

Investigating Developers' Perception on Software Testability and its Effects

Tushar Sharma¹✉, Stefanos Georgiou², Maria Kechagia³, Taher A. Ghaleb⁴, Federica Sarro³

¹Dalhousie University, Canada; ²simpleTechs, Germany; ³University College London, United Kingdom; ⁴University of Ottawa, Canada

✉ For correspondence:
tushar@dal.ca

Data availability: Replication package can be found on GitHub - <https://github.com/SMART-Dal/testability>

Funding: Maria Kechagia and Federica Sarro are supported by the ERC grant no. 741278 (EPIC).

Competing interests: The authors declare that they have no conflicts of interest.

Abstract

The opinions and perspectives of software developers are highly regarded in software engineering research. The experience and knowledge of software practitioners are frequently sought to validate assumptions and evaluate software engineering tools, techniques, and methods. However, experimental evidence may unveil further or different insights, and in some cases even contradict developers' perspectives. In this work, we investigate the correlation between software developers' perspectives and experimental evidence about testability smells (*i.e.*, programming practices that may reduce the testability of a software system). Specifically, we first elicit opinions and perspectives of software developers through a questionnaire survey on a catalog of four testability smells, we curated for this work. We also extend our tool `DESIGNITEJAVA` to automatically detect these smells in order to gather empirical evidence on testability smells. To this end we conduct a large-scale empirical study on 1,115 Java repositories containing approximately 46 million lines of code to investigate the relationship of testability smells with test quality, number of tests, and reported bugs. Our results show that testability smells do not correlate with test smells at the class granularity or with test suit size. Furthermore, we do not find a causal relationship between testability smells and bugs. Moreover, our results highlight that the empirical evidence does not match developers' perspective on testability smells. Thus, suggesting that despite developers' invaluable experience, their opinions and perspectives might need to be complemented with empirical evidence before bringing it into practice. This further confirms the importance of data-driven software engineering, which advocates the need and value of ensuring that all design and development decisions are supported by data.

Keywords: Software testability; software test quality; testability smells; developers' opinions and perspectives; software quality.

1. Introduction

The opinions and perspectives of software developers matter significantly in software engineering research. The research in the domain relies on the experience and knowledge of software practitioners to validate assumptions and evaluate the tools, techniques, and methods addressing software engineering problems. Numerous studies reveal the importance of the practitioners' perspectives [14, 55, 53, 1]. However, empirical evidence may not always agree with developers' perspective, opinions, or beliefs. For example, Devanbu *et al.* [13] show that software developers

40 have very strong beliefs on certain topics, but are often based on their personal experience; such
41 beliefs and corresponding empirical evidence may be inconsistent. Similarly, despite developers'
42 fairly common negative perspectives about code clones, Rahman *et al.* [54] did not find sufficient
43 empirical evidence to prove a strong correlation between code clones and bug proneness; how-
44 ever, there could be effects of code clones other than bugs. Along the similar lines, a study by
45 Murphy *et al.* [43] casts doubts on practitioners' and researchers' assumptions related to refactor-
46 ing. Finally, Janes *et al.* [28] challenge the practitioners' perspective on agile process by bringing
47 out the dark side of agility. In this paper, we present a case study of *testability smells*, i.e., *program-*
48 *ming practices that may negatively affect testability of a software system*, to investigate whether the
49 software developers' perspectives about testability smells is backed up by empirical evidence .

50 Researchers and practitioners have proposed various, yet inconsistent, definitions of software
51 testability [21]. The most common definition refers to *the degree to which the development of test*
52 *cases can be facilitated by the software design choices* [9, 73, 7]. Specifically, several researchers [9, 73,
53 50] define testability as the *ease of testing*. In addition, some researchers [9, 7, 75] emphasized that
54 testability is not a binary concept but must be expressed in *degree* or *extent*. Additionally, other
55 researchers [9, 75, 52] explicitly connect software *design choices* with the definition of software
56 testability. Furthermore, some studies [9, 50, 75] identify the degree of *effectiveness* by which test
57 development is facilitated as another characteristic of testability definition.

58 There has been a significant amount of work on test smells and their effects [30, 60, 4]. Test
59 smells are bad programming practices in test code that negatively affect the quality of test suites
60 and production code [23]. Though looks similar, practices that affect testability are completely
61 different than test smells. Test smells occur in test code while issues affecting testability arise in
62 production code. Also, test smells indicate the poor quality of test cases whereas testability impacts
63 the ability to write tests.

64 Several researchers have proposed frameworks for measuring and empirically evaluating soft-
65 ware testability by considering (1) software design choices, including programming language fea-
66 tures [71, 59], (2) software quality metrics, including cohesion, depth of inheritance tree, coupling,
67 and propagation rate of methods [39], metrics including the number of method calls, dependen-
68 cies, and attributes of a class with testability [41], and other software metrics [32, 59, 8], as well as
69 (3) testing effort [64]. However, there are still various aspects related to software testability that
70 remain unexplored, including the extent to which specific programming practices (e.g., handling of
71 dependencies, coding style, and access of modifiers) impact software testability.

72 The *goal* of this study is to find experimental evidence to validate current practitioners' perspec-
73 tives about testability smells. To achieve this goal, we first curate a catalog of four programming
74 practices which can affect software testability, referred to as *testability smells*. We then gather soft-
75 ware developers' perspectives through a questionnaire survey on testability in general, and our
76 proposed testability smells in particular. Finally, we conduct a large-scale empirical study guided
77 by three research questions to explore the effect of testability smells on test cases, their quality,
78 and on reported bugs at different granularity levels. To support the detection of testability and
79 test smells, we develop a tool named DESIGNITEJAVA. To answer the research questions, we curated
80 a dataset of 1, 115 Java software projects which are publicly available in GITHUB, and analyzed them
81 using our DESIGNITEJAVA tool.

82 Our survey shows that software developers consider testability as a factor that impedes soft-
83 ware testing and overwhelmingly acknowledged the proposed testability smells. Our results sug-
84 gest that testability smells show a low positive correlation with test smells at the repository granu-
85 larity; however, at the class-level, testability smells and test smells do not correlate. Our explo-
86 ration of the relationship between testability smells with test density reveals no correlation at
87 repository and class granularity. Finally, our observations from our experiment indicate that testa-
88 bility smells do not contribute to bugs. Therefore, developers' opinions and perspectives might
89 need to be complemented with empirical evidence before bringing it into practice.

90 This study makes the following contributions to the field.

- 91 1. We investigate the extent to which developers' perspectives is in line with the empirical evi-
92 dence we found in the context of *i.e.*, testability smells.
- 93 2. We consolidate a set of programming practices that affect testability of a software system in
94 the form of a catalog of testability smells. This catalog provides a vocabulary for researchers
95 and practitioners to discuss specific programming practices potentially impacting the testa-
96 bility of software systems.
- 97 3. We extend `DESIGNITEJAVA` to detect the proposed testability smells and eight test smells. The
98 tool facilitates further research on the topic of testability smells. Also, interested software
99 developers may use this tool to detect testability smells in their source code to better under-
100 stand the impact of design choices on testing.
- 101 4. We explore the relationships between testability smells and several aspects relevant to test
102 development and bugs. Such an exploration improves our understanding, both as software
103 developers and researchers, of testability and lays the groundwork for devising new tools
104 and techniques to improve test development.

105 We have made publicly available our `DESIGNITEJAVA` to identify testability smells as well as a
106 replication package at <https://github.com/SMART-Dal/testability>. We hope this facilitate other re-
107 searchers to replicate, reproduce and extend the presented study.

108 The rest of this paper is organized as follows. First, we present related work in Section 2. Sec-
109 tion 3 provides overview of the methods. Section 4 presents the initial catalog of testability smells,
110 our questionnaire survey targeting to software practitioners and obtained results, and tool imple-
111 mentation to detect testability and test smells. We present the mechanism followed to select and
112 download repositories from `GITHUB` in Section 5. We discuss results in Section 6 and their implica-
113 tions in Section 7. Threats to validity are discussed in Section 8. Finally, we conclude in Section 9.

114 2. Related work

115 **Software testability.** From existing studies, we found that testability was initially considered for
116 hardware design [36, 72, 37]. The concepts of hardware testability were then used for software
117 testability [44, 34]. Afterwards, a great deal of studies has been conducted exploring various as-
118 pects of software testability. To measure testability for data-flow software, Nguyen *et al.* [45] sug-
119 gested an approach that uses the `SATAN` method, which transforms the source code into a static
120 single assessment form. The form is then fed into a testability model to detect source code parts
121 with testability weaknesses. Bruntink *et al.* [8] collected a large number of source code metrics
122 (*e.g.*, depth of inheritance tree, fan out, and lack of cohesion of methods) and test code metrics
123 to explore the relationship with testability. The analysis focused on open-source Java applications.
124 The results suggest that there is a significant correlation between class-level metrics (most notably
125 fan out, `LOC` per class, and response for class) and test-level metrics (`LOC` per class and the num-
126 ber of test cases). Vincent *et al.* [69] investigated software components testability written in C++
127 and Java in workstations and embedded systems. Moreover, the authors have suggested an ap-
128 proach named *built-in-test* for run-time-testability which can provide more testable, reliable, and
129 maintainable software components. Filho *et al.* [18] used ten testability attributes, proposed in pre-
130 vious studies, to examine their correlation with source code metrics and test specification metrics
131 (*e.g.*, number of test cases, test coverage) on two Android applications. They found that testability
132 attributes are correlated with several source code metrics and test specification metrics. Chowd-
133 hary [9] presented experiences while applying testability concepts and introduced guidelines to
134 ensure that testability is taken under consideration during software planning and testing. Based
135 on these findings, the authors introduced a testability framework called `SHOCK`. Furthermore, var-
136 ious resources [27, 75] discussed their interpretation of testability and impact of smells affecting
137 testability.

138 **Assessing testability.** Voas [29] surveyed the factors that affect software testability, arguing that
139 a piece of software that is likely to reveal faults within itself during testing is said to have high

140 testability. According to this work, *information loss* is a phenomenon that occurs during program
141 execution and increases the likelihood that a fault will remain undetected. Finally, Voas [71] com-
142 pared the testability of both object-oriented and procedural systems, as well as whether testability
143 is affected by programming language characteristics.

144 **Surveys on testability.** Various literature surveys on testability have been carried out. Freed-
145 man [19] investigated the testability of software components. Freedman argued that the concept
146 of domain testability of software is defined by applying the concepts of observability and control-
147 lability to software. Garousi *et al.* [21] examined 208 papers on testability (published between
148 1982 and 2017) and also found that the two most commonly referred factors affecting testabil-
149 ity are observability and controllability. Furthermore, their survey argues that common ways to
150 improve testability are testability transformation, improving observability, adding assertions, and
151 improving controllability. Similarly, Hassan *et al.* [25] conducted a systematic literature review on
152 software testability to investigate to what extent it affects software robustness. Results show that
153 a variety of testability issues are indeed relevant, with observability and controllability issues be-
154 ing the most researched. They also found that fault tolerance, exception handling, and handling
155 external influence are prominent robustness issues.

