

Reporting Case Studies in Systematic Literature Studies – An Evidential Problem

Austen Rainer^a, Claes Wohlin^b

^a*Queen’s University Belfast, 18 Malone Road, Computer Science Building, BT9 5BN,
Belfast, Northern Ireland, UK*

^b*Blekinge Institute of Technology, SE-371 79 Karlskrona, Sweden*

Abstract

Context: The term and label, “case study”, is not used consistently by authors of primary studies in software engineering research. It is not clear whether this problem also occurs for systematic literature studies (SLSs). **Objective:** To investigate the extent to which SLSs in/correctly use the term and label, “case study”, when classifying primary studies. **Method:** We systematically collect two sub-samples (2010-2021 & 2022) comprising a total of eleven SLSs and 79 primary studies. We examine the designs of these SLSs, and then analyse whether the SLS authors and the primary-study authors correctly label the respective primary study as a “case study”. **Results:** 76% of the 79 primary studies are misclassified by SLSs (with the two sub-samples having 60% and 81% misclassification, respectively). For 39% of the 79 studies, the SLSs propagate a mislabelling by the original authors, whilst for 37%, the SLSs introduce a new mislabel, thus making the problem worse. SLSs rarely present explicit definitions for “case study” and when they do, the definition is not consistent with established definitions. **Conclusions:** SLSs are both propagating and exacerbating the problem of the mislabelling of primary studies as “case studies”, rather than – as we should expect of SLSs – correcting the labelling of primary studies, and thus improving the body of credible evidence. Propagating and exacerbating mislabelling undermines the credibility of evidence in terms of its quantity, quality and relevance to both practice and research.

Keywords: Systematic Mapping Study, Systematic Review, Systematic Literature Review, Case Study, Credible Evidence

1. Introduction

For almost twenty years, the software engineering (SE) research community has noted that both the definition of the term “case study” and the label “case study” are not used consistently in software engineering research. Of particular concern is the observation that the label “case study” is not used in accordance with the established definitions for case study. This observation was made by Zannier et al. [1], more than 15 years ago, in their analysis of papers published at the International Conference on Software Engineering: “... our sample indicated a large misuse of the term case study.” Not long after, in their guidelines, Runeson and Höst [2] conclude, “... the presented studies range from very ambitious and well organised studies in the field, to small toy examples that claim to be case studies.”

Similar observations have been made, much more recently, by Wohlin [3], and Wohlin and Rainer [4]. These articles conclude that approximately 50% of so-called case studies are not actually case studies, again according to the established definitions. Whilst observations of mislabelling case studies are based on different definitions – Zannier et al. [1] and Runeson and Höst [2] base their observations on the definition by Yin [5], while the more recent studies use the definition by Wohlin [3] – the problem remains the same *across these definitions*. Thus, despite the problem being highlighted over the years, the problem of mislabelling *primary* studies as case studies is still far too common. Furthermore, this problem is not simply a matter of terminological precision or nicety: primary studies that are incorrectly labelled set erroneous expectations for the reader, and are both an *indicator* that the respective primary study may not have been designed or conducted properly, i.e., the primary study is unreliable, and also a *cause* to challenge the findings arising from that study, because they are based on incorrect foundations or at least an incorrect labelling.

To respond to these issues, Wohlin and Rainer [4] developed a checklist and a case study smell indicator to help ensure and assure a primary study as a case study. Rainer and Wohlin [6] report an evaluation of the smell indicator.

But as far as we are aware, there has been no investigation of the extent to which *secondary* studies (or indeed tertiary studies) misclassify primary studies as case studies. As with primary studies, secondary studies that incorrectly label primary studies set erroneous expectations for the reader, and are both an *indicator* that the secondary study may not have been

38 designed or conducted properly – i.e., again, a threat to reliability – and also
39 a *cause* of invalid findings arising from that secondary study – i.e., again,
40 a threat to the validity of the secondary study. There are more significant
41 implications for software engineering research when a secondary study is
42 unreliable or invalid, compared to a primary study. This is because the
43 secondary study is attempting to aggregate or synthesis the state of research
44 in an area. An unreliable or invalid synthesis misrepresents the body of
45 knowledge and, as one implication, can provide incorrect recommendations
46 for practice. We discuss the implications to reliability and validity in more
47 detail later in the article.

48 We use the term, “systematic literature study” (SLS; [7, 8, 9]) as a col-
49 lective term to cover different types of secondary study, such as systematic
50 literature reviews (SLRs; [10]) and systematic mapping studies (SMSs; [11]).
51 It is important to look at both SLRs and SMSs, and other types of system-
52 atic literature study too, because all such reviews and studies should label
53 primary studies correctly and, if they do not, would suffer the problem of
54 unreliability and invalidity stated above. We use the word “label” to refer to
55 a primary study author’s description of their own work or when discussing
56 authors’ descriptions of their studies in more general terms, including both
57 primary study and SLS authors. We use the word “classify” to refer to an au-
58 thor’s assignment of a label to someone else’s work, based on a classification
59 scheme, e.g., an SLS author classifies a primary study.

60 Investigating the extent to which SLSs misclassify primary studies as case
61 studies is important because of the central contribution that SLSs make to
62 evidence-based software engineering and to credible evidence. For example,
63 practitioners are particularly interested in studies conducted in contempo-
64 rary, real-life contexts because, amongst other reasons, the studies’ findings
65 are more likely to be relevant to practitioners’ goals, challenges and decisions.
66 Basili et al. [12] acknowledge the disconnect between practice and research:
67 when research is not conducted in a real-life context, the research output is
68 neither applicable nor scalable. An SLS that misclassifies primary studies
69 as case studies can, however unintentionally, distort the perceived body of
70 credible evidence and misinform practitioners on the applicability and scala-
71 bility of recommendations arising from that evidence. Furthermore, an SLS
72 that misrepresents the primary studies’ research methods undermines our
73 confidence in the SLS overall and in the body of evidence presented by the
74 SLS.

75 Based on the above, we ask the following research question (RQ):

76 RQ For those SLSs that report the empirical research methods of primary
77 studies, do those SLSs correct mislabelled case studies, propagate al-
78 ready mislabelled case studies, or exacerbate the problem by introduc-
79 ing new mislabels?

80 We consider four alternatives to how case studies are handled by SLS
81 authors:

- 82 1. The SLS correctly accepts the labelling by the authors of the primary
83 study, i.e., the primary study is a case study and the SLS correctly
84 classifies it.
- 85 2. The SLS incorrectly accepts the labelling by the authors of the primary
86 study, i.e., the primary study is not actually a case study, but the SLS
87 accepts the respective authors' incorrect labelling of the primary study.
- 88 3. The SLS incorrectly classifies a primary study as a case study even
89 though the respective authors' labelled their study as something other
90 than a case study.
- 91 4. The SLS corrects an incorrectly labelled primary study, now correctly
92 classifying the primary study as a case study. This situation can arise
93 when the respective authors of the primary study had incorrectly la-
94 belled the primary study as something other than a case study. By
95 making this correction, the SLS helps to reduce the problem of misla-
96 belled case studies.

97 One consequence of alternative #2 is the propagation of an incorrect use
98 of the label "case study", whilst alternative #3 exacerbates the situation
99 by introducing further incorrect uses of the label "case study". Conversely,
100 alternatives #1 and #4 result in the correct use of the label "case study",
101 though sometimes alternative #4 may be a fortuitous outcome. We return to
102 these alternatives in Section 6.

103 The remainder of this article is structured as follows: Section 2 presents
104 background and related work; Section 3 explains how we investigated the RQ,
105 including some preparatory analyses to address the RQ; Section 4 presents
106 our analysis of the SLSs; Section 5 summarises our classification of primary
107 studies from across the set of SLSs; Section 6 directly answers the RQ, and
108 considers implications arising from our answer; Section 7 discusses the limi-
109 tations of our research; finally, Section 8 concludes the article.

110 2. Background and related work

111 In this section, we review background and related work. In Section 2.1,
112 we review definitions, guidelines and checklists for case study research, both
113 generally and specific to software engineering. Then, in Section 2.2, we review
114 problems of mislabelling constructs in research, and then the problems of the
115 propagation and exacerbation of such mislabels. In Section 2.3, we suggest
116 five actions to help address these problems, and position our current article
117 relative to those actions. Finally, in Section 2.4, we summarise our review.

118 2.1. Definitions, guidelines and checklists

119 Many disciplines use case studies as a research method to investigate a
120 phenomenon in its real-world context. Also, general guidelines for case study
121 research have been published by several authors, the best known source prob-
122 ably being the book by Yin [13], with its first edition published in 1984.
123 According to Runeson et al. [14], the definition of case study by Yin [13]
124 fits particularly well in software engineering. Admittedly, researchers have
125 different opinions on what constitute a case study, i.e., “one definition does
126 not fit everybody”. However, it is essential to use a definition and provide an
127 appropriate reference. Given the applied nature of software engineering re-
128 search, case study research has become an accepted method in the “research
129 toolbox” for software engineering researchers. Discipline-specific guidelines
130 for case study research in software engineering include the early guidelines
131 by Kitchenham et al. [15] and, more recently, the detailed guidelines from,
132 for example, Runeson and Höst [2], Verner et al. [16] and Runeson et al. [14].
133 Unfortunately, the SE-specific guidelines are primarily focused on the design,
134 conduct and reporting of case studies, and do not emphasise the importance
135 of ensuring the study is, in fact, a case study. Thus, although it may be pos-
136 sible to report a study as a case study, because the study “fits” the guidelines
137 for *reporting* a case study, the study may not actually *be* a case study.

138 Considering the SE-specific guidelines more closely, we see that Runeson
139 and Höst [2] refer to several prior definitions of case study research, including
140 those by Benbasat et al. [17], Robson [18] and Yin [5]. Runeson and Höst [2]
141 stress the real-world context: “Case studies are by definition conducted in
142 real world settings, and thus have a high degree of realism, mostly at the
143 expense of the level of control.” In their book on case study research in
144 software engineering, Runeson et al. [14] provide the following definition of
145 case study for software engineering, which is based on the above sources:

146 “Case study in software engineering is an empirical inquiry that
147 draws on multiple sources of evidence to investigate one instance
148 (or a small number of instances) of a contemporary software engi-
149 neering phenomenon within its real-life context, especially when
150 the boundary between phenomenon and context cannot be clearly
151 specified.”

