

Demystifying the Adoption of Behavior-Driven Development in Open Source Projects

Fiorella Zampetti, Andrea Di Sorbo, Corrado Aaron Visaggio, Gerardo Canfora,
Massimiliano Di Penta*

Department of Engineering, University of Sannio, Italy

Abstract

Context: Behavior-Driven Development (BDD) features the capability, through appropriate domain-specific languages, of specifying acceptance test cases and making them executable. The availability of frameworks such as Cucumber or RSpec makes the application of BDD possible in practice. However, it is unclear to what extent developers use such frameworks, and whether they use them for actually performing BDD, or, instead, for other purposes such as unit testing.

Objective: In this paper, we conduct an empirical investigation about the use of BDD tools in open source, and how, when a BDD tool is in place, BDD specifications co-evolve with source code.

Method: Our investigation includes three different phases: (i) a large-scale analysis to understand the extent to which BDD frameworks are used in 50,000 popular open-source projects written in five programming languages; (ii) a study on the co-evolution of scenarios, fixtures and production code in a sample of 20 Ruby projects, through the Granger's causality test, and (iii) a survey with 31 developers to understand how they use BDD frameworks.

Results: Results of the study indicate that $\simeq 27\%$ of the sampled projects use BDD frameworks, with a prevalence in Ruby projects (68%). In about 37% of the cases, we found a co-evolution between scenarios/fixtures and production code. Specifically, changes to scenarios and fixtures often happen together or

*Corresponding author

Email address: dipenta@unisannio.it (Massimiliano Di Penta)

after changes to source code. Moreover, survey respondents indicate that, while they understand the intended purpose of BDD frameworks, most of them write tests while/after coding rather than strictly applying BDD.

Conclusions: Even if the BDD frameworks usage is widespread among open source projects, in many cases they are used for different purposes such as unit testing activities. This mainly happens because developers felt BDD remains quite effort-prone, and its application goes beyond the simple adoption of a BDD framework.

Keywords: Behavior-Driven Development, Acceptance Testing, Empirical Study, Co-Evolution

1. Introduction

Behavior-Driven Development (BDD) entails the creation and usage of executable scenarios, written in natural language during the software analysis phases, to support software acceptance testing [1]. BDD is derived from Test-Driven Development (TDD) and, similarly to it, involves the creation of test cases before coding. As mentioned by North [2], BDD is based on the conjecture that: “*A story’s behavior is simply its acceptance criteria —if the system fulfills all the acceptance criteria, it is behaving correctly; if it does not, it is not.*” Differently from canonical testing approaches and TDD, BDD focuses more on “the specifications of the behavior of the class” (*i.e.*, specifications of the products’ behavior) rather than on the “unit tests of a class”. Also, TDD is enacted through test cases, written for example using frameworks like xUnit (*e.g.*, JUnit for Java), aimed at executing source code units. In BDD, test cases are scenarios—hereby also referred to as *specifications*— written using a natural language-based Domain Specific Language (DSL) in a form very similar to user stories. Such test cases describe the system’s behavior from the outside. To make them executable, some glue code — hereby referred to as *fixture* — must be written to link scenarios to production code implementing the functionality to be tested.

20 The introduction of frameworks such as BEHAVE¹, CUCUMBER², JASMINE³,
21 JBEHAVE⁴, MOCHA⁵, or RSPEC⁶ has made the use of BDD possible in practice.
22 However, the level of adoption of BDD is still unclear. Also, it is unclear whether
23 the aforementioned frameworks are actually used for BDD, or whether they are
24 used for different purposes.

25 Researchers have previously studied the co-evolution of tests and code [3],
26 as well as the extent to which developers are adopting TDD properly [4]. Also,
27 early work indicated that the availability of acceptance tests (written in form
28 of Framework for Integrated Test – Fit tables [5]) helps developers in better
29 understanding requirements [6, 7, 8, 9] and improving the correctness of the
30 developed programs [10]. However, to the best of our knowledge, there are no
31 studies investigating the extent and the way BDD-related technology is being
32 used in practice. Previous research mainly investigated the (i) strengths and
33 limitations of the BDD tool support [1], (ii) the specific maintenance challenges
34 deriving from the BDD adoption [11, 12], and (iii) possible solutions to reduce
35 human efforts in maintaining specification files [13]. In summary, a question left
36 unanswered by previous work is:

37 *To what extent are BDD tools used, and, when they are used, do*
38 *developers leverage them for actually performing BDD?*

39 The main motivation of this empirical study is to understand whether there
40 is a gap between a technological hype — in this case, the idea, and its imple-
41 mentation in several frameworks, or making scenarios executable — and its
42 adoption. That is, while BDD technology is out of there, and while previous
43 research highlighted the benefits of such a technology, also in the context of
44 acceptance testing [6, 7, 8, 9], the state-of-the-practice may provide a different

¹<https://github.com/behav/behave>

²<https://cucumber.io>

³<https://jasmine.github.io/>

⁴<https://jbehave.org>

⁵<https://mochajs.org>

⁶<http://rspec.info>

45 picture. First, it is unclear the extent to which projects regularly adopt BDD
46 tools and, when they do, whether this adoption is confined to a few developers,
47 or is quite broad across the projects. Second, using BDD frameworks may or
48 may not mean “scenario first” and not even “test first”, because in the end such
49 frameworks could be used as regular (unit) testing tools, and do not prevent
50 developers for writing/modifying tests after production code. Last, but not least,
51 developers may perceive benefits in using BDD tools, but may also perceive
52 challenges in properly adopting them. In summary, understanding the *extent*
53 and the *way* BDD tools are adopted and what benefits and challenges such
54 adoption brings, requires a thorough empirical investigation.

55 To the best of our knowledge, this paper reports the first study aimed at
56 investigating the adoption of BDD in practice. Following the motivations stated
57 above, the study is carried out along three different dimensions, reflecting the
58 three study research questions detailed in Section 3, and summarized in Figure 2.
59 Namely, we study (i) the extent to which BDD tools are used by a large pool
60 of 50,000 open source projects and by their developers, (ii) how developers
61 change BDD scenarios, tests, and production code over time, and (iii) what is
62 the developers’ perception about BDD practices and technology. Note that the
63 study is limited to the open-source because of the large availability of projects’
64 data.

65 Results of the study indicate that BDD tools are used in 27.25% of 50,000
66 studied open-source projects developed in different programming languages, and
67 that the percentage is high for Ruby projects (68%), and somewhat substantial
68 for Javascript ($\simeq 39\%$) and Python ($\simeq 26\%$), while it is below 3% for the other
69 languages. Except for Ruby, and in particular of developers using the RSpec
70 framework, a very small percentage of developers ($< 7\%$) change BDD-related
71 files. An in-depth analysis of a sample of 20 Ruby projects indicates that
72 scenarios/source code co-evolve for about 40% of the analyzed pairs. In about
73 37% of the cases, we found a co-evolution between BDD scenarios (and related
74 test fixtures) and production code traced onto them. Finally, 31 respondents of a
75 survey we conducted indicate that, while they understood the intended purpose

76 of BDD tools, most of them write tests while/after coding rather than strictly
77 applying BDD.

78 **Paper structure.** Section 2.1 provides background notions about BDD
79 frameworks, explaining in particular how RSPEC works, while Section 2.2 dis-
80 cusses the related literature. Section 3 describes the study definition and planning.
81 Results are reported and discussed in Section 4, while the threats to its valid-
82 ity are discussed in Section 5. Section 6 discusses the study implications for
83 researchers, practitioners, and educators, while Section 7 concludes the paper
84 and outlines directions for future work.

85 **The study dataset** is available for replication purposes [14].

86 2. Background Notions and Related Work

87 This section first provides some background notions about Behavior-Driven
88 Development (BDD). Then, it overviews related literature.

89 2.1. An overview of Behavior-Driven Development with RSPEC

90 Similarly to Test-Driven Development (TDD), Behavior-Driven Development
91 (BDD) is a test-first approach. In particular, both techniques leverage the as-
92 sumption that generating test cases before implementing the software itself could
93 help test engineers and developers in finding inconsistencies in the requirements
94 specification and design documents earlier in the development process [15]. In
95 BDD, the software design flow starts from the design and the implementation
96 of the acceptance tests [16]. Such tests represent the means for communicating
97 between software designers, developers and the other (non-technical) stakehold-
98 ers. Indeed, for enabling the communication between the different types of
99 actors, BDD augments TDD by forcing the use of natural language for describing
100 automated acceptance tests. Acceptance tests collectively specify the user be-
101 haviors that the system has to fulfill [17] through the definition of a collection of
102 *scenarios* and natural language represents the *ubiquitous language* [1] to describe
103 such scenarios. Each scenario is described by defining at least three elements:

104 (i) a starting position (*e.g.*, navigate to a specific interface), (ii) the execution of
105 a user action (*e.g.*, clicking a button), and (iii) the assertion to verify whether
106 the consequent action has produced the expected effect.

107 Pragmatically speaking, BDD puts together the specification of an outcome
108 relevant for an end-user (*e.g.*, role-feature matrix) and the examples and/or
109 scenarios expressed with a single notation *given-when-then* implemented in DSLs
110 (for example, CUCUMBER⁷ uses the Gherkin DSL). That is, *given* some initial
111 contexts, *when* an event occurs, *then* it ensures some outcomes. The latter
112 (i) makes the behavior descriptions readable and easily accessible to domain
113 experts, testers and developers, (ii) allows developers to associate the scenarios
114 to the production code, and (iii) makes them executable (*i.e.*, provided that the
115 acceptance criteria are executable). This is achieved through some glue code
116 (*fixtures*), typically written using the same notation used for xUnit test cases, that
117 links scenarios to production code under test. As a consequence of the test-first
118 principle, when efficiently maintained, BDD tests also represent an accurate and
119 up-to-date documentation source for the entire software system [18].

120 As described in the introduction, various BDD frameworks exist. Since for
121 addressing one of our research questions we inspected projects (written in Ruby
122 language) that use RSPEC for test automation, in the following it is convenient
123 to specifically focus on RSPEC, to briefly describe how it works, though its
124 metaphor is relatively similar to the other frameworks, with some differences
125 (*e.g.*, CUCUMBER keeps scenarios separated from fixtures).

126 RSPEC allows developers to encapsulate the behavior under test via the
127 **describe** block, to contextualize a set of tests by means of the **context** block,
128 while tests are written using the **it** block. Moreover, RSPEC provides a way to
129 set up test state using the **before** and **after** directives that will apply to all
130 the tests included in the **describe** block.

131 Figure 1 shows a very simple RSPEC example aimed at testing the function-
132 ality of the *Bowling* class containing the method *hit*. The **describe** keyword

⁷<https://cucumber.io>

```

describe Bowling do
  context "with no strikes or spares" do
    it "sums the pin count for each roll" do
      bowling = Bowling.new
      20.times { bowling.hit(4) }
      expect(bowling.score).to eq 80
    end
  end
end

class Bowling
  attr_reader :score

  def initialize
    @score = 0
  end

  def hit(pin_count)
    @score += pin_count
  end
end
end

```

Figure 1: RSpec example.