156 **Test smells.** A wide variety of studies explored test smells and their effect and relationship on
157 various aspects of software development, including change and bug proneness [60], maintainabil-
158 ity [5, 33, 67], and test flakiness [16]. Specifically, Spadini *et al.* [60] investigated the relationship
159 between the presence of test smells and the change-and defect-proneness of test code, as well
160 as the defect-proneness of the tested production code. Among their findings, they observed that
161 tests with smells are indeed more change- and defect-prone. Regarding maintainability, Bavota *et*
162 *al.* [5] presented empirical studies on test smells, and showed that test smells have a strong nega-
163 tive impact on program comprehension and maintenance. They, also, found that comprehension
164 is 30% better in the absence of test smells. Furthermore, Kim *et al.* [33] conducted an empirical
165 study to study the evolution and maintenance of test smells. They found that the number of test
166 smells increases as a system evolves, and through a qualitative analysis they revealed that most
167 test smell removal is a maintenance activities. Additionally, Tufano *et al.* [67] showed that test
168 smells are usually introduced when the corresponding test code is committed in the repository for
169 the first time. Then, those test smells tend to remain in a system for a long time, hindering soft-
170 ware maintenance. Fatima *et al.* [16] developed an approach called *Flakify*, which is a black-box,
171 language model-based predictor for flaky test cases.

172 Despite extensive work on testability, the existing literature does not translate high-level prin-
173 ciples such as observability and controllability into actionable programming practices. Due to that
174 though the high-level testability principles have been known to the community for a long time,
175 there has been no tool support to detect them. We provide a tool that supports the detection of
176 testability smells. Furthermore, we explore the relationship between testability practices and test
177 quality and size, which is our other contribution.

178 3. Overview of the Methods

179 In the pursuit of weighing developers' perceptions with empirical evidence in the context of testa-
180 bility smells, we formulate the following research questions.

181 **RQ1.** *To what extent do testability smells and test smells correlate?*

182 Testability smells refer to bad programming practices that are believed to make test case
183 design and development difficult. Developers may choose to follow non-optimal practices
184 when it is not easy to write tests, leading to poor-quality test cases. Test smells refer to bad
185 programming practices in unit test code, compromising test code quality by violating the best
186 practices recommended for writing test code [12]. This research question explores whether
187 and to which extent testability smells and test smells correlate.

188 **RQ2.** *Do testability smells correlate with test suite size?*

189 By definition, testability smells make the design and development of test cases difficult. With
190 this research question, we aim to empirically evaluate whether the presence of testability
191 smells can hinder test development and consequently lead to a fewer number of test cases.
192 By answering this research question, we can inform developers about testability smells that
193 might impede a smooth test development.

194 **RQ3.** *Do testability smells cause more bugs?*

195 Testability smells make it harder for developers to test a software system. This implies that
196 the software under test lacks appropriate testing, leaving more bugs uncovered during soft-
197 ware development. This research question aims to investigate whether and to which extent
198 testability smells can lead to a higher number of reported bugs.

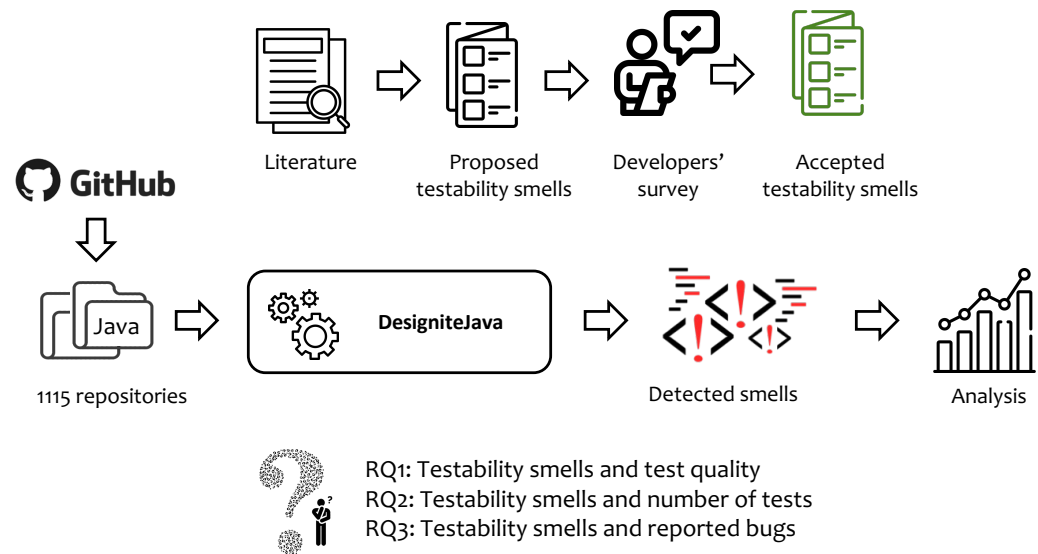


Figure 1. Overview of the study

199 Towards the goal of the study, as outlined in Figure 1, we first prepared a set of potential testa-
200 bility smells based on the available literature and recommended practices. We then carried out
201 an online survey to understand developers' perspectives on software testability and to gauge the
202 extent to which they agree that those smells really impact testability negatively. We extended our
203 tool (DESIGNITEJAVA) to detect testability and test smells. We analyzed 1,115 Java repositories down-
204 loaded from GITHUB. After identifying the smells, we reported our observations and findings with
205 respect to each research question.

206 4. Testability smells

207 We define testability smells as the programming practices that reduce the testability of a software
208 system. This section presents an initial catalog of testability smells, validates them by carrying out
209 an online developers' survey, and discusses the implementation details of our tool.

210 4.1 Initial Catalog of Testability Smells

211 In order to identify specific programming practices that negatively affect testability, we carried out
212 a light-weight multi-vocal literature (MLR) review, which surveys writings, views, and opinions in
213 diverse form and format [22]. The review process has three main stages: *search*, *selection*, and
214 *information extraction*.

215 In the search stage, two of the co-authors searched for a set of search terms (including testa-
216 bility, ease of testing, and software design+test) on Google Scholar and Google Search. For each

217 search term, we manually searched minimum seven pages of search results. After the minimum
218 threshold of seven pages, we continued the search until we get two continuous search pages with-
219 out any new and relevant articles. Adopting this mechanism avoided missing any relevant articles
220 in the context of our study.

221 We applied inclusion and exclusion criteria to filter out irrelevant sources. The main inclusion
222 criterion was that the source's content must relate to testability. Examples of exclusion criteria
223 include dropping gray literature that is too short, written in language other than English, or pre-
224 sented without objectivity in presentation. These examples map to the *Objectivity* requirement of
225 MLR process guidelines [22].

226 In the last stage, we read or observed the selected resources and extracted information rele-
227 vant to our study. Specifically, we strived for concrete recommendations in terms of programming
228 practices that influence testability from the selected sources. We grouped the practices based on
229 similarity and assigned an appropriate name reflecting the rationale. We identified four potential
230 testability smells discussed by more than one selected source. We present the consolidated set of
231 smells below. It is the first attempt, to the best of our knowledge, to document specific program-
232 ming practices as testability smells. It is by no means a comprehensive list of testability smells; we
233 encourage the research community to further extend this initial catalog of testability smells.

234 4.1.1 Hard-wired dependencies

235 This smell occurs when a concrete class is instantiated and used in a class resulting in a *hard-wired*
236 *dependency* [9, 75, 26]. A *hard-wired dependency* creates tight-coupling between concrete classes
237 and reduces the ease of writing tests for the class [9]. Such a *hard-wired dependency* makes the
238 class difficult to test because the newly instantiated objects are not replaceable with test doubles
239 (such as stubs and mocks). Hence, the test will check not only the `CUT` (class under test) but also its
240 dependencies, which is undesirable.

241 In Listing 1, the `parse`¹ method creates an object of the `BindingOperation` class (line 4) and
242 calls a few methods (lines 6 and 7). The object cannot be replaced at testing execution due to the
243 concrete object creation and its use within this method. Hence, the hard-coded dependency is
244 reducing the ease of writing tests for the class.

```
245 1 private void parse(String name, String namespace, WsdlParser parser) throws WsdlParseException  
246 {  
247 2     if (WSDL_NS.equals(namespace)) {  
248 3         if (OPERATION.equals(name)) {  
249 4             BindingOperation operation = new BindingOperation(definitions);  
250 5             operation.read(parser);  
251 6             operations.put(operation.getQName(), operation);  
252 7         }  
253 8     }  
254 9     //rest of the method  
255 10 }
```

Listing 1. Example of hard-coded dependency

256 4.1.2 Global state

257 Global variables are, in general, widely discouraged [40, 63]. This smell arises when a global vari-
258 able or a Singleton object is used [26, 62, 17, 65]. Global variables create hidden channels of
259 communication among abstractions in the system even when they do not depend on each other
260 explicitly. Global variables introduce unpredictability and hence make tests difficult to write by
261 developers.

262 The `Builder`² class in Listing 2 is accessible, and hence can be read/written, within the entire
263 project. Such practice makes it difficult to predict the state of the object in tests.

¹<https://github.com/forcedotcom/wsc/blob/master/src/main/java/com/sforce/ws/wsdl/Binding.java>

²<https://github.com/forcedotcom/wsc/blob/master/src/main/java/com/sforce/async/JobInfo.java>

```

264 1 public static class Builder {
265 2     //class definition
266 3 }

```

Listing 2. Example of global state

267 4.1.3 Excessive dependency

268 This smell occurs when the class under test has excessive outgoing dependencies. Dependencies
 269 make testing harder; a large number of dependencies makes it difficult to write tests for the class
 270 under test in isolation [62, 39, 74]. For example, the `Error`³ class in the open-source project `wsc`
 271 refers to nine other classes within the project --- `Bind`, `BulkConnection`, `TypeInfo`, `StatusCode`,
 272 `XmlOutputStream`, `XmlInputStream`, `ConnectionException`, `TypeMapper`, and `Verbose`. Such a high
 273 number of dependencies on other classes increases the effort to write tests for this class to be
 274 tested in isolation.