152 Based on problems identified concerning the mislabelling of case study
153 research, Wohlin [3] proposed a refined definition of *case study*, intended to
154 *clarify* aspects of case study research.

155 “A case study is an empirical investigation of a case, using multi-
156 ple data collection methods, to study a contemporary phenomenon
157 in its real-life context, and with the investigator(s) not taking an
158 active role in the case investigated.”

159 Wohlin’s [3] definition clarifies the need to have more than a single data
160 collection *method* in case study research. This clarification is added to avoid
161 the possible misinterpretation of the concept of *multiple sources of evidence*,
162 since multiple sources may simply imply, for example, the conduct of mul-
163 tiple interviews. Moreover, Wohlin’s definition also differentiates case study
164 research from action research through considering the role of the investigator.

165 In summary, the definition presents five essential components. These
166 components were subsequently used to formulate a checklist [4], which we
167 use in our analyses later in this article. The five components are ordered
168 as follows based on the definition: the empirical investigation of a case;
169 the use of multiple data collection methods; the study of a contemporary
170 phenomenon; the real-life context; and the role of the investigator.

171 Furthermore, for our study, we define “case” as follows:

172 “A case is a series of events or actions in a specific context, typ-
173 ically between different entities (e.g., people, teams, or organisa-
174 tional units) to accomplish an objective.”

175 2.2. *Mislabeleds, propagation and exacerbation*

176 The problems of mislabelling, and of the propagation and exacerbation of
177 mislabels, is not unique to the discipline of software engineering research or to
178 the case study research method. We consider four complementary examples
179 below.

180 In terms of mislabelling in other disciplines, Baker et al. [19] refer to
181 “method slurring” in nursing, where researchers “...blur distinctions be-
182 tween the various qualitative approaches and combine their methodologi-
183 cal prescriptions eclectically.” Baker et al. contrast grounded theory and
184 phenomenology along several dimensions of research, e.g., sources of data,
185 sampling, and validity. They conclude that failure to align the chosen re-
186 search method with the research question, and failure to define the method
187 in the reporting, will severely undermine the credibility of the (in their case,
188 qualitative) research.

189 Haslam [20] points to two forms of “concept creep” in psychology: “hori-
190 zontal” and “vertical” concept creep. Haslam reviews six psychological con-
191 cepts, synthesises from them, and then considers several benefits and draw-
192 backs of concept creep. One concrete instance to demonstrate the benefit
193 of concept creep is the expansion of the concept of *refugee* from the original
194 conception of refugees displaced by conflict, to a refined conception of people
195 also displaced by environmental catastrophe. One set of adverse effects of
196 concept creep is that concepts might become confusing, semantically diluted,
197 redundant, debased, or no longer meaningful. A second set of adverse effects
198 concerns public perception of the discipline; in our context, this concerns
199 software engineering professionals’ perception of, for example, the relevance
200 of research to practice.

201 Mislabelling occurs for other research methods in software engineering
202 too. Ayala et al. [21] observe a similar situation with research that mines
203 software repositories (MSR) and that uses the label “experiment”. Ayala et
204 al. recognise that experimental design decisions impact the respective study’s
205 ability to detect cause-effect relationships, which then has implications for
206 internal validity and reliability. One consequence is that mislabelling a study
207 misrepresents what the study is capable of doing, e.g., suggesting the study
208 is capable of detecting cause-effect relationships when it is not.

209 Finally, Stol et al. [22] investigate the extent to which 98 articles, pub-
210 lished across nine prominent SE journals, used any of the three main vari-
211 ants of Grounded Theory (GT), or techniques from those variants. Later in
212 their article, they consider several implications of “method slurring” of GT:
213 it undermines the legitimacy of GT; it undermines the legitimacy of *other*
214 methods; it misrepresents the state of current research; and it violates a key
215 principle of communicating science, i.e., accurately describing how data was
216 collected and analysed.

217 *2.3. Addressing the problems*

218 We suggest five actions are needed to help address and improve the sit-
219 uation. First, we need appropriate definitions of different research methods.
220 These definitions are, to a large extent, already available to the software
221 engineering community. Second, authors of *primary* studies – and, as we
222 will show, authors of *SLSs* too – should clearly cite the definitions used in
223 their study, e.g., [13, 17, 18, 14, 3], when it comes to case study research.
224 Third, authors of primary studies need to label their studies in accordance
225 with the definitions cited. Fourth, authors of SLSs should ensure that they
226 report research methods correctly, i.e., in accordance with correct definitions
227 of different empirical research methods. Finally, reviewers and editors need
228 to more carefully review manuscripts for these issues. Currently, the review
229 process tacitly endorses low/er standards of quality assessment and misla-
230 belling of research methods.

231 *2.4. Summary*

232 In summary, researchers in software engineering should correctly define
233 and label the research methods they use in their primary studies, and when
234 referring to others’ studies. Furthermore, because of the standards expected
235 of SLSs and of their potential impact on further research and practice, au-
236 thors of SLSs have a particular responsibility to report empirical research
237 methods correctly, including correcting mislabelled primary studies. In this
238 article, we investigate how SLSs handle the labelling of primary studies as
239 case studies.

240 **3. Research approach**

241 This section explains the approach we took to analysing the SLSs and the
242 primary studies selected from those SLSs. We explain how we identified can-
243 didate SLSs (Section 3.1), categorised and selected SLSs (Section 3.2), anal-
244 ysed the SLSs (Section 3.3), selected and analysed primary studies from those
245 SLSs (Section 3.4) and then answered our research question (Section 3.5).

246 *3.1. Identifying candidate systematic literature studies*

247 Given our focus on research methods and, in particular, on case study
248 research, we wanted SLSs with results explicitly related to case studies. Fur-
249 thermore, to mitigate the risk that the search strategy might affect the re-
250 sults, we used two complementary research strategies to collect two sub-
251 samples: Sub-sample I and Sub-sample II. The two sub-samples cover two

252 different time spans as explained below. For both sub-samples, we used Sco-
253 pus for our searches. The choice of Scopus was based on recommendations
254 in the literature, e.g., [23, 24, 25]. We chose to use one database since our
255 objective was to demonstrate a concern, and not to conduct a complete as-
256 sessment of the literature. Moreover, given our objective to demonstrate a
257 concern, rather than comprehensively assess the extent of the concern, we
258 did not apply both search strings to both time spans. The searches for both
259 sub-samples were limited to the “computer science” area and the document
260 types “article” and “review” in Scopus.

261 Sub-sample I was conducted in May 2022 and covers SLSs published in
262 the period 2010-2021. We used 2010 as the starting year since the guidelines
263 of Runeson and Höst [2] and of Verner et al. [16] were published in 2009.
264 After the publication of the guidelines for software engineering, SLS authors
265 ought to be more aware of what constitutes a case study.

266 For Sub-sample I, we looked at those SLSs that explicitly refer to “case
267 study” in their abstract, keywords or title. Our assumption was that explicit
268 reference to “case study” in these ways would give an appropriate sample to
269 investigate how case study research is treated in SLSs.

270 We used the following procedure to identify Sub-sample I:

- 271 1. We conducted three searches using Scopus to capture both systematic
272 literature reviews and mapping studies, as well as studies describing
273 structured reviews of the literature but using alternative phrasings.
- 274 2. We used the following three search strings, all limited to the “computer
275 science” area and the document types “article” and “review”:
276 #1 “systematic literature review” AND software AND (engineering
277 OR development) AND “case stud*”,
278 #2 “systematic review” AND software AND (engineering OR devel-
279 opment) AND “case stud*”,
280 #3 “mapping study” AND software AND (engineering OR develop-
281 ment) AND “case stud*”.
- 282 3. All searches were undertaken by the second author.
- 283 4. The searches returned 169 articles. The search strings resulted in an
284 overlap of articles identified, i.e., the same articles appear in more than
285 one of the searches. Thus, duplicates were removed from searches #2
286 and #3. Also, many articles contained the words used in the search
287 strings, but were not necessarily SLSs. As the next step, therefore, non-
288 SLSs were removed based on reviewing the abstracts. The filtering

289 of articles is summarised in Figure 1. 56 articles remain for further
290 analysis.

291 Sub-sample II was collected in April 2023 and covers SLSs published
292 only in 2022. We chose 2022 as a cutoff to include the full calendar year.
293 For this sub-sample, we adopted a complementary strategy: we looked at
294 all SLSs having variations of formulations related to “literature reviews” in
295 their title.

296 We used the following procedure to identify Sub-sample II:

- 297 1. The search was limited to four general software engineering research
298 journals ranked in the top 10 of Google Scholar when looking for top
299 venues for “Software systems”¹. The four journals investigated were:
300 *IEEE Transactions on Software Engineering*, *Journal of Systems and*
301 *Software*, *Journal of Information and Software Technology*, and *Journal*
302 *of Empirical Software Engineering*.
- 303 2. We used the following search string, limited to the “computer science”
304 area and the document types “article” and “review”: “systematic liter-
305 ature review” OR “mapping study” OR “systematic review” OR “lit-
306 erature review” in the title of the articles.
- 307 3. All searches were undertaken by the second author.
- 308 4. The search returned 40 SLSs published in 2022.

309 In total we have 96 SLSs across the two sub-samples. As noted above, for
310 Sub-sample I, only articles published in software engineering journals (see
311 Table 1) are included, whilst for Sub-sample II, only articles published in
312 four journals were included. A listing of the 96 SLS articles, with links to
313 the articles, is available in Supplement 1 in the online supplementary material
314 linked in Appendix A.

315 3.2. *Categorising and selecting SLS*

316 The articles identified through the search procedure were then categorised
317 according to four criteria: whether the article is in fact an SLS; whether the
318 article reports counts, or percentages, of research methods of the primary
319 studies (so that we can examine the frequency of case studies reported in the

¹https://scholar.google.com/citations?view_op=top_venues&hl=en&vq=eng-software&systems

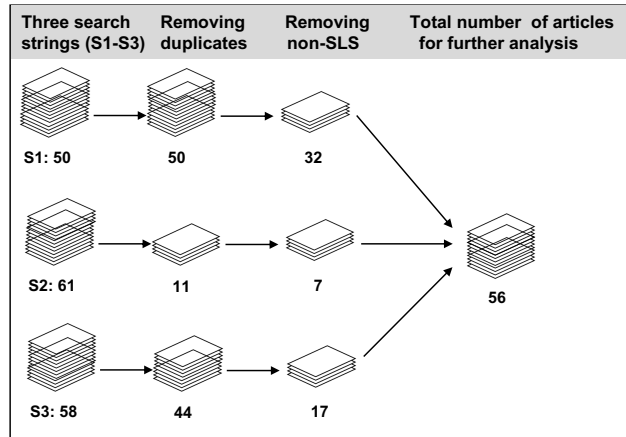


Figure 1: Search and selection process for Sub-sample I.