133 takes the class name and requires a block argument containing the individual
 134 tests (*i.e.*, examples). The block argument is designated by the Ruby `do/end`
 135 keywords. The `context` provides a common setup to multiple tests, *e.g.*, those of
 136 the *Bowling* class in our case. After that, it is required to specify the `it` keyword
 137 aimed at defining a specific example that takes a string argument highlighting
 138 the expected behavior. The block argument belonging to the `it` clause contains
 139 a concrete specification of the test. Specifically, after having instantiated an
 140 object belonging to the *Bowling* class (*bowling*), the *hit* method is invoked 20
 141 times. Finally, the `expect` statement defines the “expectation” checking whether
 142 a specific expected condition is satisfied. In our case, the code expects that the
 143 *score* attribute is equal to 80.

144 2.2. Studies on Behavior-Driven Development

145 While the characteristics of BDD, the related frameworks, and the strengths
 146 and limitations of the existing tool support have been first investigated by Solis

147 and Wang [1], to the best of our knowledge, no prior work studied how developers
148 of open-source projects use BDD frameworks and the extent to which changes
149 on specifications and the related production code might influence each other.

150 Recent research efforts have been devoted to the investigation of the specific
151 challenges deriving from the BDD adoption in software projects. For instance,
152 Binamungu *et al.* [11], by surveying 75 practitioners, investigated the challenges
153 that may arise in maintaining the BDD specifications. They found that the
154 management of BDD specifications over the long term can be challenging, as
155 different team members could update such specifications multiple times and this
156 could result in redundant specifications that are more costly to maintain. For
157 this reason, in a different work, Binamungu *et al.* [19] presented an automated
158 approach for detecting duplicate examples in BDD specifications. Rahman
159 and Gao [12], instead, highlighted some of the maintenance issues that teams
160 might face when adapting BDD tests to changing environments and proposed a
161 reusable automated acceptance testing architecture to address reusability, and
162 maintainability concerns.

163 Finally, Yang *et al.* [13] proposed a text-based approach to recommend
164 when specification files should be modified, as a consequence of changes in the
165 production code, for reducing the developers' effort in maintaining up-to-date
166 specifications. Differently from Yang *et al.*, we rely on structural dependencies
167 between spec files and production code, as explained in Section 3.1.

168 2.3. Studies on Acceptance Testing

169 Research on acceptance testing is relevant to our work because BDD tools
170 can be used to automate acceptance testing, using test cases derived from use
171 case scenarios or user stories.

172 Several researchers have studied the use of executable acceptance tests in
173 software projects and in particular the use of Fit [5]. Melnik *et al.* [7] investigated,
174 through controlled experiments, the use of Fit user acceptance tests for specifying
175 functional requirements, highlighting that the use of Fit tables and the possibility
176 to execute them improve the comprehension of requirements. In a follow-up

177 study, Melnik *et al.* [6] investigated whether acceptance tests can be authored
178 effectively by customers of agile projects. Their results did not support the
179 hypothesis that the quality of executable acceptance tests is positively correlated
180 with the quality of production code. At the same time, results indicated how
181 acceptance tests help customers to specify functional requirements.

182 From a different perspective, Ricca *et al.* [9] assessed the impact of Fit tables
183 on the clarity of requirements, but also, in a different work, Ricca *et al.* [10]
184 highlighted how the use of Fit tables during software development helps to
185 improve the correctness of the produced code, without significantly affecting the
186 productivity. Results consistent to those of Ricca *et al.* were also obtained by
187 Park and Maurer [8] in a case study with custom-built executable acceptance
188 test cases. Specifically, their study showed how acceptance test cases are useful
189 to communicate requirements and to improve developers' productivity.

190 Using automated acceptance testing can be effort-prone for testers, as also
191 pointed out by respondents to our survey. A study by Monden *et al.* [20]
192 investigated, through a simulation model, the cost-effectiveness of testing (and
193 in particular acceptance testing) activities, indicating that the testing effort can
194 be reduced of 25% if accurate defect prediction is being used.

195 Our study substantially diverges from previous studies on acceptance testing.
196 While we do not specifically look for benefits of executable acceptance testing,
197 our study shows that since the early adoption of executable acceptance testing
198 (described in the aforementioned related work), the available technology has
199 gained maturity and adoption. At the same time, it is not only used to realize
200 the BDD paradigm, but also as a support for traditional unit testing. Our result
201 is consistent with results of previous studies showing how acceptance tests can
202 help improving productivity [8] and even code correctness [10].

203 *2.4. Studies on the Co-evolution of testing and production code*

204 Zaidman *et al.* [3] studied the co-evolution of production and test code, by
205 mining software repositories to better understand the testing processes that are
206 followed in practice. In a subsequent study, Zaidman *et al.* [21] experimented

207 the use of three views, namely *Change History View*, *Growth History View* and
208 *Test Coverage Evolution View* in two open-source and one industrial software
209 project, to study how test and production code co-evolve over time. Their results
210 indicated that, through the analysis of these complementary views, it is possible
211 to recognize different co-evolution scenarios (*i.e.*, synchronous and phased).

212 Lubsen *et al.* [22] presented an approach based on association rule mining
213 to determine whether production and test code co-evolve synchronously, and
214 evaluated the proposed approach on an open-source and an industrial project,
215 showing its effectiveness in recognizing the testing approach followed by software
216 developers. Association rule mining has been also exploited by Marsavina
217 *et al.* [23], who investigated co-evolution patterns of production and test code,
218 and identified six patterns reflecting the presence (or the lack of) co-evolution.

219 Beller *et al.* [4] monitored the development activity of 416 software engineers
220 over five months, observing that the majority of developers do not practice testing
221 actively. Moreover, they evaluated the co-evolution between test and production
222 code to demonstrate that TDD is not widely practiced, and developers who
223 claim to apply it in many cases do not properly follow TDD principles.

224 Canfora *et al.* [24] performed a TDD experiment with professionals, finding
225 that while TDD improves the quality of tests, it significantly slows down the
226 whole testing process.

227 Our study on BDD adoption is inspired by previous research. While we share
228 with Beller *et al.* [4] the somewhat negative results about the proper application
229 of BDD (TDD in their case), our study is focused on the co-evolution between
230 BDD specifications and production code. Also, and differently from previous
231 work, we conduct an analysis in the large for assessing the extent to which
232 BDD frameworks are used in open-source projects and, from a methodological
233 perspective, employ the Granger causality test [25] to assess the test-production
234 code pairs that statistically exhibit co-evolution patterns.

235 3. Study Design

236 The *goal* of this study is to investigate the adoption of BDD-related tech-
237 nology in open-source projects. The *perspective* is of researchers interested in
238 understanding whether and the extent to which BDD is adopted in open-source
239 software development, and whether BDD-related technology is used for different
240 purposes. The study *context* consists of open-source projects hosted on GitHub
241 (*i.e.*, 50,000 for the first study and a sample of 20 Ruby projects for the second
242 study). Moreover (in RQ₃), to empirically validate the achieved results, we con-
243 duct a survey — collecting a total of 31 responses — in which we ask developers
244 the kinds of BDD frameworks they use, and the main purpose for using them,
245 *i.e.*, performing BDD, or rather using the frameworks for TDD or generally for
246 other testing purposes.

247 Following the general goal of this study, and also what was discussed in
248 the introduction, to understand BDD practices, we first need to investigate the
249 precondition to any possible BDD practice, *i.e.*, whether BDD tools are used by
250 projects, or not. Also, while BDD tools are being used by projects, this may or
251 may not reflect a broad adoption of BDD practices by all developers (*i.e.*, only
252 some of them use BDD tools). Therefore, our first research question is:

253 **RQ₁** *To what extent do open source projects use BDD-related frame-*
254 *works?*

255 Of course, using a BDD framework does not necessarily mean that developers
256 of such projects were strictly following BDD. It is important to understand
257 whether BDD specifications are created before, together or after their related
258 production code, and how they co-evolve. To this aim, our second research
259 question is stated as follows:

260 **RQ₂** *How do RSpec specifications co-evolve with source code?*

261 The first two research questions provide a view of BDD practice adoption as it
262 is reflected from the projects' repositories. However, it does not tell much about

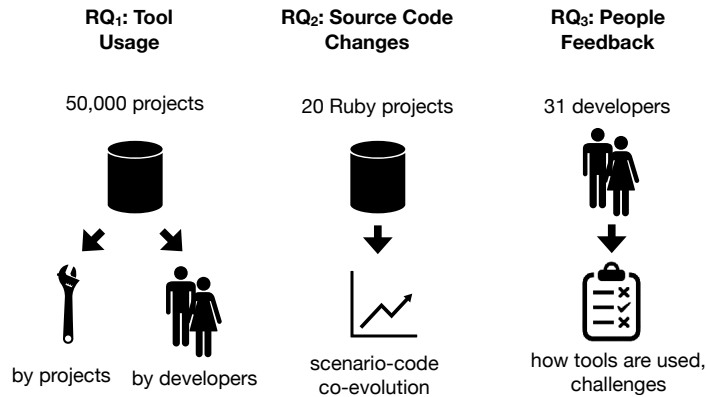


Figure 2: Study Overview.

263 developers’ perception of BDD usefulness, about the way they are using BDD
 264 tools, challenges they encounter, and limitations of the used tools as perceived
 265 by them. For these reasons, we would like to hear directly from developers about
 266 their direct experience in using BDD tools. To this aim, we ask our third and
 267 last research question:

268 **RQ₃** *How do developers use BDD-related frameworks?*

269 Figure 2 provides an overview of how the three research questions, together,
 270 contribute to the understanding of the BDD practice, analyzing the tool usage
 271 (RQ₁), changes to BDD-related resources and production code over time (RQ₂),
 272 and developers’ perception (RQ₃).

273 3.1. Context Selection and Analysis Methodology

274 In the following, for each research question, we describe the context selection
 275 and the analysis methodology. Note that the context slightly varies because, for
 276 instance, while in **RQ₁** we aim at obtaining a broad overview of the adoption of
 277 BDD tools, in **RQ₂** we study the phenomenon with a deeper degree of detail
 278 (on a fairly-limited number and kinds of projects we could analyze), and, finally,
 279 in **RQ₃** we are interested in getting feedback from developers.

280 To address **RQ₁**, we require a large set of projects, written in different
 281 programming languages, that may or may not use BDD tools. While performing

Table 1: BDD tools considered.