275 4.1.4 Law of Demeter violation

276 This smell arises when the class under test violates the law of Demeter *i.e.*, the class is interacting
 277 with objects that are neither class members nor method parameters [38, 31, 65]. In other words,
 278 the class has a chain of method calls such as `x.getY().doThat()`. Violations of the law of Demeter
 279 create additional dependencies that a test has to take care of. For example, lines 4 and 5 of the
 280 snippet¹ given in Listing 3 call a method to obtain an object that in turn calls another method on
 281 the obtained object. Such method chains introduce indirect dependencies that reduce the ease of
 282 writing tests for the class.

```

283 1 public Iterator<Part> getAllHeaders() throws ConnectionException {
284 2     HashSet<Part> headers = new HashSet<Part>();
285 3     for (BindingOperation operation : operations.values()) {
286 4         addHeaders(operation.getInput().getHeaders(), headers);
287 5         addHeaders(operation.getOutput().getHeaders(), headers);
288 6     }
289 7     return headers.iterator();
290 8 }

```

Listing 3. Example of the law of Demeter violation

291 4.2 Developer Survey

292 We carried out an anonymous online questionnaire survey targeting software developers to under-
 293 stand their perspectives on software testability. Specifically, we aimed to consolidate developers'
 294 perspectives *w.r.t.* the definition of testability as well as the relevance of our identified testability
 295 smells. We divided our survey into three sections. In the first section, we collected information
 296 about developers' experience. In the second section, we asked developers how they define testa-
 297 bility. The final section presented our initial catalog of testability smells and asked the respondents
 298 whether and to what extent they agree that the presented practices negatively affect testability. All
 299 the questions in this section were Likert-scale questions. The questionnaire that we used is avail-
 300 able online [57].

301 Before rolling out the survey to a larger audience, we ran a pilot for the survey, collected feed-
 302 back, and improved the survey. We shared the survey on all online professional social media chan-
 303 nels (such as Twitter, LinkedIn, and Reddit) and sought participation from the software develop-
 304 ment community. We kept the survey open for six weeks. We received 45 complete responses.

305 4.2.1 Findings from the survey

306 Figure 2 presents the demographic distribution of participants in terms of years of experience
 307 classified by their roles. We asked them to check all applicable roles and hence the total number

³<https://github.com/forcedotcom/wsc/blob/master/src/main/java/com/sforce/async/Error.java>

308 of responses shown in the figure is more than the number of participants. It is evident that most
 309 of the participants belong to the “software developer” role; a significant number of the participants
 310 belong to the highly experienced group (11-20 years).

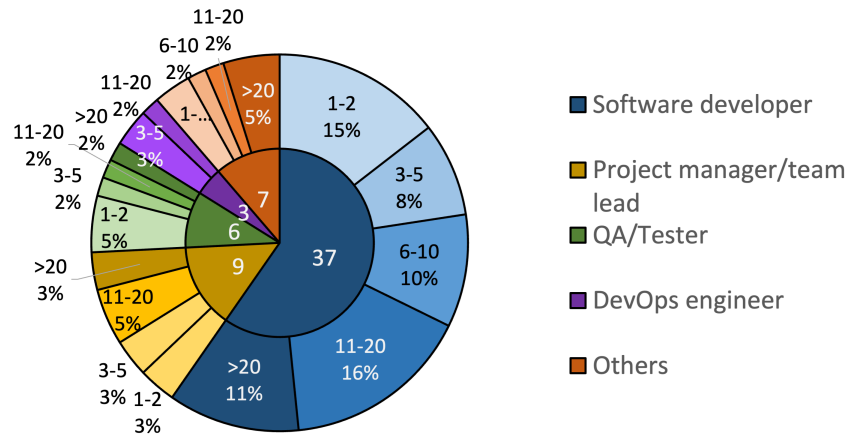


Figure 2. Demographics (role and experience in number of years) of participants

311 **Definition of software testability:** We asked the participants a question to elicit the definition of
 312 software testability. Most of the responses point to **the degree of ease with which automated**
 313 **tests can be written**. Some of the actual responses are: “The extent to which a software component
 314 can be tested”, “software testability is the degree that software artifacts support testing”, “easy testing”,
 315 and “a measure of how easy it is for the code to be tested through automated tests”. An interested
 316 reader may look at the raw anonymized responses in our replication package [57].

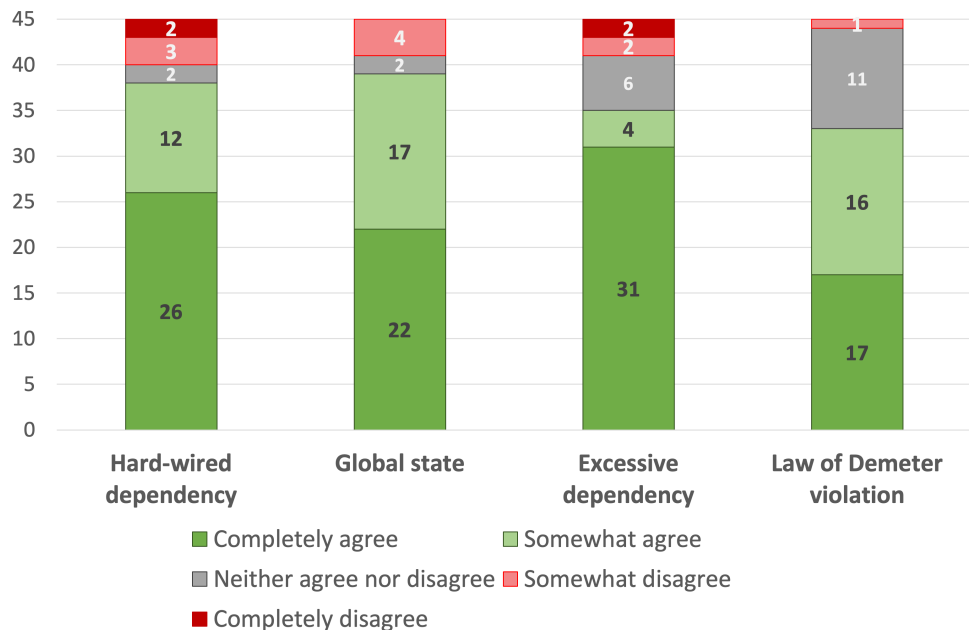


Figure 3. Respondents' perspective of considered testability smells

317 **Programming practices affecting testability:** The next set of questions presented four pro-
 318 gramming practices corresponding to each potential testability smell and asked the respondents
 319 whether and to what extent these practices negatively affect software testability. We also asked

320 about the rationale for their choice. Figure 3 presents the consolidated responses for all four con-
321 sidered smells. A very large percentage (84%, 87%, 78%, and 73% respectively for the four considered
322 smells) of the responses agreed (either completely or somewhat agree) to mark the presented prac-
323 tices as testability smells.

324 We looked into the rationale provided by respondents for other options (*i.e.*, neither agree nor
325 disagree, somewhat and completely disagree). Specifically, for *hard-wired dependency*, one of the
326 respondents who marked *completely disagree* did not offer any justification; another respondent
327 suggested to use mocking. One respondent with a *somewhat disagree* option for the same smell
328 basing his/her answer on the assumption that dependencies are trivial (*i.e.*, internal class or trivial
329 class from a library) most of the time. Respondents who opted for the option "*neither agree nor*
330 *disagree*" either expressed their ignorance about the specific question or left the rationale question
331 unanswered.

332 The respondents of "*somewhat disagree*" option for *global state* smell justify their selection by
333 providing a workaround to test a unit with global variables; for example, one of the respondents
334 provided the following rationale: "Logging where and when the global state is altered is usually
335 good enough for testing the code I work with".

336 Whereas, those who chose the "*completely disagree*" option for the *excessive dependency* smell,
337 seem, however, to agree that it is difficult to test source code containing this smell based on their
338 open answers (for example, one such an answer states "If designed properly then testing won't be
339 difficult but yes more dependencies need extra setups").

340 For the *law of Demeter violation* smell, a considerable number of respondents chose the "*nei-*
341 *ther agree nor disagree*" option; however, they did not provide any fruitful rationale towards this
342 testability smell.

343 **Additional programming practices affecting testability:** We also enquired about other pro-
344 gramming practices that may negatively influence the testability of a software system. The re-
345 sponses provided us with additional practices such as poor separation of concern (mixing `UI` and
346 non-`UI` aspects), interaction with external resources, such as sockets, files, and databases, time de-
347 pendencies, asynchronous operations, reflection, invoking command line from code, methods that
348 do not return anything but changes internal state, and requirements for authentication credentials
349 that hinder testability. In addition, the respondents mentioned spaghetti code, highly tangled code,
350 large methods, and non-standard environments as practices that reduce testability. Some of the
351 indicated practices are covered by the proposed smells. For example, interaction with external
352 resources has been captured by the *hard-wired dependency* smell since external resources such as
353 a network connection need to be instantiated.

354 The results from the survey not only suggest that the investigated smells are indeed considered
355 practices that affect testability negatively but also provide indicators for the community to extend
356 the proposed catalog.

357 4.3 The DesigniteJava Tool

358 We extended our tool `DESIGNITEJAVA` [56], to support for testability and test smells detection.⁴ We
359 select `DESIGNITEJAVA` to extend because the tool detects a variety of code smells and it has been
360 used in various studies [46, 58, 15, 68, 2]. Architecturally, `DESIGNITEJAVA` is structured in three layers
361 as shown in Figure 4. Eclipse Java Development Toolkit (`JDT`) forms the bottom layer. `DESIGNITEJAVA`
362 utilizes `JDT` to parse the source code, prepare `ASTs`, and resolve symbols. The source model is the
363 middle layer. The model invokes `JDT` and maintains a source code model from the information
364 extracted from an `AST` with the help of `JDT`. The top layer of the tool contains the business logic
365 *i.e.*, the smell detection and code quality metrics computation logic. The layer accesses the source
366 model, identifies smells and computes metrics, and outputs the generated information in either
367 `.CSV` or `.XML` files. Due to the existing support to detect smells and compute metrics, various fea-

⁴<https://www.designite-tools.com/blog/understanding-testability-test-smells>

368 tures (such as the source model) can be reused in our context. To support testability and test smell
 369 detection, we added code in the code smell detection layer. We also modified the source model
 370 layer to extract additional information required for our purpose. The extended version of the tool
 371 can be downloaded from its website.⁵

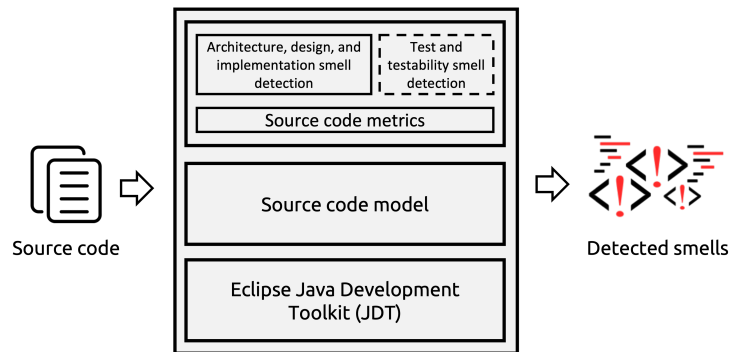


Figure 4. Architecture of DesigniteJava tool

372 Existing explorations have proposed a few tools to detect test smells. We first tried to utilize
 373 existing tools, specifically JNose [70] and TsDetect [51]. We were able to use JNose after taking help
 374 from its authors and developing a wrapper to use the tool as a console application. However, a
 375 quick analysis of the produced results showed a considerable number of false positive and false
 376 negative smell instances. Similarly, we were unable to use TsDetect because it is not suitable to
 377 analyze a large number of repositories due to a manual step requiring a mapping of test files and
 378 corresponding production files. Finally, we decided to develop our own test smell detector to iden-
 379 tify the following eight test smells---*Assertion roulette*, *Conditional test logic*, *Constructor initialization*,
 380 *Eager test*, *Empty test*, *Exception handling*, *Ignored test*, and *Unknown test*. We selected these smells
 381 because these were commonly known test smells and both the tools, *i.e.*, TsDetect and JNose, sup-
 382 port them. We implemented the support to detect test smells in DESIGNITEJAVA along with testability
 383 smells.