Table 1: Distributions of identified SLSs across journals.

| Journals for Sub-sample I | Count |
|--|-------|
| Information and Software Technology (IST) | 24 |
| Journal of Systems and Software (JSS) | 11 |
| Journal of Software Evolution and Process (JSEP) | 5 |
| Transactions on Software Engineering (TSE) | 5 |
| Software Engineering and Knowledge Engineering (SEKE) | 3 |
| Empirical Software Engineering (EMSE) | 2 |
| IET Software (IETS) | 2 |
| Requirements Engineering Journal (REJ) | 2 |
| Software Quality Journal (SQJ) | 1 |
| International Journal of Software Engineering Applications (SEA) | 1 |
| Sum | 56 |
| Journals for Sub-sample II | Count |
| Journal of Systems and Software (JSS) | 17 |
| Information and Software Technology (IST) | 14 |
| Empirical Software Engineering (EMSE) | 5 |
| Transactions on Software Engineering (TSE) | 4 |
| Sum | 40 |

320 literature); whether the article supports traceability from the SLS back to
 321 the source primary study (so that we can independently check the label of
 322 “case study”); and whether the article makes any explicit comment on the
 323 problem of labelling studies as “case studies”. The definition of categories is
 324 summarised in Table 2.

325 The two authors independently categorised all articles for Sub-sample I,
 326 shared their independent categories and discussed them, arriving at a consen-
 327 sensus decision for the category of each article. The categorisation of articles,
 328 and the process of filtering them, is summarised in Table 3. In the filtering

329 process, one article was identified in the search that did not meet the criteria
 330 in the search string. In other words, one article was incorrectly retained from
 331 the original searches, which was subsequently removed when cross-checking
 332 the identified articles with the search criteria.

333 Having categorised the articles, we removed articles in Categories A and
 334 B as they provide no meaningful information for our subsequent analysis. For
 335 Categories C, D and E, the lack of traceability means we cannot independ-
 336 dently check how the SLSs classified each primary study, however articles in
 337 these categories will be useful for some of our analyses, concerning frequency
 338 counts and comments about the misuse of labels. Articles in Categories F
 339 and G will be the most valuable articles as they provide traceability to the
 340 source primary studies.

341 We calculated unweighted and weighted Kappa statistics for Sub-sample
 342 I to evaluate the agreement between the two authors' categorisations. The
 343 unweighted Kappa statistic assumes the categories are entirely independent
 344 whilst the weighted Kappa statistic takes account of the ordering of the
 345 categories. The Kappa statistics are reported in Table 3. Overall, we have
 346 substantial to near-perfect agreement for Sub-sample I. Given the high-level
 347 of agreement, the categorisation for Sub-sample II was done only by the
 348 second author, except for two SLSs where a second opinion was perceived as
 349 needed. The categorisation of the SLS, is available in Supplement 1 in the
 350 online supplementary material linked in Appendix A.

Table 2: Decision criteria for the categories of SLSs.

| Cat. | Excluded | SLS | Counts | Traceability | Labelling |
|------|----------|-----|--------|--------------|-----------|
| A | Yes | No | N/A | N/A | N/A |
| B | Yes | Yes | No | N/A | No |
| C | | Yes | No | N/A | Yes |
| D | | Yes | Yes | No | No |
| E | | Yes | No | No | Yes |
| F | | Yes | Yes | Yes | No |
| G | | Yes | Yes | Yes | Yes |

Notes:

Category A articles are not actually SLSs.

Category B articles are SLSs but lack necessary information.

SLS = Systematic literature study

Counts = SLS reports frequencies or percentages of research methods of the primary studies.

Traceability = SLS reports classification of research method for each primary study.

Labelling = SLS makes a comment on the misuse of the label "case study".

351 During our categorisation of Sub-sample I, we identified three articles

352 that are actually tertiary studies, rather than SLSs. Two of these articles
353 are Category B articles and are removed almost immediately during the
354 funneling process (see Table 3). The third article [26] is a Category G article.
355 This article is separated out in Step 5 of Table 3. A second article [27] in
356 Category G is also removed. These two articles are removed because they
357 both comment on the misuse of the label “case study”, but do not take this
358 problem into account in their subsequent analysis. For a study in Category G
359 to be included in our analysis, the article must comment on the mislabelling
360 of case studies and also adapt their analysis accordingly. Thus, the two
361 studies in Category G are removed from the analysis. We further discuss the
362 exclusion of these two articles in Supplement 2 in the online supplementary
363 material linked in Appendix A. The effect of these exclusions is that we did
364 not identify a Category G article we could include in our analysis in the two
365 sub-samples. Note that because the Category F articles do not comment on
366 the misuse of labels, we cannot determine whether the authors know there
367 is a problem, but do not report it, or do not know there is a problem. We
368 assume the authors of SLSs in Category F do not know there is a problem.

Table 3: Identification and filtering of articles in different categories.

| Step | Description | Categories | Articles remaining for Sub-sample I | Articles remaining for Sub-sample II |
|---|--|------------|-------------------------------------|--------------------------------------|
| 1 | Input from search and screening | N/A | 56 | 40 |
| 2 | Remove Categories A and B | A & B | 24 | 29 |
| 3 | Remove Category D | D | 16 | 7 |
| 4 | Remove Categories C and E | C & E | 11 | 1 |
| 5 | Separate out remaining tertiary studies | N/A | 10 | 0 |
| 6 | Remove articles in Category G not needing a detailed analysis | N/A | 9 | 0 |
| 7 | After cross-check with search string, remove articles incorrectly included | N/A | 1 | 0 |
| | SLSs remaining for detailed analyses (all Category F) | N/A | 8 | 3 |
| Kappa statistics for Sub-sample I: Unweighted Kappa statistic: 0.796 (substantial agreement) Weighted Kappa statistic: 0.869 (near-perfect agreement) | | | | |

369 Table 4 presents summary information on the selected SLSs from the two
370 sub-samples. Given the removal of two articles in Category G, all eleven
371 SLSs selected for a detailed analysis are in Category F.

Table 4: Summary of SLSs selected for analysis.

| ID | Cit. | Jnl | Title | Type |
|----------------------|-------------|------------|--|-------------|
| <i>Sub-sample I</i> | | | | |
| SLS P2 | [28] | IST | Software engineering research for computer games: a systematic review | SLR |
| SLS P11 | [29] | IST | Past and future of software architectural decisions – a systematic mapping study | SMS |
| SLS P16 | [30] | JSS | Feature extraction approaches from natural language requirements for reuse in software product lines: a systematic literature review | SLR |
| SLS P21 | [31] | IST | The effects of test driven development on internal quality, external quality and productivity: a systematic review | SLR |
| SLS P33 | [32] | IST | What software reuse benefits have been transferred to the industry? A systematic mapping study | SMS |
| SLS P36 | [33] | IST | Empirical evidence in follow the Sun software development: a systematic mapping study | SMS |
| SLS P44 | [34] | IETS | Impact of design patterns on software quality: a systematic literature review | SLR |
| SLS P45 | [35] | IST | Empirical software product line engineering: a systematic literature review | SLR |
| <i>Sub-sample II</i> | | | | |
| SLS P59 | [36] | JSS | Revealing the state of the art of large-scale agile development research: A systematic mapping study | SMS |
| SLS P81 | [37] | IST | Ambiguity in user stories: A systematic literature review | SLR |
| SLS P85 | [38] | IST | Software security patch management - A systematic literature review of challenges, approaches, tools and practices | SLR |

372 3.3. Analysis of the SLS articles

373 Before analysing the primary studies, we wanted to assess the SLSs. We
374 are interested in two aspects of the SLSs:

- 375 1. The designs of the (retained) SLSs, e.g., what did the SLS consider
376 when searching for, selecting and analysing the primary studies? To
377 perform this analysis we used the eleven SLSs in Category F listed in
378 Table 4. This analysis is presented in Section 4.1.
- 379 2. Comments made by the authors of the SLSs about the mislabelling
380 of “case studies” in software engineering. To perform this analysis
381 we used SLSs categorised as Category C and E according to Table 2.
382 There are no SLSs available to consider from Category G, because we
383 removed these from our analyses. Our analysis of the comments made
384 by Category C and E articles is presented in Section 4.2.

385 3.4. Analysis of samples of primary studies from several SLSs

386 Having selected eleven SLSs and assessed them, we then turned to in-
387 vestigating a sample of primary studies from each SLS. We wanted to check

388 whether the SLS authors were classifying primary studies correctly. We in-
389 vestigated three aspects:

- 390 1. We compared the sample from each SLS’s classification of primary stud-
391 ies against our classification of those studies. For our classification, we
392 used the case study checklist briefly discussed in Section 2.1 and used in
393 our previous analysis of primary studies [4]. As described in the check-
394 list, the five components are investigated one at the time until one
395 component does not fulfil the definition for being a case study. Each
396 author of the current article first independently reviewed the primary
397 studies. We then discussed our views and resolved disagreements. For
398 those primary studies that we agreed were borderline “case studies”,
399 we took a generous attitude, i.e., to classify a study as a “case study”.
400 This analysis is presented in Section 5.1.
- 401 2. We examine those (few) studies correctly assessed as being case studies
402 and present information about these studies in Section 5.2.
- 403 3. We also checked how the authors of the primary studies labelled their
404 own study. This allows us to undertake a three-way comparison of the
405 SLS authors’ labels with the primary-study authors’ labels with our
406 labels. This analysis is most directly relevant to our research question
407 and is presented in Section 6.

408 The analysis of primary studies reported later in this paper, e.g., in Sec-
409 tion 5, is based on sampling six batches of primary studies from across our
410 two sub-samples, with three batches per sample. For each batch, we sought
411 to sample three primary studies per SLS, however as not all SLSs had nine
412 primary studies classified as case studies we sometimes had to settle for less
413 than nine case studies per SLS per batch. In total, 79 primary studies are
414 analysed. For simplicity, we refer to the two samples rather than the six
415 batches. A fuller explanation of the batches is provided in Supplement 3 in
416 the online supplementary material linked in Appendix A.