Name	Website	Supported languages
BEHAVE	https://github.com/behavetools/behavetools	Python
CUCUMBER	https://cucumber.io	All programming languages
JASMINE	https://jasmine.github.io/	Javascript
JBEHAVE	https://jbehave.org	Java
RSPEC	http://rspec.info	Ruby
MOCHA	https://mochajs.org	Javascript

282 such an analysis also for closed source is desirable, doing it on a large scale, and
 283 on a diverse set of projects (in terms of programming language, size, domain) is
 284 only feasible for open-source software projects. For this reason, we decided to
 285 conduct the study on open-source projects hosted on GitHub.

286 We analyze data from the top 50,000 projects — ranked in terms of the
 287 number of stars — hosted on GitHub for the five most popular programming
 288 languages on GitHub, *i.e.*, Java, Javascript, PHP, Python, and Ruby. As a
 289 further means to avoid toy projects, we only consider projects having at least
 290 one fork. Also, we exclude forked projects (to avoid duplicate data), and inactive
 291 projects (no changes after December 31st 2017). The selected projects have a
 292 median (i) size of 4,6 KLOC, (ii) number of stars of 677 and (iii) number of forks
 293 of 256.5. We considered 10,000 projects for each language, because a similar
 294 sample allows us to obtain high representativeness (*i.e.*, a margin of error less
 295 than 1% with a 95% confidence level) even for very large populations. Then,
 296 we identify, for each programming language, what are the BDD frameworks
 297 typically being used. This analysis has been performed by means of web searches.
 298 In the end, we identified six different tools described in Table 1, *i.e.*, BEHAVE,
 299 CUCUMBER, JASMINE, JBEHAVE, MOCHA, and RSPEC.

300 Clearly, it is possible that during our investigation, we could have missed
 301 some relevant BDD-related tools, as we based our analysis on those mentioned
 302 in previous research [1] and white papers about BDD. However, the selected
 303 tools are in-line with those used by our survey respondents (see RQ₃).

304 For each analyzed project, we look at imports/external dependencies into
305 source code or configuration files, to determine the adoption (or not) of one (or
306 more) of the identified BDD tools. In particular, as BDD frameworks need to
307 be imported as libraries into software projects to use their testing automation
308 features, we implemented a script able to (i) inspect the external dependencies
309 (that could be defined in configuration files, *e.g.*, `pom.xml`) and the import state-
310 ments directly occurring into source code files (*e.g.*, the statement `from behave`
311 `import` occurs in Python programs using BEHAVE), and (ii) automatically check
312 the BDD tools used.

313 We report, for each programming language, the number and percentage of
314 projects using at least one BDD-related tool, as well as the number of projects
315 using the considered tools (or even more than one tool).

316 After that, we identify the developers changing BDD-related files. To this
317 aim, we clone the repositories of all projects using at least one BDD framework,
318 and analyze their change history, including all branches while excluding merge
319 commits (to avoid analyzing the same change twice). Then, we look at (i) commit
320 authors and (ii) file paths being modified in each commit. From the latter, we
321 use regular expressions to determine whether a commit (and consequently its
322 author) has changed BDD-related files. For instance, for CUCUMBER and Java
323 projects we have looked at the main programming practices used to write testing
324 features and we have defined different regular expressions aimed at verifying
325 whether or not a file modified is a BDD-related file. Specifically, we check (i)
326 whether the file has extension `.feature`, or (ii) the file is in a `test` folder and
327 its name ends with `Steps.java`, or (iii) the file is a Java file inside a `steps` or
328 `step_definition` folder. Regarding RSPEC, instead, we check (i) whether the
329 file is a Ruby file inside a `spec` or `test` folder, and (ii) whether its name contains
330 the `spec` keywords. Since an author could use different emails (or IDs) within a
331 project, we use a simple heuristic to unify them: if two entries have the same git
332 ID, but different email, or the same email but different IDs, they belong to the
333 same identity. As a result, we report, for each programming language and for
334 each BDD tool, the average percentage of developers that change BDD-related

335 files.

336 To address **RQ**₂, we look at how specifications and production code co-evolve
337 to determine the extent to which tools are actually used to perform BDD. Since
338 this is a deeper analysis (which is also technology- and language-dependent) we
339 decided to focus on a sample of Ruby projects only, based on the results of **RQ**₁
340 (see Section 4.1), which show that BDD is more adopted in Ruby than in other
341 languages, and that `RSpec` is by large the most adopted BDD tool for Ruby.

342 We therefore randomly select a sample of 20 Ruby projects (detailed in Ta-
343 ble 2), using `RSpec` and Travis-CI as Continuous Integration (CI) platform. For
344 each project, we report the number of commits, commit authors and specification
345 files. The reason behind the inclusion of projects using Travis-CI is twofold.
346 On the one side, we wanted to guarantee that developers use BDD-frameworks
347 to fail the build process in the presence of test failures. On the other side, as
348 reported by North [2], BDD has grown to encompass the wider picture of agile
349 analysis and automated acceptance testing. However, it is important to note
350 that the use of CI does not necessarily imply a co-evolution of BDD scenarios
351 and test fixtures. Indeed, a BDD scenario could be added before it is linked to
352 production code through a fixture. At the same time, a developer could add
353 a new feature in the production code without having created any tests for it.
354 Finally, the only case for which we expect that test and production code always
355 co-change is when a feature changes its behavior or a bug is fixed.

356 The selected projects (see Table 2) cover different application domains, from
357 frameworks for building and releasing (*e.g.*, `capistrano` and `fastlane`) to support
358 frameworks for testing activities (*e.g.*, `capybara` and `vcr`) to utilities frameworks
359 like `httparty`, a HTTP library for Ruby, or `spree`, an e-commerce solution for
360 Ruby on Rails. Moreover, the involved projects are different in terms of the
361 number of specification files varying from 7 for `middleman` to 1,230 for `ruby`.

362 Due to the specificity of the analysis, we intendedly focus on a fairly limited
363 set of projects, also for showing detailed results obtained in each studied project.
364 Also, as Table 2 shows, the analysis involves nearly 200k commits, over 3,600 spec
365 files and, as reported in Table 7 and Table 8, 8,780 spec-code pairs. Indeed, for

Table 2: Characteristics of projects selected to address **RQ₂**.

Repository	Commits	Contrib.	Spec Files	Files (*.rb)
Homebrew/brew	17,513	634	318	1397
capistrano/capistrano	1,624	217	30	86
teamcapybara/capybara	3,420	248	41	186
carrierwaveuploader/carrierwave	2,214	309	29	84
freeCodeCamp/devdocs	2,748	82	28	546
diaspora/diaspora	19,737	341	228	753
discourse/discourse	30,808	664	382	2661
thoughtbot/factory_bot	990	175	42	124
fastlane/fastlane	14,794	905	316	901
guard/guard	2,396	151	52	119
jnunemaker/httparty	1,153	169	14	62
huginn/huginn	3,072	168	120	419
middleman/middleman	3,294	181	7	221
thoughtbot/paperclip	1,966	376	48	143
rspec/rspec-rails	2,650	263	34	89
ruby/ruby	53,662	51	1,230	7,389
sinatra/sinatra	4,032	362	31	139
spree/spree	18,932	763	317	1,220
hashicorp/vagrant	11,746	866	344	1,216
vcr/vcr	1,997	107	27	89
Min	990	51	7	62
Median	3,183	256	45	204
Avg	9,937	352	182	892
Max	53,662	905	1,230	7,389
Total	198,748	7,032	3,638	17,844

366 each project, we first clone the repository and extract the change history log for
367 each branch, excluding merge commits. After that, we build a traceability matrix
368 looking at the dependencies between test fixture (*i.e.*, spec files) and production
369 code. To this aim, we exploit the (i) Ruby library structure⁸, and (ii) the RSpec
370 naming conventions⁹ to statically extract traceability links between “spec” files
371 and production code. It is worth to point out that a manual inspection of the
372 selected projects allowed us to verify that each of such projects complies with
373 the aforementioned conventions. Specifically, for each project, we automatically
374 inspect all the spec (or test) files to collect:

- 375 • The relative path of the spec (or test) file: the module under test will
376 reside in a file having the same sub-path structure of the spec (or test) file
377 and the same name without the `__spec` (or `__test`) modifier; if such a file
378 exists in the project, a traceability link between the spec file and such a
379 file is added;
- 380 • The file names (and the relative paths) included in the require or load
381 clauses (*e.g.*, `require “path/filename.rb”`): if a similar file (along with the
382 specific path) exists in the project, a traceability link between the spec (or
383 test) file and such file is added;
- 384 • The defined namespaces: for a given namespace defined in the spec (or test)
385 file, *e.g.*, `A::B::Class`, we check if a file with the `.rb` extension and a similar
386 sub-path structure, *e.g.*, `a/b/class.rb`, exists in the project (also considering
387 snake case or camel case variants) and, in this case, a traceability link
388 between the spec (or test) file and such file is added.

389 Then, based on the availability of traceability matrices and the change history,
390 we identify, for each pair RSpec-production code files:

- 391 • Which one has been added into the system first, or whether they have
392 been added together. In particular, for each pair P_i we analyze the commit

⁸<https://ryanlue.com/posts/2017-03-03-how-to-structure-a-library>

⁹<https://relishapp.com/rspec/rspec-rails/v/3-0/docs/directory-structure>

393 history to establish whether the specification file, S_i , has been added to
394 the code base in the same commit (*i.e.*, *simultaneously*) or in a prior (*i.e.*,
395 *before*) or subsequent (*i.e.*, *after*) commit with respect to the commit in
396 which the production code file, C_i , has been added.

- 397 • When such files change, we look at the type of change made on the test
398 file distinguishing between (i) changes to specifications only (*i.e.*, only a
399 “describe” or “it” clauses has been modified), (ii) changes to fixture code only
400 (*i.e.*, modifying only the code leaving unchanged the behavior expressed
401 in the specification), and (iii) changes occurred to both specifications and
402 fixtures.

403 Based on the change history of the sampled projects, we report the number
404 and percentage of commit authors that modify BDD-related files (*i.e.*, following
405 the same approach used for RQ₁), and we verify as well how those values vary
406 while looking at those commits that modify both specifications and fixtures.

407 To (statistically) analyze the co-evolution of each files’ pair, we use the
408 Granger’s causality test [25], *i.e.*, a technique for determining whether a time
409 series is useful in forecasting another one. The basic idea is that a cause cannot
410 come after the effect implying that if a variable x affects a variable z , the former
411 variable should help to improve the predictions of the latter.