384 4.3.1 Detection rules for testability smells

385 We summarize below the detection strategies used for the testability smells.

386 **Hard-wired dependency:** We first detect all the objects created using the `new` operator in a class.
 387 Then, if the functionality of the newly created object is used (*i.e.*, at least one method is called) in
 388 the same class, we detect this smell.

389 **Global state:** If a class or a field in a class is declared with `public static` modifiers, we detect this
 390 smell.

391 **Excessive dependency:** We compute *fan-out* (*i.e.*, total number of outgoing dependencies) of a class.
 392 If the fan-out of the class is more than a pre-defined threshold, we detect the smell. The literature
 393 [47, 48, 76] suggests a threshold value for fan-out between 5 and 15 with a varying compliance rate.
 394 We adopted 7 as the threshold value as suggested by Arar et al. [76]. We ensure that the threshold
 395 value is configurable; hence, future studies may change any of the thresholds used.

396 **Law of Demeter violation:** We detect all the method invocation chains of the form `aField.get-`
 397 `Object().aMethod()`. We detect this smell when method calls are made on objects that are not
 398 directly associated with the current class.

399 4.3.2 Detection rules for test smells

400 The tool uses the definition of test smells and their detection strategies from existing studies [70,
 401 51, 3]. We present a summary of the detection strategies for the considered test smells below.

⁵<https://www.designite-tools.com/designitejava/>

402 **Assertion roulette:** We detect this smell when a test method contains more than one assertion
403 statement without giving an explanation as a parameter in the assertion method.
404 **Conditional test logic:** We detect this smell when there is an assertion statement within a control
405 statement block (e.g., `if` condition).
406 **Constructor initialization:** We detect this smell when a constructor of a test class initializes at least
407 one instance variable.
408 **Eager test:** We detect this smell when a test method calls multiple production methods.
409 **Empty test:** We detect this smell when a test method does not contain any executable statement
410 within its body.
411 **Exception handling:** We detect this smell when a test method asserts within a catch block or throws
412 an exception, instead of using `Assert.Throws()`.
413 **Ignored test:** We detect this smell when a test method is ignored using the `Ignore` annotation.
414 **Unknown test:** We detect this smell when a test method does not contain any `assert` call or ex-
415 pected exception.

416 4.3.3 Validation

417 We curated a *ground truth* of smells in a Java project to manually validate the tool, as explained
418 below.

419 **Subject system selection:** We used the `REPOREAPERS` dataset [42] and applied the following criteria
420 to select a subject system.

- 421 1. The repository must be implemented mainly in the Java programming language
- 422 2. The repository must be of moderate size (between 10K and 15K) to avoid toy projects on one
423 side and excessive manual effort on the other
- 424 3. The repository must have a unit-test ratio > 0.0 (number of `sLoC` in the test files to the number
425 of `sLoC` in all source files)
- 426 4. The repository must have a documentation ratio > 0.0 (number of comment lines of code to
427 the number of non-blank lines of source code)
- 428 5. The repository must have a community size > 1 (more than one developer).

429 We applied the criteria and sorted the list by the number of stars. We obtained *j256/ormlite-*
430 *jdbc*, *paul-hammant/paranamer*, and *forcedotcom/wsc* as the top three projects satisfying our criteria.
431 The majority of the source code belonging to *j256/ormlite-jdbc* and *paul-hammant/paranamer* was
432 in test cases. Hence, we selected *j256/ormlite-jdbc*,⁶ as our subject system for test smells valida-
433 tion. However, such repositories were not suitable for validating testability smells, since we detect
434 testability smells in non-test code. Hence, we selected *forcedotcom/wsc*,⁷ a project that offers a
435 high performance web service stack for clients, as our subject system for the manual validation of
436 testability smells.

437 **Validation protocol:** Two evaluators manually examined the source code of the selected subject
438 systems and documented the testability and test smells that they found. Both the evaluators hold a
439 PhD degree in computer science and have more than 5 years of software development experience.
440 Before carrying out the evaluation, they were introduced to testability and test smells. They were
441 allowed to use IDE features (such as “find”, “find usage” (of a variable) and “find definition” (of a class))
442 and external tools to collect code quality metrics to help them narrow their search space. Both
443 evaluators carried out their analyses independently. It took approximately three full work days
444 to complete the manual analysis. After their manual analysis was complete, they matched their
445 findings to spot any differences. We used *Cohen's Kappa* [6] to measure the inter-rater agreement
446 between the evaluators. The obtained result, 89% and 93% respectively for testability and test
447 smells, shows a strong agreement between the evaluators. The evaluators discussed the rest of
448 their findings and resolved the conflicts.

⁶<https://github.com/j256/ormlite-jdbc>

⁷<https://github.com/forcedotcom/wsc>

Table 1. Results of manual validation for testability smells

Testability Smells	Manually Verified Instances	TP	FP	FN
Hard-wired dependencies	64	63	2	1
Global state	22	22	0	0
Excessive dependencies	20	19	0	1
Law of Demeter violation	66	57	0	9
Total	172	161	2	11

Table 2. Results of manual validation for test smells

Testability Smells	Manually Verified Instances	TP	FP	FN
Assertion roulette	214	212	0	2
Conditional test logic	11	11	0	0
Constructor initialization	0	0	0	0
Eager test	13	13	0	0
Empty test	0	0	0	0
Exception handling	3	2	0	1
Ignored tests	2	2	0	0
Unknown test	58	58	0	0
Total	301	298	0	3

449 **Validation results:** We used our tool, `DESIGNITEJAVA`, on the subject systems and identified testability
 450 and test smells. We manually matched the ground truth prepared by the evaluators and the results
 451 produced by the tool. We classified each smell instance as true positive (TP), false positive (FP), and
 452 false negative (FN). We computed precision and recall metrics using the collected data.

453 Table 1 presents the results of the manual evaluation for testability smells. The tool identified
 454 161 instances of testability smells out of a total of 172 manually verified smell instances. The tool
 455 produced two false positive instances and eleven false negative instances. The false positive in-
 456 stances were detected mainly because the tool identified the *hard-wired dependency* even when an
 457 object was instantiated in a method call statement. Similarly, the tool reported false negatives due
 458 to an improper resolution of enumeration types; we traced back the inconsistent behavior to the
 459 `JDT` parser library. The precision and recall of the tool for testability smells based on the analysis
 460 is $161/(161 + 2) = 0.99$ and $161/(161 + 11) = 0.94$, respectively. Similarly, Table 2 shows the results
 461 of the manual evaluation carried out for test smells. Out of 301 test smells in 428 test methods,
 462 the tool correctly detected 298 smell instances. The cause of three instances of false negative is
 463 traced back to inconsistent behavior of the parser library. The precision and recall of the tool for
 464 test smells based on the analysis is $298/(298 + 0) = 1.0$ and $298/(298 + 3) = 0.99$, respectively. An
 465 interested reader may find the detailed manual analysis in our replication package [57].

466
 467 **Generalizability of conclusions:** The above validation shows that the tool produces reliable results
 468 in almost all cases. Given that the tool has been used by many researchers and practitioners,
 469 occasional issues reported by them were promptly fixed, thus further improving the reliability of
 470 the tool. A few known issues and limitations of the tool remain. First, due to a symbol resolution
 471 issue in `JDT`, in some very peculiar cases, the tool cannot resolve the symbol that leads to issues
 472 such as inability to identify the type of a variable. Also, the tool can identify test smells only when
 473 the tests are written in the JUnit framework.

Table 3. Characteristics of the analyzed repositories

Characteristics	Count
Total number of repositories	1,115
Total lines of code	46,176,914
Total number of classes	691,481
Total number of methods	4,031,216
Total number of test cases	415,527
Total number of testability smells	637,118

5. Mining GitHub repositories

We use the following mechanism to select and download repositories from GITHUB.

1. We use REPOREAPERS [42] to filter out low-quality and too small repositories on GitHub. We use quality characteristics provided by the REPOREAPERS to define a suitable criteria for repository selection. REPOREAPERS assesses repositories on eight characteristics and assigns a score typically between 0 and 1. We select all Java repositories in the REPOREAPERS dataset where architecture (as evidence of code organization), community and documentation (as evidence of collaboration), unit tests (as evidence of quality), history and issues (as evidence of accountability) scores are greater than zero. Further, we filter out repositories containing less than 1,000 lines of code (LOC) and having less than 10 stars.
2. We obtain 1,500 repositories after applying the above selection criteria.
3. We analyze all the selected repositories using the DESIGNITEJAVA tool that we developed to identify testability and test smells.

Table 3 presents the characteristics of the analyzed repositories. We attempted downloading and analyzing all the selected repositories; however, we could not download (either due to deleted or made private) and analyze (due to missing tests developed using JUnit framework) some of the repositories. Specifically, we did not find JUnit tests in 300 repositories. We successfully analyzed 1,115 repositories containing approximately 46 million LOC. Our replication package [57] includes the initial set of repositories, the names of all the successfully analyzed repositories along with the raw data generated by the employed tool, DESIGNITEJAVA.

6. Results

6.1 RQ1. To what extent do testability smells and test smells correlate?

6.1.1 Approach

The goal of this RQ is to explore the degree of correlation between test smells and testability smells in a repository. To achieve the above-stated goal, we first detect all the considered testability and test smells using DESIGNITEJAVA in the selected repositories. We calculate the sum of all testability smells and test smells per repository. Then, we compute smell density [58] to normalize the total number of smells to eliminate the potential confounding factor of project size. Testability smell density is defined as the total number of testability smells per one thousand lines of code (*i.e.*, (number of testability smells \times 1,000)/total lines of code). Test smell density is defined as the total number of test smells in each test method (*i.e.*, number of test smells/total number of tests). We use the Spearman's correlation coefficient [61] to measure the degree of association between these two smell types.