417 *3.5. Investigating the research question*

418 As noted earlier, we collected three sets of labels: the SLS authors’ la-
419 belling of primary studies, the primary-study authors’ own labels, and our
420 labels. We did this for primary studies from the two sub-samples. The three
421 sets of labels allow us to assess the four alternative answers to our research
422 question. We present the results of this analysis in Section 6, including a
423 consideration of the implications of our results.

424 4. Analysis of the selected SLSs

425 In this section we first consider the designs of the eleven SLSs, as these
426 designs provide insights into how the respective SLSs analysed the primary
427 studies (Section 4.1). We then briefly discuss comments made by some of the
428 SLSs on the way that the label “case study” has been misused in primary
429 studies (Section 4.2).

430 4.1. The designs of the selected systematic literature studies

431 Our analysis of the designs of the eleven SLSs we studied is presented
432 in Table 5. Additional information, summarising sources cited by the eleven
433 SLSs, is presented in Table 6.

434 Table 5 analyses each of the eleven SLSs according to thirteen criteria.
435 The criteria are numbered down the left-hand side of the table and sum-
436 marised in the notes at the end of the table. The “answer” to each criterion
437 may take several values, which are also summarised in the table’s notes.

438 For each sub-sample, the right-most column presents a proportion. This
439 proportion gives an *indication* of the extent to which the SLSs in the respec-
440 tive sub-sample considered the criterion. Then, in the last row of the main
441 table, labelled “**Sum**”, a total is given for most of the criteria in each col-
442 umn. This total gives an *indication* of the overall quality of the study design
443 for that SLS. We emphasise these totals and proportions are approximations
444 based our two sub-samples.

445 The proportions in Table 5 *suggest* that just over a half of the SLSs in
446 Categories F (see item #2; 6/11 SLSs) had an explicit research question
447 about research methods; that most SLSs explicitly analysed the research
448 methods (item #3; 10/11); and that most SLSs commented on the most
449 common research method (item #4; 8/11). Only *three* SLSs had an explicit
450 definition for “case study” (item #5), though none of the three were entirely
451 consistent with established definitions; just over a half of the SLSs (item #7;
452 6/11) used an existing classification (see Table 6 for more information) and
453 for those that did not use an existing classification, no SLS explained how
454 they developed their classification (item #8). Also, the SLSs varied in the
455 number of research methods they considered, from two methods to ten (item
456 #9).

457 In a previous study [4], we found that about 40% of studies reported as
458 case studies by their authors were in fact better understood as small-scale
459 evaluations. We base our definition of *small-scale evaluation* on Robson [39]:

460 Small-scale evaluations aim to demonstrate, illustrate, or show
461 the feasibility of a proposed solution, for example, a concept or
462 tool. The evaluation is conducted with a single researcher or a
463 small team of researchers. It is run over a short period, most
464 often in a single site, for example, a laboratory.

465 Of the eleven SLSs, four SLSs have a method that appears similar to the
466 small-scale evaluation (item #10). About a quarter of the SLSs explicitly
467 cite guidance on the design of case studies in software engineering (item #11;
468 see Table 6 for more information) and only three of the SLSs cite guidance
469 on the design of empirical studies in software engineering (item #12' see
470 Table 6 for more information). For items #11 and #12 in Table 5, there is no
471 straightforward way to analyse the SLSs because the relevance of the citation
472 can depend on context. For example, in Table 6, the sources cited for existing
473 classifications include citations to generic guidance on software engineering
474 research. For items #11 and #12, we therefore looked for citations made in
475 the context of the SLS paper discussing the design or conduct of software
476 engineering research. Finally, item #13 suggests that over a half (7/11) of
477 the eleven SLSs explicitly included a Quality Assessment as part of the SLS.

478 Separate from Table 5, SLS P45 [35] is of particular interest since it pro-
479 vides a quality assessment of the reporting of the empirical methods used,
480 including case study. The authors assess the quality of the case study report-
481 ing using ten criteria that case studies should fulfil. This in-depth analysis
482 is valuable, although it focuses on what the authors report and not necessar-
483 ily what the authors of the primary studies did. Furthermore, the authors
484 also report the case studies in relation to characteristics that case studies
485 ought to report. Their analysis is done based on the characteristics provided
486 by Robson [18]. The assessments illustrate that the authors are concerned
487 with how case studies are reported, although they do not directly discuss
488 the case studies in relation to the definitions of case study research, and do
489 not therefore explicitly consider whether a study has been *mislabeled* as a
490 case study. The issue of the lack of use of (established) definitions may be
491 exemplified with the authors of SLS P45 [35] reporting both *industrial* and
492 *academic* case studies.

493 4.2. Comments on the case studies

494 Separate to our analysis of the 11 SLSs we selected for our main analysis
495 (see Table 3 and Table 4), we wanted to see whether *other* articles from our

Table 5: Analysis of the SLS designs for the eleven Category F SLSs.

| # | Sub-sample I SLS # | | | | | | | | | Sub-sample II SLS # | | | |
|------------|--------------------|------|-------|------|------|------|------|------|-----|---------------------|------|------|-----|
| | 2 | 11 | 16 | 21 | 33 | 36 | 44 | 45 | Prp | 59 | 81 | 85 | Prp |
| 1 | 2010 | 2014 | 2015 | 2016 | 2018 | 2018 | 2020 | 2020 | NA | 2022 | 2022 | 2022 | NA |
| 2 | Y | N | Y | N | Y | Y | N | Y | 5/8 | N | N | Y | 1/3 |
| 3 | Y | Y | N | Y | Y | Y | Y | Y | 7/8 | Y | Y | Y | 3/3 |
| 4 | N | Y | Y | Y | Y | Y | N | Y | 6/8 | Y | N | Y | 2/3 |
| 5 | Y | N | N | N | N | N | N | N | 1/8 | Y | N | Y | 2/3 |
| 6 | N | NA | NA | NA | NA | NA | NA | NA | 0/8 | P | NA | P | 2/3 |
| 7 | Y | Y | N | N | Y | Y | N | N | 4/8 | Y | N | Y | 2/3 |
| 8 | NA | NA | NS | NS | NA | NS | NS | NS | 0/8 | NA | NS | NA | 0/3 |
| 9 | 3 | 4 | 2 x 2 | 5 | 4 | 9 | 3 | 3 | NA | 8 | 6 | 10 | NA |
| 10 | N | Y | N | N | Y | P | N | N | 3/8 | N | N | Y | 1/3 |
| 11 | N | Y | N | N | N | N | Y | Y | 3/8 | N | N | N | 0/3 |
| 12 | Y | N | N | N | N | N | N | Y | 2/8 | Y | N | N | 1/3 |
| 13 | Y | N | Y | Y | N | N | Y | Y | 5/8 | N | Y | Y | 2/3 |
| Sum | 6 | 5 | 3 | 3 | 5 | 5 | 3 | 6 | | 6 | 2 | 8 | |

Notes:

Y = Yes; N = No; P = Potentially present; Prp = Proportion of Ys & Ps; NA = Not applicable
 NS = Not clearly stated in the paper

Summary of criteria

1. Year published: In what year was the SLS published?
2. Specific RQ: Is there a research question (RQ) in the article that specifically asks about research methods?
3. Research methods analysed: Are research methods explicitly analysed?
4. Common method: Do the authors of the SLS explicitly discuss the most common, or predominant, research method?
5. Case study definition: Is “case study” explicitly defined in the paper?
6. Consistent definition: Is the case study definition consistent with established definition, e.g., [14]?
7. Existing classification: Do the authors use an existing, published classification of studies?
 The published classifications are summarised in Table 6.
8. Classification dev.: If the answer to Q#7 is “No”, or NA, how was it developed?
9. Research methods count: Number of research methods used in the SLSs’ classification?
10. Small-scale evaluation: Is there any type of suitable method related to small-scale evaluation?
11. SE case study design: Is a standard citation to case study design in SE used?
 Citations are summarised in Table 6.
12. Empirical studies: Is a citation to empirical studies used for research methods?
 Citations are summarised in Table 6.
13. Quality assessment: Does the paper report the quality assessment of primary studies?

Sum: Sum of Ys and Ps

496 broader dataset had recognised the problem of mislabelling. We identified
 497 six SLSs, all from Categories C and E and published between 2011 and 2022,
 498 that explicitly comment on the (mis)use of the label, “case study”.

499 These six SLSs all demonstrate that, over the years, the problem of mis-
 500 labelling primary studies as case studies, when they are not actually case
 501 studies, is well recognised and far too common. These six SLSs are in addi-
 502 tion to other sources cited in Section 1. Yet although these SLSs recognise
 503 problems with the use of the label “case study” these SLSs, *themselves*, misla-
 504 bel primary studies. One of these six SLSs (SLS P18) does however explicitly

Table 6: Sources cited by the SLSs.

| ID | Sources |
|--|---------------------|
| Sources cited for Q7: Existing classifications | |
| SLS P2 | [40] [15] |
| SLS P11 | [41] [42] [43] [44] |
| SLS P33 | [45] |
| SLS P36 | [46] |
| SLS P59 | [47] [48] [49] |
| SLS P85 | [50] |
| Sources cited for Q11: Case study design | |
| SLS P11 | [51] [52] [53] |
| SLS P44 | [14] [53] |
| SLS P45 | [14] [2] |
| Sources cited for Q12: Design of empirical studies | |
| SLS P2 | [54] [15] [40] |
| SLS P45 | [18] [40] |

505 introduce a distinct type, “Example”, to help distinguish true case studies
506 from other types of study. Comments on the mislabelling of case studies
507 are further discussed in Supplement 4 in the online supplementary material
508 linked in Appendix A.

509 5. Summary of analysis of primary studies

510 Having considered the SLSs in the previous section, we now turn to con-
511 sider the primary studies. For conciseness we present a summary of our
512 analysis in this section. A listing of the primary studies can be found in
513 Supplement 5 and the full analysis of the primary studies is provided in Sup-
514 plement 6. Links to the online supplementary material are provided in Ap-
515 pendix A. We first consider the classification of primary studies (Section 5.1)
516 and then consider the set of studies we identified as correctly classified as
517 case studies (Section 5.2).