412 Let $fs_i(t)$ be the time series of the cumulative changes occurred to the
413 specification file fs_i at commit t . In other words, $fs_i(t) = 0$ if the file does
414 not exist at t , $fs_i(t) = 1$ if t adds it, and then increases of 1 for each change.
415 Similarly, let $fp_j(t)$ be the time series of a production code file traced onto
416 $fs_i(t)$. The simplest bivariate Granger test between two time series, $fs_i(t)$ and
417 $fp_j(t)$, uses the autoregressive specification [26], which consists of the following
418 univariate (restricted model):

$$fp_j(t) = c_1 + \gamma_1 fp_i(t - 1) + \gamma_2 fp_i(t - 2) + \dots + \gamma_p fp_i(t - p) + e(t) \quad (1)$$

419 and bivariate (unrestricted model):

$$\begin{aligned}
fp_j(t) &= c_1 + \gamma_1 fp_i(t-1) + \gamma_2 fp_i(t-2) + \dots \\
&+ \gamma_p fp_i(t-p) + \beta_1 fp_i(t-1) + \beta_2 fp_i(t-2) + \dots + \beta_p fp_i(t-p) + e(t) \quad (2)
\end{aligned}$$

420 solved by estimating the ordinary least squares. In the equations, p is the
421 lag length (*i.e.*, the delay, in terms of number of commits, between spec and
422 production code update, or vice versa), which can be estimated with various
423 criteria [26], and $e(t)$ is an independent random variable. To test whether $fs_j(t)$
424 Granger-causes $fp_j(t)$, the null hypothesis to reject (H_0 : changes to fs_i do not
425 Granger-cause changes to fp_j) is defined as:

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_p = 0$$

426 In our work, we use the Akaike Information Criterion (AIC) [27] to estimate
427 the lag p for each pair of time series (details about the obtained lags are in our
428 replication package [14]). The AIC rewards models that achieve high goodness of
429 fits, but penalizes the increase in the number of parameters, to avoid over-fitting.

430 In summary, for each pair (fs_i, fp_j) for which there exists a traceability
431 link, we test the Granger's hypothesis. Consider that in presence of transitive
432 dependencies (*i.e.*, multiple production files covered by the same specification
433 file) we look at each pair, instead of considering the cumulative changes among
434 them. Also, we test the Granger's causality in both directions. Since we have
435 to perform multiple Granger's tests, we adjust the obtained p -values using the
436 Benjamini-Hochberg procedure [28], which ranks p -values, keeping the largest
437 p -value unaltered, multiplying the second largest by the number of p -values
438 divided by the rank (*i.e.*, two), and treating the remaining ones similarly to the
439 second.

440 Granger's test has been performed using the *lmtest* R package [29], while the
441 AIC procedure used to estimate the lag belongs to the *vars* R package [30, 31].

442 Once having the time series, for each type of change we tested the Granger's
443 hypothesis with the time series belonging to the production code. In other

444 words, we test (in both directions) the causality of changes (i) between specifica-
445 tion+fixture changes and production code changes, and (ii) between fixture-only
446 changes and production code changes. Note that we did not consider the co-
447 evolution of specification-only changes with source code, as this happened in a
448 very limited number of cases (a median of zero times, mean=0.3, and max< 10
449 per specification file).

450 Finally, we use statistical procedures to determine (i) whether the co-evolution
451 between tests and production code occurs more when the whole specification
452 (description plus fixture) is changed or when only the fixture is changed, and
453 (ii) whether changes to production code follow changes to test (*i.e.*, BDD- or
454 TDD-like), or, instead, the other way around happens, or either they co-change.
455 More specifically, we compare the number of co-evolution instances using the
456 paired Wilcoxon rank-sign test [32] assuming a significance level $\alpha = 5\%$, and
457 using the Cliff's d effect size [33].

458 To address **RQ**₃, we conducted a survey questionnaire. Previous work
459 [34, 35, 36, 37, 38, 39] provided relevant guidelines on how to design and conduct
460 survey studies in software engineering. Also, there exist broader guidelines for
461 survey research in general [40].

462 Following these guidelines, we (i) designed *specific and measurable* goals, (ii)
463 targeted subjects *able to answer* the questions posed (*i.e.*, subjects having testing
464 experience), (iii) grouped questions into topics, and (iv) when possible, used
465 *standardized response formats*. However, we are unable to assess whether the set
466 of respondents to our survey is a *representative subset* of the target population,
467 as we do not have prior knowledge of such a population and collected feedback
468 from developers enrolled in online forums.

469 The questionnaire consists of two sections. The first section asks demograph-
470 ics' questions, *i.e.*, the respondent's maximum education degree, her experience
471 in programming, testing, and the programming languages being used.

472 The second section features four questions. The first one asks which BDD
473 frameworks are being used (if any). The second question asks whether they
474 write BDD specifications to model the system behavior before starting to code,

475 whether they are used in a test-first manner (like TDD), along with coding,
476 or after. The third question asks the main purpose of using BDD frameworks.
477 Finally, the last question is about further comments the respondent may have.

478 We created the questionnaire using a Google Form, and distributed it through
479 different channels, namely Twitter, LinkedIn, and Reddit channels about testing
480 and Ruby, namely *BDD*, *TDD*, *softwaretesting*, *test*, and *Ruby*. Moreover, as the
481 scraping of personal information for spamming purposes explicitly violates the
482 GitHub’s terms of service¹⁰, we have also looked for the development mailing
483 lists/forums related to the Ruby projects sampled for answering our **RQ₂** and
484 posted our questionnaire onto the mailing lists/forums related to 14 of these
485 projects, intending to also receive feedback from developers of such projects.
486 For the remaining 6 projects, we were not able to find any mentions of similar
487 development communication channels. 31 replies from respondents with (self-
488 rated) experience in testing activities from moderate to high were collected and
489 analyzed. It is worth to point out that all the respondents to our survey own some
490 college degree (Bachelor’s, Master’s, PhD or others) and the majority of them
491 (*i.e.*, 16) are professional developers having more than 5 years of programming
492 experience.

493 **4. Study Results**

494 In this section we report and discuss the main results obtained during our
495 exploratory investigation.

496 *4.1. To what extent do open source projects use BDD-related frameworks?*

497 Table 3 reports the number and percentage of projects using BDD frameworks.
498 Note that the “Overall” column is different from the sum of the other columns
499 since it indicates the number of projects using at least one BDD framework,
500 indeed there are projects using more than one BDD framework.

¹⁰<https://help.github.com/en/articles/github-terms-of-service#5-scraping>

Table 3: Usage of BDD tools in 50,000 projects (10,000 for each programming language) in GitHub.

Language	BEHAVE	CUCUMBER	JASMINE	JBEHAVE	MOCHA	RSPEC	Overall
Java	✗	78 0.78%	✗	57 0.57%	✗	✗	125 1.25%
Javascript	✗	39 0.39%	1,206 12.06%	✗	2,915 29.15%	✗	3,902 39.02%
Python	2,522 25.22%	108 1.08%	✗	✗	✗	✗	2,578 25.78%
PHP	✗	226 2.26%	✗	✗	✗	✗	226 2.26%
Ruby	✗	706 7.06%	✗	✗	✗	6,735 67.35%	6,794 67.94%
Overall	2,522 5.04%	1,157 2.31%	1,206 2.41%	57 0.11%	2,915 5.83%	6,735 13.47%	13,625 27.25%

501 As the table shows, the percentage is low for Java and PHP, 1.25% and
502 2.26% respectively. Instead, the percentage grows substantially for Python
503 projects (25.78%, mostly using BEHAVE), Javascript projects (\simeq 39%, mostly
504 using MOCHA) and, above all, for Ruby (67.94%). In this case, not only we
505 can notice a large majority of projects using RSPEC (67.35%), but, also, we
506 notice that 7.06% of the analyzed projects also use CUCUMBER (the percentage
507 of projects using CUCUMBER is 2.26% or below for the other programming
508 languages).

509 As mentioned in Section 2.1, RSPEC is a framework conceived for BDD,
510 but it can also be used as a generic xUnit framework for Ruby, although other
511 specific tools also exist for that purpose, see for example Minitest¹¹. A more
512 detailed analysis of how RSPEC is being used in the open-source community is
513 performed in **RQ**₂.

514 To deepen our analysis about BDD tool usage, we investigate how many
515 developers (commit authors) change BDD-related files. Table 4 reports, for

¹¹<https://github.com/seattlerb/minitest>

Table 4: Percentage of commit authors that changed BDD-related files.

Language	BEHAVE	CUCUMBER	JASMINE	JBEHAVE	MOCHA	RSPEC
Java	✗	2.62%	✗	1.76%	✗	✗
Javascript	✗	0.32%	5.4%	✗	2.77%	✗
Python	6.64%	4.91%	✗	✗	✗	✗
PHP	✗	4.55%	✗	✗	✗	✗
Ruby	✗	6.09%	✗	✗	✗	20.58%

516 each BDD-tool and for each language, the percentage of developers that modify
517 BDD-related files. Unsurprisingly, except for RSPEC, the percentage is always
518 lower than 7% meaning that, only a limited number of developers change BDD-
519 related files. The latter suggests that when a team introduces specific frameworks
520 for following BDD practices only a few developers have the responsibility to
521 create and modify them. For Ruby projects using RSPEC, the percentage is
522 substantially greater than for other languages and tools, and it reaches 20.58%.
523 We conjecture that this happens because, very often, Ruby developers tend to
524 use a specific BDD framework not only to follow BDD practices but rather for
525 other purposes, *e.g.*, to perform unit testing.

While the studied BDD frameworks are rarely used in Java and PHP projects, there is a substantial proportion of Python ($\simeq 26\%$) and Javascript projects ($\simeq 39\%$) using BDD frameworks, and above all a large majority of Ruby projects ($\simeq 68\%$) using RSPEC, but also CUCUMBER. Looking at developers, less than 7% of the total number of developers change BDD-related files, except for Ruby projects using RSPEC, where this percentage is of 20.58%.

527 4.2. How do RSPEC specifications co-evolve with source code?

528 Table 5 reports, for each analyzed project, the number (and percentage) of
529 times in which (i) a specification file has been added to the code base before the
530 related production code (referred as **SAB** in the table), (ii) a specification file and
531 its related production code file have been added to the code base simultaneously

Table 5: Specification files additions: Spec Added Before (SAB), Spec Added Simultaneously (SAS), Spec Added After (SAA).

Project name	SAB (#)	SAB (%)	SAS (#)	SAS (%)	SAA (#)	SAA (%)
brew	324	39.08	140	16.89	365	44.03
capistrano	2	3.64	22	40.00	31	56.36
capybara	10	18.18	9	16.67	35	64.81
carrierwave	6	10.53	10	17.54	41	71.93
devdocs	0	0	53	96.36	2	3.64
diaspora	65	19.58	148	44.58	119	35.84
discourse	101	18.57	236	43.38	207	38.05
factory_bot	13	20.31	47	73.44	4	6.25
fastlane	76	11.57	263	40.03	318	48.42
guard	55	24.12	48	21.05	125	54.82
httparty	11	37.93	10	34.48	8	27.59
huginn	5	3.52	92	64.79	45	31.69
middleman	0	0	5	45.45	6	54.55
paperclip	0	0	2	3.17	61	96.83
rspec-rails	6	14.63	7	17.07	28	68.29
ruby	588	14.04	220	5.25	3,381	80.71
spree	125	15.57	250	31.13	428	53.30
vagrant	33	6.16	121	22.57	382	71.27
vcr	0	0	0	0	54	100.00
Min	0	0	0	0	2	3.64
Median	10.50	12.81	47.50	32.81	49.50	54.69
Avg	71	12.87	84.85	33.58	283.15	53.53
Max	588	39.08	263	96.36	3,381	100

532 (SAS), and (iii) a specification file has been added to the code base after its
533 related production code (SAA).