Furthermore, we explore the relationship at the class-level. By the fine-grained analysis, we aim to see whether a class *C* that suffers from testability smells shows a high number of test smells in the test cases that primarily test the class *C*, and vice-versa. Testability smells occur in production (non-test) code and test smells arise in test code. Hence, we require a mechanism to map a production class with corresponding test classes that test the production class. We implemented the logic

512 of identifying the production class under test for each test case in DESIGNITEJAVA. For the analysis,
513 we first find out all the method calls in each test case. Then, we identify the classes of the methods
514 that are called from the test case. It is possible that a test case calls methods belonging to multiple
515 classes; in that case, we attempt to identify the primary class that is being tested by the test case.
516 To do so, we match the name of the test class and the names of candidate primary classes; typically,
517 a test class is named by appending `Test` in the class name that the test class is testing. If the test
518 class name does not follow the specified pattern and there are multiple candidate classes to be
519 designated as the primary production class, then we pick the first candidate class whose method
520 is called from the test case. Using the above information, we prepare a reverse index mapping to
521 obtain a list of test cases corresponding to each production class. We use the mapping to retrieve
522 the number of test smells corresponding to each production class. As explained above, we com-
523 pute the testability smell density and test smell density at the class level. Finally, we compute the
524 Spearman's correlation coefficient between testability smell density and test smell density.

525 6.1.2 Results

526 Figure 5 shows the scatter plot between testability smell density and test smell density in the soft-
527 ware systems under examination. We obtain the Spearman's correlation coefficient $\rho = 0.246$
528 ($p\text{-value} < 2.2e - 16$); the coefficient indicates that testability and test smells share a low positive
529 correlation.

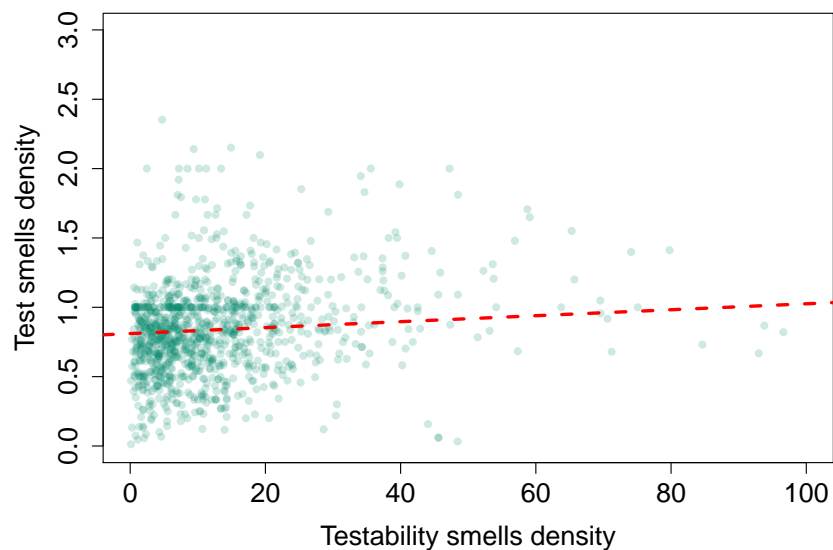


Figure 5. RQ1. Correlation between testability and test smell density

530 We extend our analysis by computing the correlation between the density of individual testabil-
531 ity smells and the test smell density per repository. We observe that *law of Demeter violation* shows
532 the highest correlation $\rho = 0.358$ ($p\text{-value} < 2.2e - 16$) with test smells. On the other hand, the *global*
533 *state* exhibits the lowest correlation $\rho = 0.076$ ($p\text{-value} < 2.2e - 16$). *Hard-wired dependency* and *ex-*
534 *cessive dependency* show correlation $\rho = 0.328$ ($p\text{-value} < 2.2e - 16$) and $\rho = 0.248$ ($p\text{-value} < 2.2e - 16$),
535 respectively.

536 We also investigate the relationship at the class-level. We identify the test cases and correspond-
537 ing test smells for each production class and compute the Spearman's correlation between the nor-
538 malized values of testability smells and test smells. We obtain $\rho = 0.050$ ($p\text{-value} = 2.903e - 08$). The
539 results indicate that testability smells and test smells do not share any correlation at the class-level

540 granularity.

Answer to RQ1. Testability smells show a low positive correlation with test smells. A fine-grained analysis at the class-level reveals that testability smells and test smells do not correlate with each other.

541

542 6.2 RQ2. Do testability smells correlate with test suite size?

543 6.2.1 Approach

544 RQ2 investigates whether and to what extent the presence of testability smells leads to fewer test
545 cases. To study this relationship, we first compute the testability smells in all the considered repos-
546 itories using DESIGNITEJAVA. In addition to smells, we use the tool to figure out the total number of
547 test cases in a repository; the tool marks each method as a test method or a normal non-test
548 method. For simplicity, we treat each test method (*i.e.*, a method with a `@Test` annotation) as a
549 test case. Next, we compute the *testability smell density* as described in RQ1 and the *test density* of
550 each repository. Test density is a normalized metric that represents the total number of test cases
551 per one thousand lines of code. We compute the Spearman's correlation coefficient between the
552 testability smell density and the test density for each repository.

553 Similar to RQ1, we explore the correlation at the class-level. For the analysis, as we explain
554 in RQ1, we first find out the production classes that a test case is testing. With this information,
555 we prepare a mapping between production classes and their corresponding set of test cases. We
556 use the mapping to retrieve the number of test smells corresponding to each production class.
557 We compute testability smell density and test case density at the class level. We measure the
558 correlation between testability smells and the number of test cases using Spearman's correlation
559 coefficient.

560 6.2.2 Results

561 Figure 6 presents a scatter plot between the testability smell density and the test density of the
562 selected repositories. We obtain $\rho = -0.033$ (p-value = 0.308), which is not statistically significant.
563 Therefore, testability and test smells do not correlate with each other.

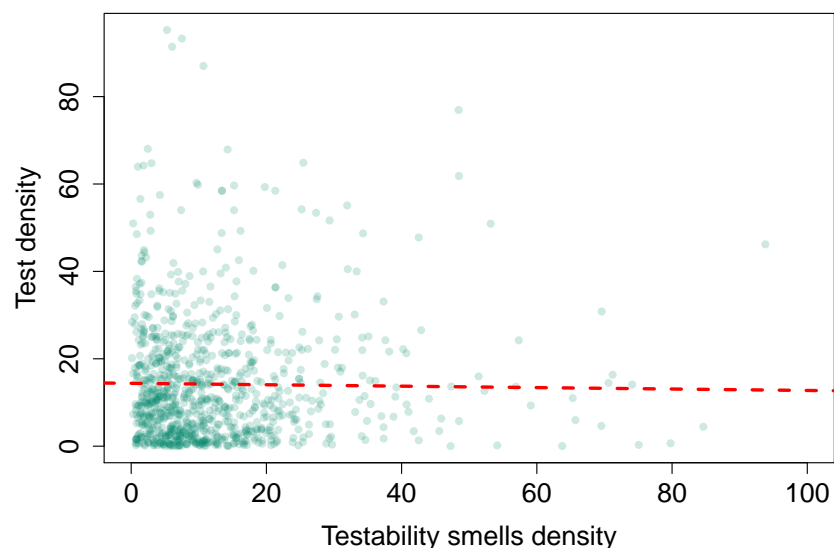


Figure 6. RQ2. Correlation of testability smells with test density

564 We extend our analysis by segregating the repositories into two categories by size. In the first
565 set, we put all the repositories that have less than 50,000 lines of code and, then, we put the rest
566 of the repositories in the second set. We carry out the same analysis on both of these sets. We
567 obtain $\rho = -0.009$ (p-value = 0.789) for the first set and $\rho = -0.050$ (p-value = 0.492) for the second
568 set between testability smells and test density. The obtained results are not statistically significant.

569 Furthermore, we observe the relationship between testability smells and the number of tests
570 at the class-level granularity. We compute the total number of testability smells for each non-test
571 class and figure out the total number of tests written for the class. In the computation, we did not
572 include the classes where the number of tests for the entire project is zero indicating that either the
573 test cases are not written for the project or the test cases are written using a framework other than
574 JUnit. We perform the above step to reduce the noise in the prepared data. We obtain $\rho = -0.179$
575 (p-value < $2.2e - 16$) as the correlation coefficient. The results clearly show that testability smells
576 show a very low correlation with the size of the test suite.

Answer to RQ2. Testability smells do not exhibit any correlation with the test density of a software system.

577

578 6.3 RQ3. Do testability smells cause more bugs?

579 6.3.1 Approach

580 RQ3 aims to investigate whether and to what extent testability smells relate to, and even cause,
581 bugs in a given software system. To answer this question, we choose five subject systems manually
582 and perform a trend analysis by extracting information from multiple commits for each of these
583 subject systems.

584 We use the following protocol to identify the subject systems for this research question. First,
585 we obtain a sorted list of repositories by their number of commits in descending order from our
586 selected initial set of repositories (see Section 5) using the GITHUB API. The intent here is to choose
587 repositories with a rich commit history to facilitate detailed trend analysis. Then, we manually
588 check these repositories one by one to assess whether a repository uses GITHUB issues and whether
589 these issues are labelled as “bugs”. In addition, we execute DESIGNITEJAVA on the latest commit of
590 each of these repositories to ensure that it does not take too long to run, as, otherwise, it might be
591 prohibitive for us to run it to analyze the entire repository containing multiple (hundreds of)
592 commits. Finally, we select the first five repositories that satisfy the above criteria, which are: Magarena,⁸
593 XP,⁹ Rundeck,¹⁰ MyRobotLab,¹¹ and Ontrack.¹²

594 In order to perform a trend analysis, it is crucial to select a suitable set of commits from each
595 of these repositories. One common way is to select commits at a fixed interval either by com-
596 mit number (for example, every 100th commit) or by commit date (for example, one commit per
597 month). However, such a mechanism may result in a skewed set of commits where either sig-
598 nificant changes in the commits are missed or commits with hardly any change are analyzed. To
599 overcome this limitation, Sharma *et al.* [58] proposed a commit selection algorithm where commits
600 are selected based on the amount of changes introduced in a commit *w.r.t.* the previous selected
601 commit. In this work, we follow this strategy to select commits for each of the five identified repos-
602 itories. Specifically, we first obtain all the commits in a repository in the main branch, and then we
603 choose the first and the last commit to get started. Then, we compute five code quality metrics
604 (*i.e.*, weighted methods per class (wmc), number of children (NC), lack of cohesion among methods
605 (lcom), fan-in, and fan-out) and identify changed classes based on the changes in any of these met-
606 rics. If the changed number of classes between two analyzed commits differ by a threshold (set

⁸<https://github.com/magarena/magarena>

⁹<https://github.com/enonic/xp>

¹⁰<https://github.com/rundeck/rundeck>

¹¹<https://github.com/MyRobotLab/myrobotlab>

¹²<https://github.com/nemerosa/ontrack>

607 to 5%), we consider the commit having significant changes [58]. We then pick the middle commit
608 (*i.e.*, the commit between the currently selected two commits) and repeat the process until we find
609 commits with non-significant changes [58].