518 5.1. Classifications of the primary studies

519 Based on our analysis of the 79 primary studies, and connecting our
520 analysis back to the criteria identified in Wohlin’s *refined* definition of a case
521 study in Section 2.1, we conclude that:

- 522 1. The most common reason for deciding that a study was not a case
523 study was that the study was not conducted in a real-life context. As
524 Wohlin [3] has suggested before: it is a case and it is a study, but it
525 is not a case study. Studies not conducted in a real-life context were

526 often small-scale evaluations in a laboratory environment, as identified
527 in [3].

528 2. The second most common reason was that the study was not contem-
529 porary. Studies that were not of a contemporary phenomenon most
530 commonly are archival studies, for example, of open source projects.

531 3. For some studies, they were not case studies, either because a *case* was
532 not being studied (as required of a case study), or the study was better
533 described as action research.

534 Across our sample of 79 primary studies, we found zero instances were
535 the SLS authors classified a primary study correctly as a “case study” when
536 the original authors mislabelled the study as not a case study. On the other
537 hand, we found 29 instances (37% of the analysed primary studies) where
538 the SLS classified a primary study as a “case study” despite the fact that
539 the authors of the primary study did not label their study that way.

540 *5.2. Investigation of the nineteen case studies*

541 The nineteen actual case studies are published between 2003 and 2020.
542 This means that the case study guidelines published in 2009 by Runeson
543 and Höst [2] and Verner et al. [16] were not available for all primary-study
544 authors, but these guidelines were available to the SLSs we studied.

545 About half of the nineteen case study articles refer to general references to
546 empirical research and metrics, and slightly fewer refer specifically to sources
547 concerning case study design.

548 Based on our comparison of primary studies between the two sub-samples,
549 we conclude that too few primary studies refer to sources concerning the
550 research method being used, even when they are correctly labelled.

551 The lack of citations to sources concerning research methodology may be
552 a reason that too many studies are mislabelled, although other reasons may
553 exist, for example, the different expectations at different publication venues.

554 Finally, it is worth noting that we, the authors of the current article, did
555 not initially agree on the classification of “case study” for six of these nineteen
556 case studies. This may appear surprising. We argue that the main reason for
557 this disagreement is that the articles are often unclear about the components
558 of the case study definitions (see Section 2.1) – in other words, authors do not
559 clearly and transparently report the design etc of their study [55] – which is
560 natural since many articles do not refer to sources for research methodology,
561 in particular not to case study methodology.

562 **6. The treatment of case studies by SLSs**

563 Having completed several analyses, we are now in a position to directly
 564 consider our research question. In this section, we first discuss the frequency
 565 of case studies in software engineering. We present reported frequencies
 566 from the SLSs and compare those with our revised frequencies. We then
 567 use those frequencies to examine the degree to which SLSs have distorted,
 568 through propagation and exacerbation, the status of case studies in software
 569 engineering research.

570 *6.1. The prevalence of case studies in SE*

571 Table 7 presents counts and percentages of the prevalence of case studies
 572 in software engineering and of our estimated correctness of this prevalence,
 573 for the two sub-samples. As the table indicates, the full sample of eleven
 574 SLSs report a frequency of case studies in the range 10% (SLS P85) to 60%
 575 (SLS P59) with a mean average of 39%. For those studies classified by
 576 the SLSs as case studies, our estimates of correctness show that the SLS
 577 authors misreport the number of case studies in the literature. The degree
 578 of misreporting ranges from 33% (SLS P59) to 100% (SLSs P2, P16, P85).
 579 Based on our re-classification of primary studies, the average over-estimate
 580 for this set of SLSs is 76%, given that we estimate that 24% of primary
 581 studies are correctly classified.

Table 7: Prevalence of case studies and an estimated correctness rate.

| ID | Jnl | Domain | # RM | Counts and % of case studies | Our estimated correctness |
|--|------|-------------------------------------|---------|---------------------------------|------------------------------|
| <i>Sub-sample I</i> | | | | | |
| SLS P2 | IST | Games development | 3 | 4/20 (20%) | 0/4 (0%) |
| SLS P11 | IST | Software architecture | 4 | 31/144 (22%) | 1/9 (11%) |
| SLS P16 | JSS | Software product lines | 2 | 2/13 (15%) | 0/2 (0%) |
| SLS P21 | IST | Test-driven development | 4 | 8/27 (30%) | 2/8 (25%) |
| SLS P33 | IST | Software reuse | 4 | 25/50 (50%) | 2/9 (22%) |
| SLS P36 | IST | Follow-the-sun software development | 9 | 18/32 (56%) | 1/9 (11%) |
| SLS P44 | IETS | Design patterns | 4 | 27/50 (54%) | 1/9 (11%) |
| SLS P45 | IST | Software product lines | 3 | 35/62 (56%) | 4/9 (44%) |
| Sub-total | | | | 150/398 (38%) | 11/59 (19%) |
| <i>Sub-sample II</i> | | | | | |
| SLS P59 | JSS | Large scale agile development | 8 | 78/129 (60%) | 6/9 (67%) |
| SLS P81 | IST | User stories | 6 | 4/17 (24%) | 2/4 (50%) |
| SLS P85 | IST | Security patch management | 10 | 7/68 (10%) | 0/7 (0%) |
| Sub-total | | | | 89/214 (42%) | 8/20 (40%) |
| Total | | | | 239/612 (39%) | 19/79 (24%) |
| Notes: RM = Research methods in the classification for each SLS. | | | | | |

582 Our overall estimate of correctly labelled case studies is 24%. This esti-
583 mate is based on the use of established definitions of case study research, and
584 is derived only from the re-analysis of a sample of primary studies reported as
585 case studies in the eleven SLSs. This estimate is *lower* than both the SLSs’
586 overall estimate in relation to other research methods and the estimates we
587 reported previously [3, 4], which were around 50% based on assessing pri-
588 mary studies claiming to present case studies. Our estimated corrections for
589 the SLSs suggests that the authors of SLSs contribute to the over-reporting
590 of the use of case study research in software engineering.

591 In 50 primary studies of the 79 primary studies in the overall sample, the
592 *primary-study* authors claim that they present case studies. Based on our
593 assessment, only 19 of these 50 primary studies present case studies. Thus,
594 the correctness rate for *primary-study authors* self-reporting is 38% (19/50),
595 which again is lower than our previous estimates in [3, 4].

596 In our previous analyses, the focus was either on case study articles pub-
597 lished in journals (i.e., [3]) or on journal articles citing the case study guide-
598 lines by Runeson and Höst [2] (i.e., [4]). But in the current article, the
599 primary studies included are those identified *by the SLS authors* when tar-
600 geting a specific area of research in software engineering. Thus, we have three
601 samples constructed using different sampling strategies.

602 The misclassifications by the SLSs is partly dependent on the classifica-
603 tion scheme used by the SLS authors and hence the frequencies of different
604 research methods. One SLS only classifies using two research methods (SLS
605 P16), whereas another SLS (SLS P85) classifies using ten research meth-
606 ods, and the other SLSs are in between these two limits. The differences in
607 classification schemes are an issue since, depending on the scheme used, all
608 primary studies are “forced” into the classes in the scheme.

609 6.2. The mislabelling of case studies in SE

610 We now return to our RQ and specifically to the four alternative ways
611 of labelling a study’s research method, as discussed in Section 1. The four
612 alternatives are concerned with how an SLS treats primary studies as case
613 studies.

614 Table 8 summarises the four alternative ways of labelling case studies for
615 the two sub-samples of 59 and 20 primary studies respectively.

616 Our overall conclusions from the analysis presented in Table 8 are av-
617 eraged over the two sub-samples. Overall, the table indicates that authors
618 of SLSs have a tendency to simply restate the label, “case study”, used by

Table 8: SLS authors handling of case studies.

| Alter- native | Description | Sub-sample counts and %s | | |
|------------------|---|--------------------------|-------------------|----------|
| | | I Batches 1-3 | II Batches 4-6 | Total |
| 1 | Restate a correctly labelled primary study | 11 (19%) | 8 (40%) | 19 (24%) |
| 2 | Restate an incorrectly labelled primary study | 27 ^a (46%) | 4 (20%) | 31 (39%) |
| 3 | Distort a correctly labelled primary study | 21 ^a (36%) | 8 (40%) | 29 (37%) |
| 4 | Correct an incorrectly labelled primary study | 0 (0%) | 0 (0%) | 0 (0%) |
| Total | | 59 | 20 | 79 |

^adue to rounding, 46% + 36% = 82% but (27+21)/59=81%

619 the authors of the respective primary study, independent of whether the la-
620 bel of “case study” is correct or not ((19+31)/79 instances; 63%). In about
621 two-fifths of instances (31/79; 39%) the authors of the SLSs propagate an *in-*
622 *correct* use of the term “case study”. Furthermore, in over a third of instances
623 (29/79; 37%), SLS authors exacerbate the situation: they incorrectly classify
624 a primary study as a “case study”, even when the primary study was *not*
625 *labelled* as “case study” by the authors of the primary study. Taken together,
626 propagation and exacerbation occur in 76% of the instances ((31+29)/79).
627 By doing this, the SLS authors are distorting the body of evidence in the area
628 of the respective SLS; or, in other words, by reporting *incorrect information*
629 the SLS authors do not report *credible evidence*, one of the main objectives
630 of SLSs.

631 Sub-sample II, drawn from 2022, does seem to present “better” results
632 than Sub-sample I. Sub-sample II has a considerably lower percentage of re-
633 stating incorrectly labelled primary studies (20% compared for 48% for Sub-
634 sample I), a higher percentage of re-stating correctly labelled primary studies
635 (40% compared to 19% for Sub-sample I) and a slightly lower percentage of
636 simply restating the primary study’s label (60% compared to 67%). Sub-
637 sample II has a slightly higher percentage of distorting a correctly labelled
638 primary study (40% to 36%). Like Sub-sample I, Sub-sample II does not
639 correct any incorrectly labelled primary studies. However, Sub-sample II
640 only includes three SLS, and one of them (SLS P59) is the SLS with the
641 highest correctness, see Table 7. Thus, it is premature to assume that the
642 situation is improving. Further studies are needed to be able to investigate
643 any potential improvement over time.