534 Although using BDD frameworks (*i.e.*, RSpec), it seems that developers
535 often avoid putting into practice the test-first principle of BDD. Indeed, for all
536 the analyzed projects, the majority of the specification files are introduced into
537 the code base with or after the related production code. Moreover, 5 out of 20
538 projects (*carrierwave*, *middleman*, *paperclip*, *sinatra*, and *vcr*) have no cases in
539 which the specification files have been introduced before the related production
540 code. Looking deeper into these cases (*i.e.*, specifications introduced before the
541 production code), we found that this frequently occurs when (i) the production

542 code files are renamed¹²¹³, and (ii) the dependency between the specification
543 file and the production code file is established when new specifications or test
544 cases are introduced into the specification file, as a result of improvements in
545 the production code¹⁴¹⁵.

546 It is important to note that by “addition before” or “addition after” we mean
547 at any time in the observed project duration. Additions of either specification or
548 production code files may occur at different development stages. In particular,
549 a file added years after a spec file may have been linked to the specification
550 long after its creation. Similarly, specification files could have been introduced
551 long after the related production code, as RSpec has been added during the
552 project, and, consequently, not used for BDD purposes on the already existing
553 production code. In summary, the addition of specification and production code
554 is not necessarily aligned. Instead, as we will observe in the rest of the section,
555 it is more relevant to analyze how they co-evolve once having been added.

556 Based on the change analysis performed, and based on the results obtained
557 in **RQ₁**, we now refine (on the 20 Ruby projects) the analysis concerning the
558 percentage of developers that used BDD tools. Table 6 reports, for each project,
559 (i) the number of developers (unique commit authors in the observed history),
560 (ii) the number and percentage of those modifying RSpec files (third column, as
561 done in **RQ₁**), and, (iii) as a fine-grained analysis, the number and percentage
562 of those who, in doing so, modified RSpec scenarios (features) at least once
563 (fourth column).

564 Concerning the percentage of developers that change RSpec files, it varies in
565 the range [1.6%-48.4%]. While there are a few projects in which RSpec files are
566 modified by a limited number of developers (*i.e.*, only 5 projects have a percentage
567 < 15%), in all the other cases that percentage increases (*i.e.*, 15 projects have
568 a percentage > 15%). The fourth column in Table 6 reports the number and

¹²<https://github.com/jnunemaker/httparty/commit/4c6514f>

¹³<https://github.com/teamcapbara/capybara/commit/6b099a4>

¹⁴<https://github.com/jnunemaker/httparty/commit/d0d88fe>

¹⁵<https://github.com/teamcapbara/capybara/commit/db1bbd2>

Table 6: RQ₂: Percentage of commit authors that change the BDD-related files.

Project Name	# Dev.	Changing RSpec files	(%)	Changing scenarios	(%)
brew	828	51	(6.2%)	48	(5.8%)
capistrano	426	106	(24.9%)	52	(12.2%)
capybara	338	124	(36.7%)	62	(18.3%)
carrierwave	410	125	(30.5%)	119	(29.0%)
devdocs	128	2	(1.6%)	2	(1.6%)
diaspora	502	118	(23.5%)	118	(23.5%)
discourse	923	147	(15.9%)	147	(15.9%)
factory_bot	247	69	(27.9%)	25	(10.1%)
fastlane	1162	308	(26.5%)	194	(16.7%)
guard	186	55	(29.6%)	13	(7.0%)
httparty	215	104	(48.4%)	90	(41.9%)
huginn	213	47	(22.1%)	47	(22.1%)
middleman	221	6	(2.7%)	6	(2.7%)
paperclip	453	188	(41.5%)	60	(13.2%)
rspec-rails	322	57	(17.7%)	54	(16.8%)
ruby	276	17	(6.2%)	10	(3.6%)
sinatra	451	112	(24.8%)	32	(7.1%)
spree	1065	245	(23.0%)	227	(21.3%)
vagrant	1113	139	(12.5%)	139	(12.5%)
vcr	140	49	(28.6%)	10	(7.1%)

569 percentage of developers that changed scenarios at least once while modifying
570 RSpec tests. As the reader can notice, in 7 projects (*capistrano*, *capybara*,
571 *factory_bot*, *guard*, *paperclip*, *sinatra*, and *vcr*) the percentage substantially
572 drops (of 50% or more, compared to the third column). In other words, only
573 a minority of developers who changed test cases also modify the scenarios.
574 Therefore, the rest of the developers changing tests are using RSpec as a unit
575 testing tool. In the following, by looking at the co-evolution it is possible to gain
576 more insights into how developers use BDD-related frameworks.

577 Tables 7 and 8 summarize the results of the Granger’s Causality test related
578 to the co-evolution of specification files and their related production code. As
579 explained in Section 3, we split the analysis between changes that involved the
580 whole specification (Table 7) and changes to fixture only (Table 8). The tables

Table 7: Granger Causality Test: summary of results (changes to both scenarios and test code).

Proj.	Pairs	Exact Match	Code→ Spec	Spec→ Code	Both	Overall	(%)
brew	829	90	110	80	190	380	45.84
capistrano	55	4	12	11	16	39	70.91
capybara	54	0	12	9	4	25	46.30
carrierwave	57	2	8	13	6	27	47.37
devdocs	55	7	7	7	11	25	45.45
diaspora	332	8	99	48	96	243	73.19
discourse	544	17	123	81	129	333	61.21
factory_bot	64	6	21	11	8	40	62.50
fastlane	657	3	124	65	52	241	36.68
guard	228	0	20	41	1	62	27.19
httparty	29	1	5	3	5	13	44.83
huginn	142	8	32	12	67	111	78.17
middleman	11	2	3	1	3	7	63.64
paperclip	63	1	7	11	3	21	33.33
rspec-rails	41	2	9	6	9	24	58.54
ruby	4,189	69	620	122	202	944	22.54
sinatra	37	2	9	6	5	20	54.05
spree	803	13	102	91	157	350	43.59
vagrant	536	23	132	53	183	368	68.66
vcr	54	0	2	1	0	3	5.56
Min	11	0	2	1	0	3	5.56
Median	63.50	3.50	16	11.50	10	39.50	46.84
Avg	439	12.90	72.85	33.60	57.35	163.80	49.48
Max	4,189	90	620	122	202	944	78.17
Total	8,780	258	1,457	672	1,147	3,276	37.31

581 report, for each project (i) the number of overall pairs analyzed (**Pairs** column),
582 (ii) the amount of pairs exhibiting a perfect co-evolution (**Exact Match** column),
583 (iii) the number of statistically significant pairs for which changes on the specifi-
584 cation file occur after the changes in the production code (**Code→Spec** column),
585 (iv) the amount of statistically significant pairs for which changes on production

Table 8: Granger Causality Test: summary of results (changes to test code, scenarios unchanged).

Proj.	Pairs	Exact	Code→	Spec→	Both	Overall	(%)
		Match	Spec	Code			
brew	829	1	71	40	23	134	16.16
capistrano	55	0	10	8	9	27	49.09
capybara	54	0	16	3	4	23	42.59
carrierwave	57	0	25	5	3	33	57.89
devdocs	55	0	1	4	0	5	9.09
diaspora	332	0	93	33	70	196	59.04
discourse	544	0	97	71	63	231	42.46
factory_bot	64	0	7	2	0	9	14.06
fastlane	657	1	59	93	29	181	27.55
guard	228	0	1	2	0	3	1.32
httparty	29	0	4	1	1	6	20.69
huginn	142	0	31	24	17	72	50.70
middleman	11	0	1	0	0	1	9.09
paperclip	63	0	10	5	2	17	26.98
rspec-rails	41	1	5	3	1	9	21.95
ruby	4,189	179	479	7	185	671	16.02
sinatra	37	0	3	3	7	13	35.14
spree	803	61	117	83	165	365	45.45
vagrant	536	0	73	60	74	207	38.62
vcr	54	0	3	2	0	5	9.26
Min	11	0	1	0	0	1	1.32
Median	63.50	0	13	5	5.5	25	27.27
Avg	439	12.15	55.30	22.45	32.65	110.40	29.66
Max	4,189	179	479	93	185	671	59.04
Total	8,780	243	1,106	449	653	2,208	25.15

586 code occur after changes in the related specification file (**Spec→Code** column),
587 (v) the number of pairs for which changes in specification and production code
588 files influence each other (*i.e.*, the null hypothesis of the Granger’s Causality test
589 is rejected in both directions, **Both** column), and (vi) the number (and percent-
590 age) of overall pairs exhibiting any statistically significant causality relationship

591 (*i.e.*, Exact Match+Code→Spec+Spec→Code+Both, **Overall** column).

592 Looking at the quantitative results achieved for changes affecting both speci-
593 fications and fixtures (Table 7), we found that, for 9 out of 20 projects, more
594 than 54% of the considered specification-production code pairs exhibit causality
595 relationships. However, a comparison of the number of co-evolution pairs (*i.e.*,
596 the **Overall** column in Table 7 and Table 8) reveals that there is significantly
597 more co-evolution when both specifications and fixtures are changed, than for
598 changes affecting fixtures only (p -value=0.0009, $d = 0.26$ - small). The latter
599 implies that even when using a BDD framework, the feature addition and/or
600 modification results in modifying both production and test code.

601 For what concerns the co-evolution of specifications and fixtures together
602 with production code (Table 7), we found that the number of co-evolving pairs
603 where specifications/fixtures are changed together (**Exact Match** and **Both**
604 columns) or after the production code (**Code→Spec** column) is significantly
605 greater than the number of co-evolving pairs where specifications/fixtures change
606 before production code (**Spec→Code** column). The paired Wilcoxon signed-
607 rank test returned a p -value=0.0005, and the Cliff's $d = 0.4$ - medium. Similar
608 results were obtained when analyzing the co-evolution of changes to fixtures only
609 with production code (p -value=0.0005, $d=0.33$ - small).

610 Although some differences can be observed across projects, the median lag
611 between changes of specification and source code (in a direction or in the other)
612 is of only 3 commits. This clearly contrasts with distances observed in the
613 additions discussed above. That is, while specifications and related source code
614 files can be added at a distance of years (for the reasons explained above), once
615 they are in, they exhibit a close co-change, either in a direction (specification
616 before code) or in the other (code before specification).