610 Once we identify the set of commits for the trend analysis, we detect testability smells in each of
611 the selected commits for each repository by using our `DESIGNITEJAVA` tool. Also, we identify the total
612 number of open and closed issues that has the tag “bugs” when the commit was made. The GitHub
613 `API` does not provide a direct way to figure out issues at the time of a specific commit. To identify
614 issues at the time of a specific commit, we first fetch the issues that are open (or closed) since the
615 time of a commit and subtract them from the total open (or closed) issues at present. It gives us the
616 total number of open (or closed) issues at the time of a specific commit. We record this information
617 along with the total detected testability smells for each selected commit. Using this information,
618 we compute the Spearman's correlation coefficient between the total detected testability smells
619 and the total number of issues (*i.e.*, the sum of open and closed for each considered commit).

620 We also carry out a causal analysis to figure out whether testability smells *cause* bugs. We use
621 Granger's causality [24] analysis for this purpose. The method has been used in similar studies [10,
622 49, 58] to explore the causal relationship within the software engineering domain. Equation 1
623 presents Granger's method mathematically.

$$a(t) = \sum_{j=1}^k f(s_{t-j}) + \sum_{j=1}^k f(b_{t-j}) \quad (1)$$

624 In our context, time series S and B represent the testability smell density *i.e.*, total number of
625 testability smells per one thousand lines of code, and reported bugs computed over a period of
626 time. Variables s_t and b_t represent testability smell instances and total reported bugs at time t . If
627 the predictions of variable b with the past values of both s and b are better than the predictions
628 using only the past values of b , then testability smells *cause* the bugs.

629 In such analysis, we must ensure the *stationary* property of a time-series before analyzing it
630 and drawing conclusions based on that. A time-series is stationary if its statistical properties, such
631 as mean, variance, and autocorrelation, are constant over time [11]. A non-stationary time-series
632 shows seasonal effects, trends, and fluctuating statistical properties changing over time. Such
633 effects are undesired for the causality analysis and thus a time-series must be made stationary
634 before we perform the causality analysis. We carried out the augmented Dickey-Fuller unit root
635 test [20] to check the stationary property of our time-series. Initially, our time-series was non-
636 stationary. There are a few techniques to make a non-stationary time-series a stationary one [35].
637 We addressed this issue by applying a *difference* transformation, *i.e.*, subtracting the previous ob-
638 servation from the present observation for all columns. Techniques such as *differencing*, that we
639 applied, help stabilize the mean of a time series by removing changes in the time series, and there-
640 fore eliminate or reduce the non-stationary nature of the series [35]. After the transformation, we
641 obtain a stationary time-series that we confirmed by performing the augmented Dickey-Fuller unit
642 root test again. Finally, we carry out the causality analysis using Equation 1.

643 6.3.2 Results

644 Table 4 presents the results of the experiment. The number of analyzed commits ranges between
645 38 (for `XP`) and 180 (for `Rundeck`). The size of the selected repositories varies between ≈ 71 `KLOC` (for
646 `Magarena`) to ≈ 181 `KLOC` (for `XP`) as measured for the most recent analyzed commit. We compute
647 the total number of testability smells in each selected commit as well as the total reported (open
648 and closed) issues marked as bugs at the time of the corresponding commits for each of the se-
649 lected repositories individually. The table shows the total number of testability smells detected in
650 the most recent analyzed commit. We compute the Spearman correlation coefficient between the
651 reported issues and testability smell density. We observe mixed results for the correlation analy-
652 sis; two repositories show strong, one repository shows moderate, and one repository shows low

Table 4. RQ3. Correlation and causation relationships between testability smell density with the number of reported bugs

Repository	#Commits	LOC	Testability Smells	Correlation Coefficient (p-value)	Causality p-value
Magarena	66	71,567	1,425	0.482(4.4e - 4)	0.119
MyRobotlab	76	118,532	2,643	0.761(< 1.4e - 15)	0.661
Ontrack	107	17,009	72	0.105(< 0.280)	0.708
Rundeck	180	81,198	1,570	0.193(< 0.009)	0.825
XP	38	181,278	2,143	0.937(< 2.2e - 16)	0.249

653 correlation. We observe that the correlation coefficient is not statistically significant for the *Ontrack*
 654 repository.

655 The last column of Table 4 presents the results of the causality test. Each cell in the column
 656 shows the p-value computed for the causal relationship of testability smells with the reported bugs.
 657 The results for all the analyzed repositories show that testability smells *do not* cause bugs as all
 658 the obtained p-values are greater than 0.05.

Answer to RQ3: The causality analysis reveals that **testability smells do not cause bugs.**

659

660 7. Implications and Discussion

661 The results of our first research question reveal that there is no correlation between testability
 662 smells and test smells. This suggests that writing good-quality test code is possible even with poor
 663 testability, at least for all testability smells considered herein. The results also indicate that either
 664 the difficulty in writing tests due to the considered testability smells is orthogonal to test smells, or
 665 existing testing frameworks, *e.g.*, mocking frameworks, make it easier to overcome the challenges
 666 posed by poor testability. Researchers may investigate further the influence of tools' features, such
 667 as mocking, to facilitate testing despite poor testability.

668 We explore the effect of testability smells on test suite size, represented by the number of
 669 test cases, in RQ2. Figure 7 shows box-plots of the categories of testability smell density with test
 670 density. We divide the repositories into four categories C1 to C4 based on the value of the testability
 671 smell density. For example, repositories with a testability smell density of less than five are put into
 672 category C1. We observe that the median test density for the first category is the highest among all
 673 the categories and the test density dips for the category C2. However, against the common belief
 674 of developers, test density rises again in categories C3 and C4. **The analysis further reaffirms
 675 that testability smells do not share a linear monotonic relationship with test density.**

676 Our experiment to investigate the correlation of testability smell density with the number of
 677 reported bugs does not show a consistent strong relationship. The strong correlation in two repos-
 678 itories and the moderate correlation in a repository show that the density of testability smell in-
 679 creases as the size of the software grows since the total number of reported bugs always increases
 680 with time. Hence, a strong correlation implies that the rate of testability smells increases as the
 681 software systems grow.

682 In the context of our study, one might wonder about the relationship of testability smells with
 683 traditional code smells. Given the definition and scope, it is likely that some code smells are also
 684 considered testability smells if they impact testability. However, this interpretation is not uniquely

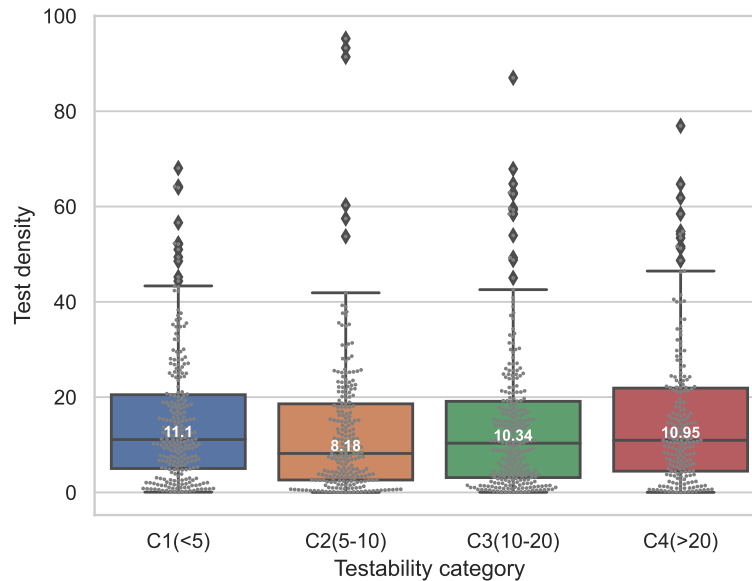


Figure 7. Box-plots of the categories of testability smell density with test density

685 applicable only in this context. For example, a violation of the ‘single responsibility principle’ may
 686 introduce incohesive class (or multifaceted abstraction) at design and ‘feature concentration’ smell
 687 at architecture granularity. Nevertheless, we perform an analysis of testability smell density with
 688 code smell density not only at the repository-level but also at a fine-grained granularity of class-
 689 level (where we compute the total number of smells for each class of the considered repositories).
 690 We use the `DESIGNITEJAVA` tool to detect code smells and testability smells. We compute the Spear-
 691 man’s correlation coefficient between the normalized total number of smells. At the repository-
 692 level, we obtain $\rho = 0.851$ (p-value $< 2.2e - 16$) as the Spearman’s correlation coefficient. Similarly,
 693 we get $\rho = 0.857$ (p-value $< 2.2e - 16$) when we compute the correlation at the class-level. The strong
 694 correlation indicates that the presence of a large population of code smells is associated with the
 695 presence of a large number of testability smells and vice versa.

696 The elicited developers’ perspective clearly emphasizes the importance of testability smells and
 697 the potential negative impact on testing aspects. However, the empirical evidence observed in
 698 the study does not agree with the perspective. We observed that testability smells, at the class-
 699 level fine-grained granularity, do not correlate with test smells. Also, the smells do not show any
 700 influence on test density. Furthermore, the results show that testability smells do not contribute
 701 to a higher number of bugs. The results suggest that despite developers’ invaluable experience,
 702 their opinions and perspectives might need to be complemented with empirical evidence before
 703 bringing it into practice.

704 8. Threats to Validity

705 This section discusses the potential threats to the validity (construct, internal, and external) of our
 706 reported results.

707 **Construct validity.** Construct threats to validity are concerned with the degree to which our
 708 analyses measure what we claim to analyze. In our study, we used our `DESIGNITEJAVA` tool to iden-
 709 tify the four testability smells. However, the strategies used to identify testability smells may not
 710 capture all testability cases. To mitigate this threat, we thoroughly tested the tool using different
 711 cases for each smell, and also fine-tuned the tool based on testing. Then, we performed a manual
 712 analysis of the four testability smells on a complete project, namely `wsc`. The results of the man-

713 ual validation show a very high recall and precision. Similarly, we also implemented support to
714 detect test smells by following detection strategies proposed in the existing literature to identify
715 test smells.

716 **Internal validity.** Internal threats to validity are concerned with the ability to draw conclusions
717 from our experimental results. We carried out an online anonymous survey targeting develop-
718 ers by posting our survey on social media professional channels (Twitter, LinkedIn, and Reddit).
719 Given the anonymity of the survey, we do not have any mechanism to verify the level of experi-
720 ence claimed by the participants. However, based on the quality of the responses provided by the
721 participants, we believe that such a threat is mild. In addition, software developers participated in
722 our online survey were not selected based on the repositories we analyzed. As a result, opinions of
723 developers could be influenced by the repositories they usually contribute to and might not agree
724 with our empirical results. To mitigate this, we did not target developers from specific repositories
725 but rather expanded our participation range by posting invitations on online professional social
726 media channels.