644 An underlying, recurring problem seems to be that many researchers have
645 a far too flexible interpretation of case study research or, alternatively, they
646 lack sufficient knowledge concerning the definition of case study research. Re-
647 searchers need to know and understand the definitions of research methods,

648 and use them correctly. Items #5 and #6 in Table 5 show that only three
649 of the eleven SLSs explicitly defined the term, “case study”, and *none* of the
650 three SLSs presented a definition consistent with the established definitions,
651 although two had partially consistent definitions. SLS authors need to be
652 even more aware of these issues, compared to authors of primary studies, as
653 SLSs attempt to gather and synthesise the available evidence in a research
654 area. SLS authors therefore need to be careful gatekeepers, who ensure that
655 research methods are correctly reported in SLSs. If not, we do not have *cred-*
656 *ible evidence*; we get the propagation of incorrect evidence or, even worse,
657 the exacerbation of incorrect evidence.

658 A related issue, contributing to the problem, is *how* SLS authors choose to
659 classify the research methods used in the primary studies. If authors develop
660 their own classification scheme and decide *a priori* the research methods in
661 their classification, then there is a risk that not all research methods reported
662 across the primary studies will be properly represented: there may situations
663 where a primary study is misclassified to fit the *a priori* classification. Item
664 #7 in Table 5 shows that six of the eleven SLSs used an existing classification,
665 however even then it is not clear that the existing classification would cover
666 all types of primary study, as indicated by the range of research methods for
667 item #9 of Table 5. For the other five SLSs, none explain how they developed
668 their classification.

669 A potential way forward is to either use an already accepted and compre-
670 hensive classification scheme (preferred option) or formulate the classification
671 based on the research methods stated by the authors of the primary study,
672 whilst also ensuring that the methods stated in the primary studies adhere
673 to the formal definitions of each research method. In other words, to take
674 a bottom-up approach to classification, with careful checks on the labels
675 allowed into the classification.

676 Relating these results back to our research question, presented in Sec-
677 tion 1, we conclude that our analysis of primary studies corroborates prior
678 research (the first item listed below), and, for SLSs, produces three novel
679 findings (the subsequent three items listed below), i.e.:

- 680 1. Too many primary studies incorrectly present themselves as being case
681 studies. This conclusion is supported by the literature, including, for
682 example, [1, 2, 3, 4].
- 683 2. SLS authors frequently simply restate the incorrect case study label
684 from the respective primary study.

- 685 3. In many cases, SLS authors incorrectly *change* the labels of primary
686 studies to “case study”. Thus, they make the incorrect labelling of case
687 studies worse than if only re-stating the original label for the primary
688 study.
- 689 4. Given our results, it seems very unlikely that SLS authors correct the
690 labelling, i.e., correctly classify a primary study as a “case study” when
691 the primary study is presented incorrectly as something else by the
692 authors of the primary study.

693 In summary, based on our two sub-samples, SLS authors distort the la-
694 labelling of primary studies as “case studies”, through propagation and exac-
695 erbation, in 76% $((27+21)/59 = 81\%$ for Sub-sample I and $(4+8)/20 = 60\%$
696 for Sub-sample II) of the instances of such labelling.

697 6.3. *Implications*

698 The mislabelling of primary studies as case studies raises several implica-
699 tions. We consider four broad implications here, supporting our discussion
700 with illustrative examples from the eleven SLSs we analysed. Given publica-
701 tion constraints, we focus here on illustrative examples that can be concisely
702 presented, which means the examples are structured or quantified. The prin-
703 ciples behind our illustrative examples are intended to apply much more
704 broadly across the different kinds of case study research. In their tertiary
705 review of research synthesis in software engineering, Cruzes and Dybå [56]
706 categorised the methods of synthesis into 15 categories, e.g., narrative, the-
707 matic, grounded, comparative, content, vote-counting, and quantitative (the
708 15th was no method). We cite Cruzes and Dybå [56] to acknowledge that
709 the broad implications we discuss here will manifest differently across differ-
710 ent kinds of case study research and different methods of research synthesis.
711 Also, our discussion here complements the consequences and implications
712 already recognised in Section 2.2.

713 6.3.1. *Labels set expectations*

714 Consider an example where an SLS incorrectly labels a small-scale evalu-
715 ation as a case study. Typically, small-scale evaluations occur in a laboratory
716 environment. To label such a study as a case study creates the erroneous
717 expectation for the reader that the study took place in a real, contemporary
718 setting and was conducted to the standards of case study research. It also
719 creates erroneous expectations about the way that the study’s findings might

720 be generalised, e.g., through an analytical, or theoretical, mechanism, rather
721 than through a statistical mechanism.

722 More generally, it is important to clearly and transparently report the
723 actual research method used in the primary study so that the reader has a
724 clear and transparent knowledge of the research design and the underlying
725 research process [55]. Or, in other words, when we label a primary study as
726 a case study, we create expectations for the reader about the aims, design,
727 conduct, interpretation of results, and implications of results for professional
728 practice.

729 The incorrect labelling of a primary study by an SLS is therefore both
730 an *indicator* that the respective SLS may not have been designed, conducted
731 etc. properly, and also a *cause* of invalid findings.

732 In terms of the indicator, the analysis and results of an SLS are dependent
733 on the quality of the primary studies that are input to that SLS. If these
734 inputs are mislabelled, this raises uncertainty about the general treatment
735 of the primary studies as input to the SLS and, more generally, about the
736 overall reliability of the SLS.

737 In terms of cause, if primary studies are mislabelled, this directly affects
738 how each primary study is subsequently treated by the SLS and, by extension,
739 how primary studies are synthesised. This raises concerns about the validity
740 of the aggregated or synthesised findings.

741 6.3.2. *Estimations and comparisons*

742 The mislabelling of primary studies leads to difficulties estimating the
743 number of genuine case studies conducted in SE in general or conducted for
744 a particular domain in SE. This leads to difficulties comparing results both
745 *within* and also *across* secondary studies and difficulties reporting trends over
746 time, e.g., a supposed growth in the number of case studies being conducted.
747 As one example, SLS P11 states that, "... surveys have similar distributions
748 to case studies, but the number of surveys is about half the number of case
749 studies..." Our findings challenge both SLS P11's observation on similar
750 distributions as well as its observation on proportions. As another example
751 of the problem of comparison, in Table 7 we sum and average the counts and
752 percentages of case studies reported by SLSs, yet comparing across the SLSs
753 is problematic because of the misclassification by the SLSs (and also because
754 of the variation in number of research methods used by SLSs).

755 In terms of reporting trends, SLS P59 presents a bubble grid showing the
756 number of studies per research approach over a thirteen year period. The

757 publication of case studies per year has grown from one in 2007 to seven
758 in 2017 and then jumps to nineteen in 2018. Our estimates of correctness
759 in Table 7 suggest that SLS P59 is the most correct, of the eleven SLSs,
760 in its classification of case studies, though we estimate a 1/3 of SLS P59’s
761 studies are still incorrectly classified as “case study”. This misclassification
762 might have a particularly significant impact on the reported “growth” be-
763 tween 2007-2017.

764 6.3.3. *Quality and quantity of evidence*

765 Some SLSs organise their findings according to research methods. SLSs
766 P11, P44 and P45 are particularly good examples. We use one main ex-
767 ample, from SLS P44 [34], to illustrate the impact of organising findings by
768 research method, and complement that example with brief comments about
769 other examples. Before doing that, we first consider the nature of evidence
770 presented in case studies.

771 The nature of evidence reported in case studies, even the “best” case stud-
772 ies, makes it hard for an SLS to perform any kind of *quantitative* synthesis
773 other than vote-counting. It is also hard to perform analytical (also called
774 theoretical) synthesis, particularly with case studies that do not present any
775 kind of theoretical framework. A consequence is that evaluating the actual
776 impact of the mislabelling of primary studies is affected by the limits of the
777 (so-called) case studies. Incorporating studies that are not actually case stud-
778 ies further complicates the evaluation, since these studies are likely to have
779 less rigorous evidence and theories. Mixing case studies and non-case studies
780 further reduces/limits the kinds of qualitative or quantitative synthesis that
781 might be properly conducted. Mislabelling therefore distorts the body of
782 evidence, e.g. for vote-counting, “confuses” the quality of that evidence, and
783 complicates how this mixture of evidence might be synthesised.

784 Table 9 presents a simplified version of a table, Table 5, found in SLS P44
785 [34]. The original table enumerates a list of case studies that each contain
786 a dataset relating to the respective six systems presented in Table 9. We
787 simplify our version of the table to only show the number of case studies
788 reported in the SLS (see column, *SLS #CS*). We add into the table the
789 number of studies that we found to not be case studies (see *Not CS (us)*),
790 the number that the original authors of the primary studies considered were
791 not case studies (see *Not CS (PS)*) and the resultant reduction in genuine
792 case studies and, consequently, datasets (see *Reduction*). Because we only
793 *sampled* from SLS P44, and did not assess all of the primary studies in that

794 SLS, the estimates we present in Table 9 are minimum estimates. In other
795 words, *at least* four of the nine supposed case studies for JHotDraw are not
796 actually case studies. JHotDraw also provides a concrete example where SLS
797 authors are “distorting” the evidence: for four of the primary studies, their
798 authors did not consider their studies to be case studies, and we conclude the
799 same, yet the SLS authors have re-labelled these four studies as case studies.
800 For JFreeCart, the sample of nine primary studies that we drew from SLS
801 P44 did not include either of the two studies presented for JFreeCart, so we
802 cannot assess this row of the table. Finally, for JRefractory, we consider one
803 of the two case studies to not be a case study whilst the original authors
804 presented it as a case study. We recognise this difference of opinion through
805 approximating the reduction.

806 Overall, the table illustrates the distorting effect of mislabelling case stud-
807 ies on the body of evidence available for some item (in this example, software
808 systems) both in quantitative terms (see *Reduction*), but also in qualitative
809 terms, e.g., as the number of independent datasets reduces, so we have less
810 opportunity for independent corroboration or for synthesis.

Table 9: Re-analysis of Table 5 from SLS P44 [34].

| System | SLS #CS | Not CS (us) | Not CS (PS) | Reduction |
|-------------|---------|-------------|-------------|---------------------|
| JHotDraw | 9 | 4 | 4 | 4/9 (44%) |
| Eclipse | 5 | 3 | 3 | 3/5 (60%) |
| ArgoUML | 5 | 1 | 1 | 1/5 (25%) |
| Xerces-J | 3 | 2 | 2 | 2/3 (66%) |
| JFreeCart | 2 | NA | NA | NA |
| JRefractory | 2 | 1 | 0 | $\approx 1/2$ (50%) |

811 There are two other tables in SLS P44 (Table 4 and Table 10) that also
812 present findings specific to case studies. The findings of these two tables are
813 therefore also challenged by our analysis.