617 A manual inspection of pairs exhibiting statistically significant co-evolution
618 values indicates that very often the additions and/or modifications of test cases
619 are clearly driven by changes applied to production code, hence violating the
620 BDD principles. As an example, we report a case occurring in the *middleman*
621 project, in which a new production file is added to address a pull request aimed at

622 implementing new piece of functionality. Before merging the pull request, there
623 is a comment asking for properly testing the implemented functionality¹⁶. Once
624 created, the production and test files are changed together as a consequence
625 of feature improvements, changes, or refactoring operations. Similar results
626 have been identified for the *httparty* project where the modifications to the
627 production code file `lib/httparty/request.rb`¹⁷¹⁸ are very often followed by
628 changes (in subsequent commits) introducing new testing scenarios within the
629 spec file `spec/httparty/request_spec.rb`¹⁹²⁰.

630 In other cases, it happens that after some changes occurred on both test
631 and production files, developers start modifying only the fixtures by applying
632 refactoring operations or test code improvements. For instance, in the *huginn*
633 project, a commit “*to increase test stability*”²¹ can be noticed. Finally, by looking
634 at cases in which developers change the fixtures only, we found that this is done
635 in order to make tests passing on Travis-CI by (i) accounting for implementation
636 details not contemplated in the test case yet²², or (ii) by simply skipping the
637 test case, *e.g.*, when it failed for unknown reasons²³.

638 Even if it seems that developers do not follow the BDD practice, we found a
639 few cases in which the specifications included in the test files have been used to
640 guide functional implementation. For instance, for the *middleman* project we
641 report a commit²⁴ in which the addition of a set of example groups is aimed
642 at testing a specific feature previously implemented. After that, the developers
643 realize that a conditional branch included in the test specifications was not
644 implemented at all in the production file, and, thus, the production code file is

¹⁶<https://github.com/middleman/middleman/commit/9ba1dc04>

¹⁷<https://github.com/jnunemaker/httparty/commit/50c46e2>

¹⁸<https://github.com/jnunemaker/httparty/commit/0f606b2>

¹⁹<https://github.com/jnunemaker/httparty/commit/604bfa6>

²⁰<https://github.com/jnunemaker/httparty/commit/dc14b30>

²¹<https://github.com/huginn/huginn/commit/771afe9>

²²<https://github.com/teamcapbara/capybara/commit/573668c>

²³<https://github.com/teamcapbara/capybara/commit/864fb62>

²⁴<https://github.com/middleman/middleman/commit/49a7435>

645 subsequently modified to implement the missing functionality²⁵.

646 In general, it may be possible that developers modify only the specifications
647 without altering the fixtures. However, looking at those cases we found many
648 situations in which changes were aimed at (i) addressing compatibility issues²⁶,
649 (ii) updating to a new Ruby version syntax²⁷, or (iii) fixing typos or spelling
650 issues in the description of the behavior²⁸.

BDD scenarios are almost always introduced in the code along with the related production code, or (even a long time) after the code has been added. BDD scenarios co-evolve with production code in up to 37% of the cases, but once again it is more likely that they change together (or after) production code than before (the median lag is of 3 commits). In essence, only in a minority of cases, the co-evolution analysis indicates symptoms of a likely application of BDD principles.

651

652 4.3. How do developers use BDD-related frameworks?

653 We received in total 35 responses, of which four were discarded as the
654 respondents declared to have a low experience on testing. We, therefore, analyzed
655 a total of 31 responses. Out of the 31 respondents, 19 own a Bachelor Degree (or
656 some college degree), 9 a Master's degree, and 3 a Ph.D. 14 of them have more
657 than 10 years of experience in programming, 8 between 5 and 10, and 9 less than
658 5. Most of them (21) are professional developers, but we got 2 responses from
659 academics, 2 (Ph.D.) students, 2 industrial researchers, 2 IT managers, and 2
660 QA consultants. Their level of experience in testing activity is (self)-rated from
661 moderate (14 of them) to high (17 of them).

662 The adoption of programming languages, with the exception of PHP, reflects
663 quite well the choices of our study. 17 respondents indicated to develop in

²⁵<https://github.com/middleman/middleman/commit/6ab6669>

²⁶<https://github.com/vcr/vcr/commit/00f4309>

²⁷<https://github.com/capistrano/capistrano/commit/2c14957>

²⁸<https://github.com/middleman/middleman/commit/b42e4ed>

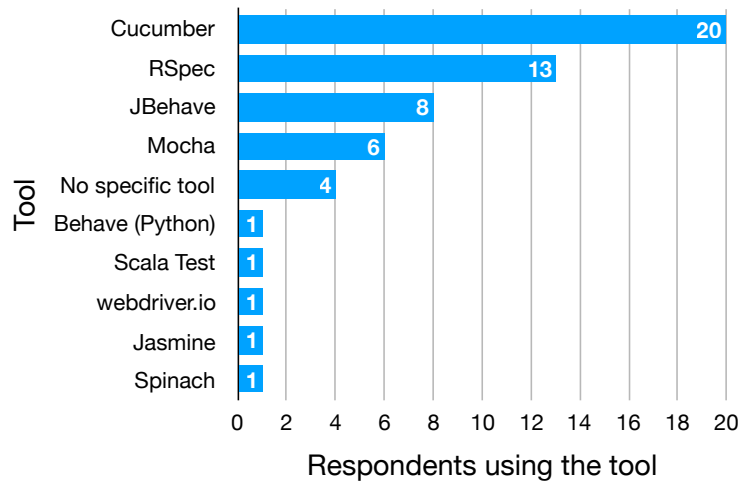


Figure 3: BDD frameworks used by respondents.

664 Javascript, 15 in Java, 15 in Ruby, 11 in Python, 2 in PHP, and one for each of
 665 the following languages: Kotlin, Scala, Objective-C, and Swift. Note that values
 666 do not sum up to 31 as respondents could indicate multiple languages.

667 As shown in Figure 3, in terms of frameworks being used, unsurprisingly, most
 668 of the respondents (20) indicated CUCUMBER, but this is also because CUCUMBER
 669 was mentioned as a framework being used by programmers working in various
 670 languages (all five we considered in our study). Nearly all Ruby programmers use
 671 RSpec (13 out of 15), while Mocha is less frequent for Javascript programmers
 672 (only 6 out of 17), and JASMINE used by only one respondent. 8 out of 15 Java
 673 programmers indicated to use Jasmine. Finally, some developers mentioned
 674 specific frameworks we did not consider in \mathbf{RQ}_1 , *i.e.*, SCALA TEST (as we did
 675 not study Scala), WEBDRIVER.IO (which however is a Web testing tool using the
 676 other frameworks we considered) and SPINACH. It is important to notice that,
 677 with the exception of the last three frameworks, all others have been considered
 678 in our \mathbf{RQ}_1 investigation. Finally, our pool of respondents also included four
 679 respondents indicating that they did not use specific frameworks for BDD, a
 680 practice used in the past for automated acceptance test [6].

681 By looking on the testing practices, our study participants in principle agreed

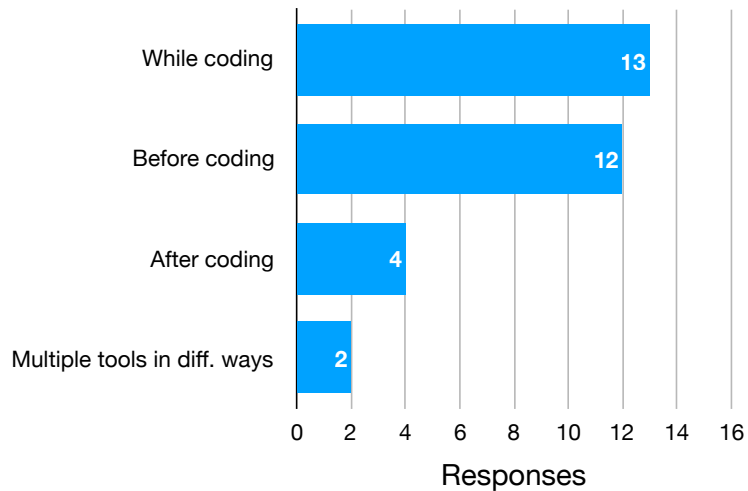


Figure 4: When test cases are written.

682 on the intended purpose of the BDD frameworks, *i.e.*, writing scenarios to model
 683 the system behavior before coding, and therefore performing BDD (17 responses)
 684 or, at least, performing TDD (5). However, there are also respondents using the
 685 frameworks for other testing purposes (6), or for multiple purposes (2).

686 However, when we asked when developers write test cases (Figure 4), only
 687 12 of them indicated that they write test cases before coding. 13 respondents do
 688 it while coding and 4 after coding.

689 In the open comments²⁹, on the one side, some respondents highlight the
 690 difficulties in applying BDD in practice: *“it’s just an overly complicated way*
 691 *to write unit tests”*. Also, and despite results of previous literature [6, 8, 9, 10],
 692 developers are skeptical about the advantages of BDD *“I have found that behavior*
 693 *and test-driven development results in less than optimal code quality and slower*
 694 *overall development with the benefit of a marginal improvement in bug reports.”*
 695 On the other side, they point out how BDD DSLs enforce cleaner and consistent
 696 write-up of scenarios *“significance of using a DSL and the defined glossary is to*

²⁹Due to the limited number of open comments, we simply report and discuss them, as coding would not make sense.

697 *develop a shared understanding of the problem.*” Also, they highlight that writing
698 specifications in a development team has, sometimes, different purpose than
699 using them for testing: it can be misleading to talk about “*BDD testing tools,*
700 *as writing and automating executable specifications is quite different to testing*”.
701 Moreover, one developer mentions how BDD is being used during evolution
702 activities, including refactoring to automatically check regressions and/or that
703 pre- and post-conditions still hold after refactoring. A different developer, instead,
704 highlights that “*I more often write tests first when refactoring code or when*
705 *working in scenarios where the programming interface is already well-defined. I*
706 *more often write tests afterwards when what I need to write is less well-defined,*
707 *requiring more exploratory coding.*” In the former case, she is applying TDD
708 instead of BDD, while in the latter case, despite the BDD approach may help in
709 defining/clarifying the system behavior under specific conditions, the surveyed
710 developer declares that she prefers starting with coding. Finally, there is a
711 respondent highlighting that the scenarios are mainly written by the product
712 owner instead of the developers “*The product owner (non-developer) writes the*
713 *scenarios.*”

While respondents understood the purpose of the BDD frameworks they use, in practice many of them write tests while (or after) coding, hence they do not perform BDD and not even TDD. Respondents felt that while BDD frameworks enforce cleaner and more consistent specifications, the effort in adopting them might not be paid back by an improved code quality.

714

715 **5. Threats To Validity**

716 Threats to *construct validity* concern the relationship between theory and
717 experimentation. In **RQ₁**, while we used several heuristics to select projects
718 that help us to avoid inactive projects, duplicates (forks), and projects without
719 source code, it is still possible that our sample contains non-development projects.
720 Considering that we also selected projects based on their number of commits, and

721 considering the large sample, we do not expect the percentage of such projects
722 to be relevant so to introduce substantial noise in our results.