727 *Agreement bias (or acquiescence bias)* refers to the participants' tendency to agree with a state-
728 ment rather than disagreeing with it [66]. We design our questionnaire in a neutral tone and pro-
729 vide options by using a Likert-scale to mitigate this threat. A similar threats to validity is partic-
730 ipants' acquaintance to the authors. To avoid this threat, we did not circulate the survey in our
731 internal organization groups. Also, we restricted the sharing to *professional* social media channels
732 and hence did not sharing the survey on our, for example, Facebook profiles or groups.

733 RQ3 investigates causality between testability smells and the number of reported bugs; the
734 analysis reveals that the testability smells do not cause bugs. There are two possible threats to the
735 conclusion. First, it is possible that the testability smells other than those considered in the study
736 have a larger impact on bugs. However, though there could be many other testability smells, the
737 considered smells are representative as shown by our developers' survey. Second, the study only
738 considered reported known bugs. It is possible that there are many more unknown bugs that may
739 influence the results and conclusion of the experiment.

740 **External validity** External threats are concerned with the ability to generalize our results. The
741 1, 115 GitHub repositories analyzed in this paper were selected using well-defined criteria from the
742 REPOREAPERS dataset. However, some repositories might have switched from *public* to *private* or no
743 longer exist on GITHUB, which might affect the criteria used to select repositories in this paper. In
744 addition, all the selected repositories contain software written in Java, which might affect the gener-
745 alizability of our findings. The major reason for focusing on Java is that the majority of research on
746 software quality analysis has been done on Java code, and hence we can leverage existing tools to
747 achieve the goals of our study. Along the same lines, the implemented test smell detection works
748 only if the tests are written using JUnit. The rationale behind this decision is that JUnit is the most
749 commonly used testing framework for Java. We encourage future research to expand the analyses
750 conducted in this paper to software written in different programming languages.

751 9. Conclusions

752 This study explores practitioners' perspectives about testability smells as well as experimental evi-
753 dence gathered via a large-scale empirical study on 1, 115 Java repositories containing approxi-
754 mately 46 million lines of code in order to better understand the relationship of testability smells
755 with test quality, number of tests, and reported bugs.

756 Specifically, the study surveyed software developers to elicit their opinions and perspectives
757 about testability smells. The survey showed that software developers consider testability a factor
758 that impedes software testing; the survey also revealed their strong agreement with the proposed
759 testability smells. Then, we conducted an extensive empirical evaluation to observe the relation-
760 ship between testability smells and test-related aspects such as test smells and test suit size. Our
761 results show that testability smells do not correlate with test smells at the class granularity and

762 with test suit size. Furthermore, we did not find evidence that testability smells cause bugs.

763 Our study has implications for both the research and industrial communities. Software de-
764 velopers often have strong opinions about software engineering concepts; however, experimental
765 evidence may not support them in general. Specifically, this study shows that developers' opinions
766 about testability do not concur with the experimental evidence. Hence, opinions and perspectives
767 must be complemented with empirical evidence before bringing into practice. This also highlight
768 the importance of data-driven software engineering, which advocates the need and value of adopt-
769 ing design and development decisions supported by data. Researchers can use our tool to detect
770 testability smells to further evaluate and confirm our observations. Also, researchers may pro-
771 pose additional testability smells and investigate their collective impact on other relevant testing
772 aspects, such as testing efforts.

773 References

- 774 [1] A. A. Al-Subaihin, F. Sarro, S. Black, L. Capra, and M. Harman. App store effects on software
775 engineering practices. *IEEE Transactions on Software Engineering*, 47(2):300--319, 2021.
- 776 [2] M. Alenezi and M. Zarour. An empirical study of bad smells during software evolution using
777 designite tool. *i-Manager's Journal on Software Engineering*, 12(4):12, 2018.
- 778 [3] W. Aljedaani, A. Peruma, A. Aljohani, M. Alotaibi, M. W. Mkaouer, A. Ouni, C. D. Newman,
779 A. Ghallab, and S. Ludi. Test smell detection tools: A systematic mapping study. In *Evalu-
780 ation and Assessment in Software Engineering*, EASE 2021, page 170–180, New York, NY, USA,
781 2021. Association for Computing Machinery.
- 782 [4] G. Bavota, A. Qusef, R. Oliveto, A. De Lucia, and D. Binkley. An empirical analysis of the dis-
783 tribution of unit test smells and their impact on software maintenance. In *2012 28th IEEE
784 International Conference on Software Maintenance (ICSM)*, pages 56--65, 2012.
- 785 [5] G. Bavota, A. Qusef, R. Oliveto, A. Lucia, and D. Binkley. Are test smells really harmful? an
786 empirical study. *Empirical Softw. Engg.*, 20(4):1052–1094, aug 2015.
- 787 [6] K. J. Berry and J. Paul W. Mielke. A Generalization of Cohen's Kappa Agreement Measure
788 to Interval Measurement and Multiple Raters. *Educational and Psychological Measurement*,
789 48(4):921--933, 1988. _eprint: <https://doi.org/10.1177/0013164488484007>.
- 790 [7] R. V. Binder. Design for testability in object-oriented systems. *Commun. ACM*, 37(9):87–101,
791 sep 1994.
- 792 [8] M. Bruntink and A. van Deursen. An empirical study into class testability. *Journal of Systems
793 and Software*, 79(9):1219--1232, Sept. 2006.
- 794 [9] V. Chowdhary. Practicing testability in the real world. In *2009 International Conference on
795 Software Testing Verification and Validation*, pages 260--268, 2009.
- 796 [10] C. Couto, P. Pires, M. T. Valente, R. da Silva Bigonha, A. C. Hora, and N. Anquetil. Bugmaps-
797 granger: A tool for causality analysis between source code metrics and bugs. 2013.
- 798 [11] D. Cox and H. Miller. *The Theory of Stochastic Process*. Chapman and Hall, London, 1 edition,
799 1965.
- 800 [12] A. V. Deursen, L. Moonen, A. Bergh, and G. Kok. Refactoring test code. In *Proceedings of the 2nd
801 International Conference on Extreme Programming and Flexible Processes in Software Engineering
802 (XP2001)*, pages 92--95, 2001.

- 803 [13] P. Devanbu, T. Zimmermann, and C. Bird. Belief & evidence in empirical software engineer-
804 ing. In *Proceedings of the 38th International Conference on Software Engineering, ICSE '16*, page
805 108–119, New York, NY, USA, 2016. Association for Computing Machinery.
- 806 [14] M. Eck, F. Palomba, M. Castelluccio, and A. Bacchelli. Understanding flaky tests: The devel-
807 oper’s perspective. In *Proceedings of the 2019 27th ACM Joint Meeting on European Software
808 Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE
809 2019*, page 830–840, New York, NY, USA, 2019. Association for Computing Machinery.
- 810 [15] A. Eposhi, W. Oizumi, A. Garcia, L. Sousa, R. Oliveira, and A. Oliveira. Removal of design prob-
811 lems through refactorings: Are we looking at the right symptoms? In *2019 IEEE/ACM 27th
812 International Conference on Program Comprehension (ICPC)*, pages 148–153, 2019.
- 813 [16] S. Fatima, T. A. Ghaleb, and L. Briand. Flakify: A black-box, language model-based predictor
814 for flaky tests. *IEEE Transactions on Software Engineering*, 2022.
- 815 [17] M. Feathers. *Working Effectively with Legacy Code: WORK EFFECT LEG CODE _p1*. Prentice Hall
816 Professional, 2004.
- 817 [18] F. G. S. Filho, V. Lelli, I. d. S. Santos, and R. M. C. Andrade. Correlations among software
818 testability metrics. In *19th Brazilian Symposium on Software Quality, SBQS'20*, New York, NY,
819 USA, 2020. Association for Computing Machinery.
- 820 [19] R. Freedman. Testability of software components. *IEEE Transactions on Software Engineering*,
821 17(6):553–564, 1991.
- 822 [20] W. A. Fuller. *Introduction to Statistical Time Series*. John Wiley and Sons New York, 1 edition,
823 1976.
- 824 [21] V. Garousi, M. Felderer, and F. N. Kılıçaslan. A survey on software testability. *Information and
825 Software Technology*, 108:35–64, 2019.
- 826 [22] V. Garousi, M. Felderer, and M. V. Mäntylä. Guidelines for including grey literature and con-
827 ducting multivocal literature reviews in software engineering, 2018.
- 828 [23] V. Garousi and B. Küçük. Smells in software test code: A survey of knowledge in industry and
829 academia. *Journal of Systems and Software*, 138:52–81, 2018.
- 830 [24] C. W. J. Granger. Investigating causal relations by econometric models and cross-spectral
831 methods. *Econometrica*, 37(3):424–438, 1969.
- 832 [25] M. M. Hassan, W. Afzal, M. Blom, B. Lindström, S. F. Andler, and S. Eldh. Testability and soft-
833 ware robustness: A systematic literature review. In *2015 41st Euromicro Conference on Software
834 Engineering and Advanced Applications*, pages 341–348, 2015.
- 835 [26] M. Hevery. Writing Testable Code, Aug. 2008. [https://testing.googleblog.com/2008/08/
836 by-miko-hevery-so-you-decided-to.html](https://testing.googleblog.com/2008/08/by-miko-hevery-so-you-decided-to.html).
- 837 [27] M. Human. Why You Should Be Replacing Full Stack Tests with Ember Tests, Dec. 2022. [https:
838 //www.mutuallyhuman.com/blog/why-you-should-be-replacing-full-stack-tests-with-ember-tests/](https://www.mutuallyhuman.com/blog/why-you-should-be-replacing-full-stack-tests-with-ember-tests/).
- 839 [28] A. A. Janes and G. Succi. The dark side of agile software development. In *Proceedings of the
840 ACM International Symposium on New Ideas, New Paradigms, and Reflections on Programming and
841 Software, Onward! 2012*, page 215–228, New York, NY, USA, 2012. Association for Computing
842 Machinery.
- 843 [29] V. Jeffrey M. Factors that affect software testability. Technical report, 1991.