814 As another example, SLS P45 presents nine tables that all include infor-
815 mation about case studies, of which three tables (i.e., Table 4, Table 12 and
816 Table 13 in SLS P45) are all based only on the case studies.

817 Finally, SLS P33 states of their primary studies, “Remarkably, only 10
818 (40%) of the 25 papers that used the Case Study research method reported
819 validity threats.” We have estimated that approximately 22% in our sample
820 of the supposed case studies in SLS P33 are genuine case studies. Further
821 details are available through the tables in Supplement 6 of the online supple-
822 mentary material linked in Appendix A. Thus, SLS P33’s observation about
823 the number of studies reporting validity threats is unreliable.

824 6.3.4. *Decision making*

825 The effects of mislabelling on the quality and quantity of evidence have
826 consequences for decision making by both researchers and practitioners. Tak-
827 ing Table 9 and SLS P44 as a convenient example, there would be appear to
828 be much less actual evidence to inform decision making on these six systems
829 than might originally appear. As another example, SLS P59 presents a bub-
830 ble grid that maps research approaches, contribution types and publication
831 domains. Case studies make up 57% of the research approaches, with Sys-
832 tematic Reviews ranked second at 11%, and the remaining seven approaches
833 each with 7% or less. Lessons learned constitute 78% of the contribution
834 types, with guidelines (themselves appearing to be form of lesson learned)
835 ranked second at 7%, and the other three contributions each, again, with less
836 than 7%. The bubble grid is not a correlation, of course (and also presents
837 three factors: approach, contribution and domain) but the frequency of case
838 studies, lessons learned and guidelines again demonstrates the association
839 of case studies with aspects of decision making, in this case lessons learned
840 and guidelines. More generally, a distorted body of evidence can erroneously
841 suggest there is more evidence, or more reliable evidence, for an intervention
842 in software practice.

843 7. Limitations

844 In terms of limitations to our research, we consider four research-design
845 decisions we made that introduce limitations to our study.

846 First, we did not use multiple databases in our search for SLSs. We
847 only used one database for identifying a set of SLSs. Scopus was chosen
848 given that it has good coverage of publications from different publishers.
849 Several researchers (e.g., [23, 24, 25]) highlight Scopus as a viable option. We
850 chose to use one database because our objective was to focus on a potential
851 problem, i.e., the labelling of primary studies as “case study” by SLSs. Using
852 multiple databases would most likely not affect the outcome of our research,
853 since other databases are not likely to contain substantially different kinds
854 of SLS. We chose different search strings for the two sub-samples to mitigate
855 threats related to the search strategy, search engine, and the search string
856 formulation.

857 Ultimately, we are interested in an indicative sample of SLSs in software
858 engineering reporting primary studies being classified as case study research.

859 Since our focus is on having an indicative sample to demonstrate a concern,
860 we did not apply both search strings to all years investigated.

861 Second, the designs of our search strategies may have affected the find-
862 ings. In the first search, we looked for “case study” in the title, abstract or
863 keywords for articles published between 2010-2021. The first search strategy
864 may have selected only those articles that identified specific results relat-
865 ing to case studies, for example, the frequency of case studies in software
866 engineering.

867 To investigate and mitigate this risk, we chose a second search strategy
868 when investigating studies published in 2022. There is a difference between
869 our two samples concerning the percentages of mislabelling, although the
870 problem is the same independent of the search strategy. Thus, the percent-
871 ages may be affected by our choices, but not the overall concern of case
872 studies being mislabelled in SLSs.

873 Third, we did not systematically perform any quality assessment of the
874 SLSs identified, though we have examined the *designs* of the SLSs, and thus
875 considered some aspects of quality. Our objective was not to assess the
876 SLSs of highest-quality, but instead to assess a sample of what is published
877 in software engineering. All SLS articles analysed are published in well-
878 established software engineering journals (nine are published in I&ST which
879 explicitly encourages SLSs) and ought to be indicative of how case studies
880 are treated in SLSs. Our categorisation of the SLSs worked as a filtering
881 mechanism to ensure we looked at SLSs that included specific information
882 concerning what the community refer to as being case studies. Where it was
883 feasible to do so, we have analysed all of the categories of SLSs.

884 Fourth, we based our assessment on one definition of *case study*, and one
885 definition of *case*. We recognise there are other definitions and opinions. Our
886 starting point was Yin’s definition [13], which Runeson et al. [14] view as
887 suitable for SE research. Based on that starting point, we chose to use the
888 *refined* definition of Wohlin [3], given that it distinguishes five components
889 for a case study, and given that a checklist has been formulated from the
890 definition [4], which we used in our assessment.

891 Finally, we based our assessment on consensus discussions after independ-
892 ent review by the two authors. When we had different interpretations of an
893 article, we chose to be generous and accept a primary study as presenting a
894 case study. We have reported an agreement index for the categorisation of
895 SLSs, but not for the assessment of the primary studies since the individual
896 assessments were primarily seen as a basis for the consensus discussions. One

897 consequence of our generosity is that we may be *over-reporting* the frequency
898 of actual case studies.

899 **8. Conclusions**

900 *8.1. Summary of findings*

901 It has been known for almost twenty years that authors of primary studies
902 of software engineering do not consistently use the term and label, “case
903 study”, when reporting their studies. Systematic literature studies (SLSs),
904 such as systematic literature reviews and systematic mapping studies, provide
905 the opportunity to correct such misreporting, e.g., to re-label a primary study
906 as not a “case study” when it does not conform to the established definitions
907 for case study.

908 It was not clear whether SLSs do indeed make corrections to the labelling
909 of primary studies. We therefore investigated the following research question:
910 For those SLSs that report the empirical research methods of primary studies,
911 do those SLSs correct mislabelled case studies, propagate already mislabelled
912 case studies, or exacerbate the problem by introducing new mislabels?

913 Through a systematic search and categorisation of SLSs, we identified
914 eleven SLSs and studied 79 primary studies drawn by stratified random sam-
915 pling, from all eleven SLSs.

916 For the sample of primary studies, their authors do not claim that their
917 study is a case study in 29 of the articles. However, they do claim incorrectly
918 that it is a case study for 62% of the articles (31/50). This could be compared
919 with previous studies reporting a mislabelling of about 50% [3, 4]. However,
920 based on these 79 primary studies we find that SLSs restate the correct
921 label of “case study” in 24% of instances, but misclassify studies in 76% of
922 instances, either propagating an already incorrect label used by the respective
923 primary study (39%) or introducing a new mislabel for a primary study
924 (37%). 76% is a point estimate across all 79 studies. Considering averages
925 *per sub-sample*, we find, respectively, 60% and 81% of the primary studies are
926 mislabelled. Thus, we report a point estimate of 76% with a range estimate
927 of between 60% and 81%.

928 We found eight of the eleven SLSs did not define the term, “case study”.
929 Of the three SLSs that provide a definition, none were entirely consistent
930 with the established definitions for “case study” in software engineering re-
931 search, though two were partially consistent. Thus, part of the explanation

932 for why SLSs misclassify a primary study may be that SLS authors are not
933 considering accepted standards for research methods in their analyses.

934 We also briefly considered significant implications arising from the misla-
935 belling: it leads to difficulties with estimations and comparisons; it distorts
936 the body of evidence, “confuses” the quality of evidence, and complicates
937 the synthesis of evidence; and it undermines the confidence we can have in
938 recommendations arising from SLSs, e.g., for interventions in practice.

939 *8.2. Recommendations and checklist*

940 Credible evidence, i.e., the validity and relevance of the evidence, is es-
941 sential for decision-making [57]. One vital aspect of credible evidence is how
942 the evidence was acquired, for example, which research method was used to
943 obtain the evidence.

944 On the basis of our findings and the need for credible evidence, we make
945 several recommendations and propose a simple checklist. The checklist is in-
946 tended to complement the recommendations, since the checklist can apply to
947 editors, reviewers and readers, as well as SLS authors during the preparation
948 of their articles.

949 In terms of the recommendations:

- 950 1. Those who conduct primary studies and SLSs need to better under-
951 stand research methodology, and need to apply the terminology cor-
952 rectly, i.e., according to the definitions.
- 953 2. To support future syntheses of credible evidence in SLSs, primary study
954 authors should write for synthesis [58], e.g., stating clearly their re-
955 search method.
- 956 3. SLS authors must be evidence gate-keepers. It is essential that SLS
957 authors ensure that evidence is presented in such a way that readers
958 can determine the credibility of the evidence in their context, whether
959 being for research or in practice. Thus, SLS authors should: a) check
960 and report the research method claimed by the authors of the primary
961 studies, as well as b) check and report their own assessment of the
962 research method; and c) specify, or provide references to, the research
963 method definitions they have used.
- 964 4. SLS authors should use a sufficiently comprehensive research method
965 classification (preferably an already published one) that can accommo-
966 date the diversity of primary studies, and so avoid “forcing” studies
967 into an overly simplistic classification.

- 968 5. Reviewers and editors need to more carefully review the manuscripts of
969 both primary studies and SLSs concerning the research method claimed
970 by the authors, as the current practices in the review process tacitly
971 endorse low/er standards of quality assessment and mislabelling of re-
972 search methods.
- 973 6. To reassure themselves of (but not guarantee) the reliability and valid-
974 ity of an SLS, readers of SLSs – which may include professional practi-
975 tioners – should check, perhaps using the proposed checklist, whether
976 the SLS authors have explicitly considered the problem of mislabelling.
- 977 7. Professional practitioners, and researchers working closely with indus-
978 try should be particularly aware that any of the SLS’s recommendations
979 based on aggregations or syntheses of results from so-called case studies
980 may not be reliable or valid, in particular such studies may not have
981 been conducted in a contemporary, real-world setting.

982 In terms of the proposed checklist for checking whether SLSs have prop-
983 erly labelled primary studies:

- 984 1. Does the SLS report the research methods of each primary study?
985 2. Do the SLS authors report their own classification of the primary
986 study’s research methods?
987 3. Does the SLS formally assess the correctness of the labelling of research
988 methods by the primary-study authors?
989 4. Does the SLS show a clear mapping of primary-study authors’ and the
990 SLS authors’ classification of each primary study’s research method,
991 e.g., as a table or similar?
992 5. Does the SLS use a sufficiently comprehensive classification of research
993 methods, and preferably an already-published classification?
994 6. Does the SLS provide a set of definitions of research methods, or clearly
995 cite and apply a reference set of definitions?
996 7. Does the SLS explicitly recognise that there may be a difference of
997 classification of research method, between the primary-study authors’
998 and the SLS authors’ classification of a primary study?
999 8. Does the SLS explain the reason, or reasons, for any differences in clas-
1000 sification between primary-study authors’ and the SLS authors’ classi-
1001 fication?

