0,452</b>	<b>0,494</b>	<b>0,147</b>	<b>0,500</b>	<b>0,401</b>	<b>0,000</b>	<b>0,503</b>	<b>0,490</b>	<b>0,453</b>
	Tarantula	<b>0,480</b>	<b>0,201</b>	<b>0,467</b>	<b>0,483</b>	<b>0,475</b>	<b>0,263</b>	<b>0,478</b>	<b>0,537</b>	<b>0,464</b>	<b>0,474</b>	<b>0,435</b>	<b>0,484</b>	<b>0,097</b>	<b>0,497</b>	<b>0,387</b>	<b>0,497</b>	<b>0,000</b>	<b>0,478</b>	<b>0,427</b>
Zollar	<b>0,507</b>	<b>0,268</b>	<b>0,484</b>	<b>0,507</b>	<b>0,503</b>	<b>0,291</b>	<b>0,499</b>	<b>0,553</b>	<b>0,487</b>	<b>0,496</b>	<b>0,448</b>	<b>0,509</b>	<b>0,103</b>	<b>0,510</b>	<b>0,402</b>	<b>0,510</b>	<b>0,522</b>	<b>0,000</b>	<b>0,448</b>	
Burgueno [18]	<b>0,556</b>	<b>0,329</b>	<b>0,534</b>	<b>0,557</b>	<b>0,553</b>	<b>0,343</b>	<b>0,551</b>	<b>0,600</b>	<b>0,538</b>	<b>0,546</b>	<b>0,504</b>	<b>0,559</b>	<b>0,123</b>	<b>0,547</b>	<b>0,463</b>	<b>0,547</b>	<b>0,573</b>	<b>0,552</b>	<b>0,000</b>	
ECORE2MAUDE	Arithmetic Mean	0,000	<b>0,330</b>	<b>0,575</b>	<b>0,620</b>	<b>0,500</b>	<b>0,100</b>	<b>0,575</b>	<b>0,600</b>	<b>0,575</b>	<b>0,500</b>	<b>0,575</b>	<b>0,500</b>	<b>0,080</b>	<b>0,620</b>	<b>0,430</b>	<b>0,620</b>	<b>0,550</b>	<b>0,615</b>	<b>0,480</b>
	Barinel	<b>0,670</b>	0,000	<b>0,750</b>	<b>0,790</b>	<b>0,670</b>	<b>0,160</b>	<b>0,750</b>	<b>0,790</b>	<b>0,750</b>	<b>0,670</b>	<b>0,750</b>	<b>0,670</b>	<b>0,080</b>	<b>0,790</b>	<b>0,570</b>				