- 844 [30] N. S. Junior, L. Rocha, L. A. Martins, and I. Machado. A survey on test practitioners' awareness
845 of test smells, 2020.
- 846 [31] T. Kaczanowski. *Practical Unit Testing with JUnit and Mockito*. Tomasz Kaczanowski, 2013.
- 847 [32] R. A. Khan and K. Mustafa. Metric based testability model for object oriented design (mtmood).
848 *SIGSOFT Softw. Eng. Notes*, 34(2):1–6, feb 2009.
- 849 [33] D. J. Kim, T.-H. P. Chen, and J. Yang. The secret life of test smells - an empirical study on test
850 smell evolution and maintenance. *Empirical Software Engineering*, 26(5):100, July 2021.
- 851 [34] R. Kolb and D. Muthig. Making testing product lines more efficient by improving the testability
852 of product line architectures. In *Proceedings of the ISSTA 2006 Workshop on Role of Software
853 Architecture for Testing and Analysis*, ROSATEA '06, pages 22---27, New York, NY, USA, 2006.
854 Association for Computing Machinery.
- 855 [35] D. Kwiatkowski, P. C. Phillips, P. Schmidt, and Y. Shin. Testing the null hypothesis of stationarity
856 against the alternative of a unit root: How sure are we that economic time series have a unit
857 root? *Journal of Econometrics*, 54(1):159 -- 178, 1992.
- 858 [36] Y. Le Traon and C. Robach. From hardware to software testability. In *Proceedings of 1995 IEEE
859 International Test Conference (ITC)*, pages 710--719, 1995.
- 860 [37] Y. Le Traon and C. Robach. Testability measurements for data flow designs. In *Proceedings
861 Fourth International Software Metrics Symposium*, pages 91--98, 1997.
- 862 [38] K. J. Lienberherr. Formulations and benefits of the law of demeter. *SIGPLAN Not.*, 24(3):67–78,
863 mar 1989.
- 864 [39] B. Lo and H. Shi. A preliminary testability model for object-oriented software. In *Proceedings.
865 1998 International Conference Software Engineering: Education and Practice (Cat. No.98EX220)*,
866 pages 330--337, 1998.
- 867 [40] L. Marshall and J. Webber. Gotos considered harmful and other programmers' taboos. *De-
868 partment of Computing Science Technical Report Series*, 2000.
- 869 [41] S. Mouchawrab, L. C. Briand, and Y. Labiche. A measurement framework for object-oriented
870 software testability. *Information and Software Technology*, 47(15):979--997, 2005. Most Cited
871 Journal Articles in Software Engineering - 1999.
- 872 [42] N. Munaiah, S. Kroh, C. Cabrey, and M. Nagappan. Curating GitHub for engineered software
873 projects. *Empirical Software Engineering*, 22(6):3219--3253, Dec. 2017.
- 874 [43] E. Murphy-Hill, C. Parnin, and A. P. Black. How we refactor, and how we know it. *IEEE Transac-
875 tions on Software Engineering*, 38(1):5--18, 2012.
- 876 [44] T. Nguyen, M. Delaunay, and C. Robach. Testability analysis for software components. In
877 *International Conference on Software Maintenance, 2002. Proceedings.*, pages 422--429, 2002.
- 878 [45] T. B. Nguyen, M. Delaunay, and C. Robach. Testability Analysis of Data-Flow Software. *Elec-
879 tronic Notes in Theoretical Computer Science*, 116:213--225, Jan. 2005.
- 880 [46] W. Oizumi, L. Sousa, A. Oliveira, L. Carvalho, A. Garcia, T. Colanzi, and R. Oliveira. On the density
881 and diversity of degradation symptoms in refactored classes: A multi-case study. In *2019 IEEE
882 30th International Symposium on Software Reliability Engineering (ISSRE)*, pages 346--357. IEEE,
883 2019.

- 884 [47] P. Oliveira, F. P. Lima, M. T. Valente, and A. Serebrenik. Rtttool: A tool for extracting relative
885 thresholds for source code metrics. In *2014 IEEE International Conference on Software Maintenance and Evolution*, pages 629--632, 2014.
886
- 887 [48] P. Oliveira, M. T. Valente, and F. P. Lima. Extracting relative thresholds for source code metrics.
888 In *2014 Software Evolution Week - IEEE Conference on Software Maintenance, Reengineering, and*
889 *Reverse Engineering (CSMR-WCRE)*, pages 254--263, 2014.
- 890 [49] F. Palomba, G. Bavota, M. D. Penta, F. Fasano, R. Oliveto, and A. D. Lucia. A large-scale empirical
891 study on the lifecycle of code smell co-occurrences. *Information and Software Technology*, 99:1
892 -- 10, 2018.
- 893 [50] J. E. Payne, R. T. Alexander, and C. D. Hutchinson. Design-for-testability for object-oriented
894 software. *Object Magazine*, 7(5):34--43, 1997.
- 895 [51] A. Peruma, K. Almalki, C. D. Newman, M. W. Mkaouer, A. Ouni, and F. Palomba. Tsdetect:
896 An open source test smells detection tool. In *Proceedings of the 28th ACM Joint Meeting on*
897 *European Software Engineering Conference and Symposium on the Foundations of Software Engi-*
898 *neering, ESEC/FSE 2020*, page 1650--1654, New York, NY, USA, 2020. Association for Computing
899 Machinery.
- 900 [52] B. Pettichord. Design for testability. In *Pacific Northwest Software Quality Conference*, pages
901 1--28, 2002.
- 902 [53] D. Pina, C. Seaman, and A. Goldman. Technical debt prioritization: A developer's perspective.
903 In *Proceedings of the International Conference on Technical Debt*, TechDebt '22, page 46--55, New
904 York, NY, USA, 2022. Association for Computing Machinery.
- 905 [54] F. Rahman, C. Bird, and P. Devanbu. Clones: What is that smell? In *2010 7th IEEE Working*
906 *Conference on Mining Software Repositories (MSR 2010)*, pages 72--81, 2010.
- 907 [55] D. M. Ribeiro, F. Q. B. da Silva, D. Valença, E. L. S. X. Freitas, and C. França. Advantages and
908 disadvantages of using shared code from the developers perspective: A qualitative study. In
909 *Proceedings of the 10th ACM/IEEE International Symposium on Empirical Software Engineering and*
910 *Measurement, ESEM '16*, New York, NY, USA, 2016. Association for Computing Machinery.
- 911 [56] T. Sharma. DesigniteJava, Dec. 2018. <https://github.com/tushartushar/DesigniteJava>.
- 912 [57] T. Sharma, S. Georgiou, M. Kechagia, T. A. Ghaleb, and F. Sarro. Replication Package for Testa-
913 bility Study, Nov. 2022. <https://github.com/SMART-Dal/testability>.
- 914 [58] T. Sharma, P. Singh, and D. Spinellis. An empirical investigation on the relationship between
915 design and architecture smells. *Empirical Software Engineering*, 25(5):4020--4068, 2020.
- 916 [59] P. K. Singh, O. P. Sangwan, A. P. Singh, and A. Pratap. An assessment of software testability
917 using fuzzy logic technique for aspect-oriented software. *International Journal of Information*
918 *Technology and Computer Science (IJITCS)*, 7(3):18, 2015.
- 919 [60] D. Spadini, F. Palomba, A. Zaidman, M. Bruntink, and A. Bacchelli. On the relation of test smells
920 to software code quality. In *2018 IEEE International Conference on Software Maintenance and*
921 *Evolution (ICSME)*, pages 1--12, 2018.
- 922 [61] C. Spearman. The proof and measurement of association between two things. 1961.
- 923 [62] G. Suryanarayana, G. Samarthyam, and T. Sharma. *Refactoring for Software Design Smells: Man-*
924 *aging Technical Debt*. Morgan Kaufmann, 1 edition, 2014.

- 925 [63] R. E. Sward and A. Chamillard. Re-engineering global variables in ada. In *Proceedings of the*
926 *2004 annual ACM SIGAda international conference on Ada: The engineering of correct and reliable*
927 *software for real-time & distributed systems using Ada and related technologies*, pages 29--34,
928 2004.
- 929 [64] V. Terragni, P. Salza, and M. Pezzè. Measuring software testability modulo test quality. In
930 *Proceedings of the 28th International Conference on Program Comprehension*, ICPC '20, page
931 241–251, 2020.
- 932 [65] D. Thomas and A. Hunt. *The Pragmatic Programmer: your journey to mastery*. Addison-Wesley
933 Professional, 2019.
- 934 [66] B. Toner. The impact of agreement bias on the ranking of questionnaire response. *The Journal*
935 *of Social Psychology*, 127(2):221--222, 1987.
- 936 [67] M. Tufano, F. Palomba, G. Bavota, M. Di Penta, R. Oliveto, A. De Lucia, and D. Shybyanyk.
937 An empirical investigation into the nature of test smells. In *Proceedings of the 31st IEEE/ACM*
938 *International Conference on Automated Software Engineering, ASE '16*, pages 4---15, New York,
939 NY, USA, 2016. Association for Computing Machinery.
- 940 [68] A. Uchôa, C. Barbosa, W. Oizumi, P. Blenilio, R. Lima, A. Garcia, and C. Bezerra. How does mod-
941 ern code review impact software design degradation? an in-depth empirical study. In *2020*
942 *IEEE International Conference on Software Maintenance and Evolution (ICSME)*, pages 511--522,
943 2020.
- 944 [69] J. Vincent, G. King, P. Lay, and J. Kinghorn. Principles of Built-In-Test for Run-Time-Testability
945 in Component-Based Software Systems. *Software Quality Journal*, 10(2):115--133, Sept. 2002.
- 946 [70] T. Virgínio, L. Martins, L. Rocha, R. Santana, A. Cruz, H. Costa, and I. Machado. Jnose: Java test
947 smell detector. In *Proceedings of the 34th Brazilian Symposium on Software Engineering, SBES*
948 '20, page 564–569, New York, NY, USA, 2020. Association for Computing Machinery.
- 949 [71] J. M. Voas. *Object-Oriented Software Testability*, pages 279--290. Springer US, Boston, MA, 1996.
- 950 [72] H. Vranken, M. Witteman, and R. Van Wuijtswinkel. Design for testability in hardware software
951 systems. *IEEE Design Test of Computers*, 13(3):79--86, 1996.
- 952 [73] L. Zhao. A new approach for software testability analysis. In *Proceedings of the 28th Interna-*
953 *tional Conference on Software Engineering, ICSE '06*, page 985–988, 2006.
- 954 [74] Y. Zhou, H. Leung, Q. Song, J. Zhao, H. Lu, L. Chen, and B. Xu. An in-depth investigation into
955 the relationships between structural metrics and unit testability in object-oriented systems.
956 *Science china information sciences*, 55(12):2800--2815, 2012.
- 957 [75] G. Zilberfeld. Design for Testability – The True Story , Jan. 2012. [https://www.infoq.com/articles/
958 Testability/](https://www.infoq.com/articles/Testability/).
- 959 [76] Ömer Faruk Arar and K. Ayan. Deriving thresholds of software metrics to predict faults on
960 open source software: Replicated case studies. *Expert Systems with Applications*, 61:106--121,
961 2016.