723 The heuristics used in **RQ₁** to detect the tool usages might be imprecise
724 or erroneous. We mitigated this threat through manual analysis of a few
725 projects for each programming language and BDD framework. Similarly, the
726 fine-grained analysis of changes in **RQ₂** is based on `git diff` and regular
727 expressions, and despite some tests of our scripts, imprecisions are still possible.
728 Moreover, traceability links between specification and production files have
729 been heuristically recovered by exploiting directory structures and file naming
730 conventions. However, some links could have been missed or wrongly identified.
731 To alleviate this concern we manually checked that all the selected projects
732 comply with the aforementioned conventions. As for **RQ₃**, the major threat
733 is represented by the limited understanding of BDD by respondents. This is
734 partially mitigated by the relatively high experience and expertise they declared,
735 in particular about testing.

736 In **RQ₂**, the observation of whether BDD is applied in projects is limited to
737 the analysis of Granger causality. Further investigation might be required to
738 verify whether the process is properly followed. In **RQ₃**, we based our conclusions
739 from respondents' feedback, which might deviate from their actual practice.

740 Threats to *internal validity* concern factors internal to our study that could
741 influence our conclusions. Recent research [41] studied the perils of mining
742 GitHub repositories for software engineering research, observing that the majority
743 of projects are personal and inactive. To avoid considering similar projects (*i.e.*,
744 personal and inactive) in **RQ₁**, we selected the top 10,000 projects (in terms of
745 stars) for each programming language, excluding forked projects, projects having
746 no changes in the last year (*i.e.*, 2018), and projects having a single contributor.

747 As we explained in Section 3.1, we used a simple heuristic to unify identities
748 from git repositories. Although more complex and reliable heuristics exist for
749 this purpose [42, 43], but they are needed when merging identities from different,
750 heterogeneous repositories.

751 While through Granger's test we establish a causal relationship on a statistical

752 basis, we cannot claim the presence of a cause-effect relationship between changes
753 to tests and production code, or vice versa. We mitigate this threat with a
754 qualitative analysis detailed in **RQ₂**. In addition, similarly to previous work
755 [3], to study the co-evolution of specification and production code, we rely on
756 the data that is stored in the version control system. However, especially for
757 additions/modifications of specifications and related production files that occur
758 in the same commit, we can not establish which elements are added/modified
759 first. This issue is partially mitigated by the empirical assessment carried out in
760 our **RQ₃**, in which the majority of survey respondents stated that they usually
761 perform modifications to specifications while or after coding.

762 **RQ₃** is affected by the self-selection of its participants, which could have
763 influenced the kinds of provided answers. Although we do not claim full repre-
764 sentativeness of the universe of developers using BDD practices and tools, it is
765 possible that those who responded were the ones more inclined towards BDD,
766 or at least more competent. Finally, another possible threat for **RQ₃** is that
767 projects might use multiple BDD tools over their evolutionary history. We have
768 mitigated this threat by checking that, for the studied projects, RSpec was used
769 across the entire observed history.

770 Threats to *conclusion validity* concern the relationship between experimen-
771 tation and outcome. Wherever possible, in particular in **RQ₂**, we base our
772 conclusion on statistical procedures, including Granger’s test, Wilcoxon signed
773 rank test, and Cliff’s delta effect size. Moreover, we cope with multiple tests by
774 adjusting p -values with the Benjamini-Hochberg procedure [28].

775 Finally, threats to *external validity* concern the generalization of our findings.
776 As for **RQ₁**, although the sample is large (50,000 projects) and diverse in terms
777 of programming languages, it might not fully represent the whole picture of
778 the open-source community. More importantly, it does not represent the BDD
779 practice in the industry, and in general in a context where there is a better
780 formalization of scenarios than in open source projects where, often, features
781 are added thanks to the contribution of external developers (*i.e.*, through pull
782 requests). Also, in the search for used frameworks, we may have missed some

783 relevant ones, although we based our analysis on those mentioned in white
784 papers about BDD and our choices are to a major extent in-line with our survey
785 respondents. In addition, for **RQ₂** we intendedly decided to do a focused analysis
786 on a fairly limited set of projects developed with a specific language (Ruby)
787 and using a specific BDD framework (RSpec). Therefore, future studies on
788 other languages and frameworks are needed. Finally, for **RQ₃** the sample of
789 respondents is made up of experts (as they admitted) that well represent the
790 technologies investigated in **RQ₁**. However, a larger investigation on more
791 subjects beyond our 31 respondents is surely desirable.

792 **6. Implications**

793 Results of our study have implications for developers, educators, and re-
794 searchers.

795 For what concerns **developers**, our study shows that, while a plethora of
796 BDD frameworks exist, their adoption has still a long way to go. Only a small
797 percentage of projects use such tools and, more importantly, when this happens,
798 it is done by a very small percentage of developers. In summary, most of the
799 developers are still not ready (or not willing) to use such a relatively new piece
800 of technology and its related development methodology.

801 Such low percentages also allow conjecturing that, in general, there is no
802 application of BDD practices at project-level. While we do not have specific
803 evidence of that, it could either happen that in a project only a few developers
804 act as scenario-writers (but this contradicts agile principles fostering a shared
805 knowledge and discouraging a clear separation of roles) or, more likely, BDD
806 tools are opportunistically adopted by a limited number of developers, while the
807 project has not set their usage as a recommended practice.

808 There are noticeable exceptions, like the Ruby language and its RSpec
809 framework. However, a deeper observation made through a co-evolution analysis
810 (**RQ₂**) and a survey (**RQ₃**) indicates that very often developers write/change
811 BDD specifications while coding or even after that. Truly, BDD frameworks are

812 also suited for other testing purposes, however having a clear separation between
813 acceptance tests and other (*e.g.*, unit, integration) tests is surely desirable [44,
814 45, 46].

815 As previous research has shown [6, 10], writing executable scenarios is still
816 overkill for developers. Therefore, while BDD-related technology is nowadays
817 in place, and while the answers to our survey questionnaire indicate that in
818 most cases the participants are knowledgeable about BDD, the study has shown
819 that such a methodology is not exploited as it should be yet. Moreover, there
820 could be a great potential in using BDD tools as a part of a DevOps Continuous
821 Integration and Delivery (CI/CD) Pipeline, helping the team to monitor the
822 development progress of scenarios assigned to a release. However, once again,
823 the current BDD adoption level, and the way BDD tools are used, make such
824 integration with CI/CD still limited to simply executing test cases, *i.e.*, like any
825 other unit test framework.

826 To favor a wider and proper adoption of BDD tools, **educators** should
827 properly introduce the BDD methodology when teaching agile development,
828 requirement engineering, and testing, and possibly highlight its strengths, but
829 also open challenges in its applications. This will not only increase the awareness
830 about BDD tools but also let developers be aware of how such tools could
831 properly be applied in the context of a development process (*i.e.*, not limiting
832 their use as they were simple unit test tools). For example, since the results
833 of **RQ₂** and **RQ₃** indicated that, in many circumstances, developers prefer to
834 write tests after the code, we argue that it is also important for educators to
835 clearly illustrate, also with the support of real-world examples, how the BDD
836 approach can help in clarifying/defining the system's behavior under specific
837 conditions even when it is not yet completely specified, as in the case of open
838 source development. Also, as explained above, educators could properly frame
839 BDD in the context of CI/CD, as a way to monitor development.

840 Answers provided to our survey seem to indicate that, while developers are
841 aware of how BDD works, the available tool support makes its application too
842 overkill. More specifically, survey respondents highlighted the difficulty and

843 the effort needed to write scenarios in the appropriate format, and in writing
844 fixtures to link them onto production code. As a consequence, there are still
845 many open areas on which **researchers** can put their effort to enhance the
846 applicability of BDD in practice, especially with better tool support for writing
847 scenarios (or converting existing scenarios and user stories in the format required
848 by BDD tools), and automatically-generating fixtures.

849 Based on the obtained results, there is also room for further empirical
850 investigations, for example, surveying projects and organizations encouraging
851 BDD for their development, or studying the benefits of BDD in the context of
852 CI/CD pipelines.

853 **7. Conclusion**

854 This paper has investigated the use of Behavior-Driven Development (BDD)
855 frameworks in the open-source community. The study has been conducted by
856 analyzing (**RQ₁**) the adoption of BDD frameworks in 50,000 GitHub projects
857 written in five different programming languages, (**RQ₂**) the co-evolution of BDD
858 specifications and production files in 20 Ruby projects using RSpec [47], and
859 (**RQ₃**) feedback provided by 31 developers about their knowledge and usage of
860 BDD frameworks.

861 Results of the study indicate how the adoption of BDD tools is still limited,
862 except for cases such as Ruby-RSpec where the tool is also used for unit testing.
863 The co-evolution analysis we performed makes it evident how, in practice,
864 developers do not write “scenarios-first”, but they mostly use BDD tools as any
865 other unit testing tool, and they mostly write tests together or after the related
866 production code. The feedback we collected through the survey indicated how
867 the current BDD technology is still too overkilled to be properly adopted in
868 practice.

869 As work-in-progress, we are investigating the extent to which the introduction
870 of BDD in an ongoing project has any effect on the product and process quality,
871 but also the use of BDD in combination with other practices, such as Continuous

872 Integration and Delivery. In the future, we envision to carry out a quantitative
873 and qualitative investigation on the costs and benefits deriving from the adoption
874 of BDD-like approaches. In this context, since our research empirically pointed
875 out the adoption of BDD-related frameworks such as RSPEC also for non-BDD
876 purposes, we also plan to shed some light on the motivations and practical
877 advantages of using BDD-related tools in such cases.

878 **References**

- 879 [1] C. Solís, X. Wang, A study of the characteristics of behaviour driven
880 development, in: 37th EUROMICRO Conference on Software Engineering
881 and Advanced Applications, SEAA 2011, Oulu, Finland, 2011, pp. 383–387.
- 882 [2] D. North, How to sell bdd to the business, <https://skillsmatter.com/skillscasts/923-how-to-sell-bdd-to-the-business> (last access:
883 01/10/2019).
- 884 [3] A. Zaidman, B. Van Rompaey, S. Demeyer, A. van Deursen, Mining software
885 repositories to study co-evolution of production & test code, in: First
886 International Conference on Software Testing, Verification, and Validation,
887 ICST 2008, Lillehammer, Norway, 2008, pp. 220–229.
- 888 [4] M. Beller, G. Gousios, A. Panichella, A. Zaidman, When, how, and why
889 developers (do not) test in their ides, in: Proceedings of the 2015 10th
890 Joint Meeting on Foundations of Software Engineering, ESEC/FSE 2015,
891 Bergamo, Italy, 2015, pp. 179–190.
- 892 [5] R. Mugridge, W. Cunningham, Fit for Developing Software: Framework for
893 Integrated Tests., Prentice Hall, 2005.
- 894 [6] G. Melnik, F. Maurer, Multiple perspectives on executable acceptance
895 test-driven development, in: Agile Processes in Software Engineering and
896 Extreme Programming, 8th International Conference, XP 2007, Como, Italy,
897 2007, pp. 245–249.
- 898

- 899 [7] G. Melnik, K. Read, F. Maurer, Suitability of FIT user acceptance tests
900 for specifying functional requirements: Developer perspective, in: *Extreme*
901 *Programming and Agile Methods - XP/Agile Universe 2004*, 4th Conference
902 on Extreme Programming and Agile Methods, Calgary, Canada, 2004, pp.
903 60–72.
- 904 [8] S. Park, F. Maurer, Communicating domain knowledge in executable accep-
905 tance test driven development, in: *Agile Processes in Software Engineering*
906 *and Extreme Programming*, 10th International Conference, XP 2009, Pula,
907 Sardinia, Italy, 2009, pp. 23–32.
- 908 [9] F. Ricca, M. Torchiano, M. Di Penta, M. Ceccato, P. Tonella, Using accep-
909 tance tests as a support for clarifying requirements: A series of experiments,
910 *Information & Software Technology* 51 (2) (2009) 270–283.
- 911 [10] F. Ricca, M. Di Penta, M. Torchiano, P. Tonella, M. Ceccato, C. A. Visaggio,
912 Are fit tables really talking?: a series of experiments to understand whether
913 fit tables are useful during evolution tasks, in: *30th International Conference*
914 *on Software Engineering (ICSE 2008)*, Leipzig, Germany, 2008, pp. 361–370.
- 915 [11] L. P. Binamungu, S. M. Embury, N. Konstantinou, Maintaining behaviour
916 driven development specifications: Challenges and opportunities, in: *25th*
917 *International Conference on Software Analysis, Evolution and Reengineering*,
918 *SANER 2018*, Campobasso, Italy, 2018, pp. 175–184.
- 919 [12] M. Rahman, J. Gao, A reusable automated acceptance testing architecture
920 for microservices in behavior-driven development, in: *2015 IEEE Symposium*
921 *on Service-Oriented System Engineering, SOSE 2015*, San Francisco Bay,
922 CA, USA, 2015, pp. 321–325.
- 923 [13] A. Z. H. Yang, D. A. da Costa, Y. Zou, Predicting co-changes between
924 functionality specifications and source code in behavior driven development,
925 in: *Proceedings of the 16th International Conference on Mining Software*
926 *Repositories, MSR 2019*, Montreal, Canada., 2019, pp. 534–544.

- 927 [14] F. Zampetti, A. Di Sorbo, C. A. Visaggio, G. Canfora, M. Di Penta,
928 Demystifying the use of behavior-driven development tools in open source
929 projects. Replication Package
930 <https://tinyurl.com/y9yacwyx>.
- 931 [15] G. S. Mahalakshmi, V. Vani, Theoretical verification of test cases for
932 behavior driven development, in: 2017 Second International Conference
933 on Recent Trends and Challenges in Computational Models (ICRTCCM),
934 2017, pp. 309–314.
- 935 [16] M. Soeken, R. Wille, R. Drechsler, Assisted behavior driven development using
936 natural language processing, in: Objects, Models, Components, Patterns
937 - 50th International Conference, TOOLS 2012, Prague, Czech Republic,
938 2012, pp. 269–287.
- 939 [17] G. Lucassen, F. Dalpiaz, J. M. E. M. van der Werf, S. Brinkkemper,
940 D. Zowghi, Behavior-driven requirements traceability via automated acceptance
941 tests, in: IEEE 25th International Requirements Engineering
942 Conference Workshops, RE 2017 Workshops, Lisbon, Portugal, 2017, pp.
943 431–434.
- 944 [18] M. Wynne, A. Hellesoy, S. Tooke, The cucumber book: behaviour-driven
945 development for testers and developers, Pragmatic Bookshelf, 2017.
- 946 [19] L. P. Binamungu, S. M. Embury, N. Konstantinou, Detecting duplicate
947 examples in behaviour driven development specifications, in: 2018 IEEE
948 Workshop on Validation, Analysis and Evolution of Software Tests, 2018,
949 Campobasso, Italy, 2018, pp. 6–10.
- 950 [20] A. Monden, T. Hayashi, S. Shinoda, K. Shirai, J. Yoshida, M. Barker,
951 K. Matsumoto, Assessing the cost effectiveness of fault prediction in acceptance
952 testing, IEEE Transactions on Software Engineering 39 (10) (2013)
953 1345–1357.

- 954 [21] A. Zaidman, B. Van Rompaey, A. van Deursen, S. Demeyer, Studying
955 the co-evolution of production and test code in open source and industrial
956 developer test processes through repository mining, *Empirical Software*
957 *Engineering* 16 (3) (2011) 325–364.
- 958 [22] Z. Lubsen, A. Zaidman, M. Pinzger, Using association rules to study the co-
959 evolution of production & test code, in: *Proceedings of the 6th International*
960 *Working Conference on Mining Software Repositories, MSR 2009 (Co-located*
961 *with ICSE)*, Vancouver, BC, Canada, 2009, pp. 151–154.
- 962 [23] C. Marsavina, D. Romano, A. Zaidman, Studying fine-grained co-evolution
963 patterns of production and test code, in: *14th IEEE International Work-*
964 *ing Conference on Source Code Analysis and Manipulation, SCAM 2014,*
965 *Victoria, BC, Canada, 2014*, pp. 195–204.
- 966 [24] G. Canfora, A. Cimitile, F. García, M. Piattini, C. A. Visaggio, Evaluat-
967 ing advantages of test driven development: a controlled experiment with
968 professionals, in: *2006 International Symposium on Empirical Software*
969 *Engineering (ISESE 2006) Rio de Janeiro, Brazil, 2006*, pp. 364–371.
- 970 [25] C. W. J. Granger, Investigating causal relations by econometric models and
971 cross-spectral methods, *Econometrica* 37 (3) (1969) 424–438.
- 972 [26] J. D. Hamilton, *Time Series Analysis*, Princeton University Press, 1994.
- 973 [27] H. Akaike, Information theory and an extension of the maximum likelihood
974 principle, in: *2nd International Symposium on Information Theory, 1973,*
975 pp. 267–281.
- 976 [28] Y. Benjamini, Y. Hochberg, Controlling the false discovery rate: a practical
977 and powerful approach to multiple testing, *Journal of the Royal Statistical*
978 *Society, Series B* 57 (1) (1995) 125–133.
- 979 [29] A. Zeileis, T. Hothorn, Diagnostic checking in regression relationships, *R*
980 *News* 2 (3) (2002) 7–10.

- 981 [30] B. Pfaff, Analysis of Integrated and Cointegrated Time Series with R, 2nd
982 Edition, Springer, New York, 2008.
- 983 [31] B. Pfaff, Var, svar and svec models: Implementation within R package vars,
984 Journal of Statistical Software 27 (4).
- 985 [32] F. Wilcoxon, Individual comparisons by ranking methods, Biometrics Bul-
986 letin 1 (6) (1945) 80–83.
- 987 [33] R. J. Grissom, J. J. Kim, Effect sizes for research: A broad practical
988 approach, 2nd Edition, Lawrence Earlbaum Associates, 2005.
- 989 [34] S. L. Pfleeger, B. A. Kitchenham, Principles of survey research: part 1:
990 turning lemons into lemonade, ACM SIGSOFT Software Engineering Notes
991 26 (6) (2001) 16–18.
- 992 [35] B. A. Kitchenham, S. L. Pfleeger, Principles of survey research: part 5:
993 populations and samples, ACM SIGSOFT Software Engineering Notes 27 (5)
994 (2002) 17–20.
- 995 [36] B. A. Kitchenham, S. L. Pfleeger, Principles of survey research part 4:
996 questionnaire evaluation, ACM SIGSOFT Software Engineering Notes 27 (3)
997 (2002) 20–23.
- 998 [37] B. A. Kitchenham, S. L. Pfleeger, Principles of survey research: part 3:
999 constructing a survey instrument, ACM SIGSOFT Software Engineering
1000 Notes 27 (2) (2002) 20–24.
- 1001 [38] B. A. Kitchenham, S. L. Pfleeger, Principles of survey research part 2:
1002 designing a survey, ACM SIGSOFT Software Engineering Notes 27 (1)
1003 (2002) 18–20.
- 1004 [39] B. A. Kitchenham, S. L. Pfleeger, Principles of survey research part 6: data
1005 analysis, ACM SIGSOFT Software Engineering Notes 28 (2) (2003) 24–27.
- 1006 [40] R. M. Groves, F. J. J. Fowler, M. P. Couyper, J. M. Lepkowski, E. Singer,
1007 R. Tourangeau, Survey Methodology, 2nd edition, Wiley, 2009.

- 1008 [41] E. Kalliamvakou, G. Gousios, K. Blincoe, L. Singer, D. M. Germán, D. E.
1009 Damian, An in-depth study of the promises and perils of mining github,
1010 Empirical Software Engineering 21 (5) (2016) 2035–2071.
- 1011 [42] G. Robles, J. M. González-Barahona, Developer identification methods for
1012 integrated data from various sources, in: Proceedings of the 2005 Interna-
1013 tional Workshop on Mining Software Repositories, MSR 2005, Saint Louis,
1014 Missouri, USA, 2005.
- 1015 [43] C. Bird, A. Gourley, P. T. Devanbu, M. Gertz, A. Swaminathan, Mining
1016 email social networks, in: Proceedings of the 2006 International Workshop
1017 on Mining Software Repositories, MSR 2006, Shanghai, China, 2006, pp.
1018 137–143.
- 1019 [44] P. Duvall, S. M. Matyas, A. Glover, Continuous Integration: Improving
1020 Software Quality and Reducing Risk (The Addison-Wesley Signature Series),
1021 Addison-Wesley Professional, 2007.
- 1022 [45] G. Orellana, G. Laghari, A. Murgia, S. Demeyer, On the differences between
1023 unit and integration testing in the travistorrent dataset, in: Proceedings of
1024 the 14th International Conference on Mining Software Repositories, MSR
1025 2017, Buenos Aires, Argentina, 2017, pp. 451–454.
- 1026 [46] C. Vassallo, G. Schermann, F. Zampetti, D. Romano, P. Leitner, A. Zaidman,
1027 M. Di Penta, S. Panichella, A tale of CI build failures: An open source and a
1028 financial organization perspective, in: 2017 IEEE International Conference
1029 on Software Maintenance and Evolution, ICSME 2017, Shanghai, China,
1030 2017, pp. 183–193.
- 1031 [47] RSpec. <http://rspec.info> (last access: 01/10/2019).