1002 8.3. *Further research*

1003 In terms of further research, we consider three directions here. First, our
1004 point estimate of 24% (19% and 40% for the two sub-samples) for correctly
1005 labelled case studies in software engineering is substantially lower than our
1006 previous estimate [3, 4] of 50%, where we directly sampled primary studies.
1007 The main reason for the differences in percentages is that in the current study
1008 we look at the misclassification by SLSs. The percentage of correctly labelled
1009 case studies by the primary-study authors is 38%. One direction for further
1010 research is therefore to better understand the reasons behind mislabelling by
1011 the primary-study authors, restatement of mislabelling by SLS authors and,
1012 probably most importantly, the misclassifications added by the SLS authors.
1013 A second direction is to replicate our study, both as a literal replication (do
1014 others find the same results with the same data?), and theoretical replication
1015 (do others find corroborating results with different data?). A third direction
1016 is to continue to examine whether the phenomenon of mislabelling occurs
1017 for different kinds of labels. An obvious example would be other research
1018 methods not investigated yet, e.g., Ayala et al. [21] observe a similar situation
1019 with the mis/use of the label “experiment”, but this mislabelling may recur
1020 for other kinds of labels used by SLSs, such as different types of requirement
1021 or testing, or labels for agile practices, or types of defect.

1022 8.4. *Concluding remarks*

1023 Overall, we conclude that case studies are substantially over-reported in
1024 the literature, i.e., there are far fewer case studies actually conducted than
1025 are reported as being conducted. Our analysis is based on the definition
1026 by Wohlin [3]. As mentioned in Section 5.1, the two most common reasons
1027 for a study not being a case study are: the study is not conducted in a
1028 real-life context or it is not contemporary. These two aspects are common
1029 across the different definitions of case study research. Thus, the findings
1030 are not a consequence of our choice of definition. Furthermore, SLSs are
1031 both propagating and further exacerbating the problem of the mislabelling
1032 of primary studies as “case studies”, rather than – as we should expect of
1033 SLSs – improving the labelling of primary studies, and thus improving the
1034 body of credible evidence.

1035 **Acknowledgement**

1036 We thank the anonymous reviewers and the editor handling the paper
1037 for their time and effort. We appreciate their constructive feedback, which
1038 helped us improve the article.

1039 **Appendix A. Supplementary material**

1040 The following is the supplementary material related to this article.

- 1041 • Supplement S1 – SLSs categorised (**Link to be provided by the**
1042 **publisher**)
- 1043 • Supplement S2 – Studies excluded from our SLS analyses (**Link to be**
1044 **provided by the publisher**)
- 1045 • Supplement S3 – Explanation of batches (**Link to be provided by**
1046 **the publisher**)
- 1047 • Supplement S4 – Comments on the mislabelling of case studies (**Link**
1048 **to be provided by the publisher**)
- 1049 • Supplement S5 – Listing of primary studies (**Link to be provided by**
1050 **the publisher**)
- 1051 • Supplement S6 – Full analysis of primary studies (**Link to be pro-**
1052 **vided by the publisher**)

1053 **References**

- 1054 [1] C. Zannier, G. Melnik, F. Maurer, On the success of empirical studies
1055 in the international conference on software engineering, in: Proceedings
1056 International Conference on Software Engineering, 2006, pp. 341–350.
- 1057 [2] P. Runeson, M. Höst, Guidelines for conducting and reporting case study
1058 research in software engineering, Empirical Software Engineering 14 (2)
1059 (2009) 131.
- 1060 [3] C. Wohlin, Case study research in software engineering—it is a case,
1061 and it is a study, but is it a case study?, Information and Software
1062 Technology 133 (2021) 106514.

- 1063 [4] C. Wohlin, A. Rainer, Is it a case study?—a critical analysis and guid-
1064 ance, *Journal of Systems and Software* 192 (2022) 111395.
- 1065 [5] R. K. Yin, *Case study research: Design and methods*, 3rd Edition, Sage
1066 Publications, 2003.
- 1067 [6] A. Rainer, C. Wohlin, Case study identification: A trivial indicator
1068 outperforms human classifiers, *Information and Software Technology* 161
1069 (2023) 107252.
- 1070 [7] S. Jalali, C. Wohlin, Systematic literature studies: database searches vs.
1071 backward snowballing, in: *Proceedings of the ACM-IEEE international
1072 symposium on Empirical software engineering and measurement*, 2012,
1073 pp. 29–38.
- 1074 [8] C. Wohlin, Guidelines for snowballing in systematic literature studies
1075 and a replication in software engineering, in: *Proceedings of the 18th
1076 international conference on evaluation and assessment in software engi-
1077 neering*, 2014, pp. 1–10.
- 1078 [9] D. Badampudi, C. Wohlin, K. Petersen, Experiences from using snow-
1079 balling and database searches in systematic literature studies, in: *Pro-
1080 ceedings of the 19th international conference on evaluation and assess-
1081 ment in software engineering*, 2015, pp. 1–10.
- 1082 [10] B. A. Kitchenham, D. Budgen, P. Brereton, *Evidence-based Software
1083 Engineering and Systematic Reviews*, Vol. 4, CRC press, 2015.
- 1084 [11] K. Petersen, S. Vakkalanka, L. Kuzniarz, Guidelines for conducting sys-
1085 tematic mapping studies in software engineering: An update, *Informa-
1086 tion and Software Technology* 64 (2015) 1–18.
- 1087 [12] V. Basili, L. Briand, D. Bianculli, S. Nejati, F. Pastore, M. Sabetzadeh,
1088 *Software engineering research and industry: a symbiotic relationship to
1089 foster impact*, *IEEE Software* 35 (5) (2018) 44–49.
- 1090 [13] R. K. Yin, *Case study research and applications: Design and methods*,
1091 6th Edition, Sage Publications, 2018.
- 1092 [14] P. Runeson, M. Höst, A. Rainer, B. Regnell, *Case Study Research
1093 in Software Engineering: Guidelines and Examples*, Wiley Publishing,
1094 2012.

- 1095 [15] B. Kitchenham, L. Pickard, S. Pfleeger, Case studies for method and
1096 tool evaluation, *IEEE Software* 12 (4) (1995) 52–62.
- 1097 [16] J. M. Verner, J. Sampson, V. Tasic, N. A. A. Bakar, B. A. Kitchenham,
1098 Guidelines for industrially-based multiple case studies in software en-
1099 gineering, in: *Proceedings Third International Conference on Research*
1100 *Challenges in Information Science*, 2009, pp. 313–324.
- 1101 [17] I. Benbasat, D. K. Goldstein, M. Mead, The case research strategy in
1102 studies of information systems, *MIS Quarterly* 11 (3) (1987) 369–386.
- 1103 [18] C. Robson, *Real World Research : A Resource for Social Scientists and*
1104 *Practitioner-researchers*, Wiley Publishing, 2002.
- 1105 [19] C. Baker, J. Wuest, P. N. Stern, Method slurring: the grounded the-
1106 ory/phenomenology example, *Journal of Advanced Nursing* 17 (11)
1107 (1992) 1355–1360.
- 1108 [20] N. Haslam, Concept creep: Psychology’s expanding concepts of harm
1109 and pathology, *Psychological Inquiry* 27 (1) (2016) 1–17.
- 1110 [21] C. Ayala, B. Turhan, X. Franch, N. Juristo, Use and misuse of the term
1111 “experiment” in mining software repositories research, *IEEE Transactions*
1112 *on Software Engineering* 48 (11) (2021) 4229–4248.
- 1113 [22] K.-J. Stol, P. Ralph, B. Fitzgerald, Grounded theory in software engi-
1114 neering research: A critical review and guidelines, in: *Proceedings of*
1115 *the 38th International Conference on Software Engineering*, 2016, pp.
1116 120–131.
- 1117 [23] O. Dieste, A. Grimán, N. Juristo, Developing search strategies for de-
1118 tecting relevant experiments, *Empirical Software Engineering* 14 (5)
1119 (2009) 513–539.
- 1120 [24] V. Garousi, M. V. Mäntylä, Citations, research topics and active coun-
1121 tries in software engineering: A bibliometrics study, *Computer Science*
1122 *Review* 19 (2016) 56–77.
- 1123 [25] E. Mourão, J. F. Pimentel, L. Murta, M. Kalinowski, E. Mendes,
1124 C. Wohlin, On the performance of hybrid search strategies for system-
1125 atic literature reviews in software engineering, *Information and Software*
1126 *Technology* 123 (2020) 106294.

- 1127 [26] D. Budgen, P. Brereton, N. Williams, S. Drummond, The contribution
1128 that empirical studies performed in industry make to the findings of sys-
1129 tematic reviews: A tertiary study, *Information and Software Technology*
1130 94 (2018) 234–244.
- 1131 [27] S. M. Melo, J. C. Carver, P. S. Souza, S. R. Souza, Empirical research on
1132 concurrent software testing: A systematic mapping study, *Information*
1133 *and Software Technology* 105 (2019) 226–251.
- 1134 [28] A. Ampatzoglou, I. Stamelos, Software engineering research for com-
1135 puter games: A systematic review, *Information and Software Technol-*
1136 *ogy* 52 (9) (2010) 888–901.
- 1137 [29] D. Tofan, M. Galster, P. Avgeriou, W. Schuitema, Past and future of
1138 software architectural decisions—a systematic mapping study, *Informa-*
1139 *tion and Software Technology* 56 (8) (2014) 850–872.
- 1140 [30] N. H. Bakar, Z. M. Kasirun, N. Salleh, Feature extraction approaches
1141 from natural language requirements for reuse in software product lines:
1142 A systematic literature review, *Journal of Systems and Software* 106
1143 (2015) 132–149.
- 1144 [31] W. Bissi, A. G. S. S. Neto, M. C. F. P. Emer, The effects of test driven
1145 development on internal quality, external quality and productivity: A
1146 systematic review, *Information and Software Technology* 74 (2016) 45–
1147 54.
- 1148 [32] J. L. Barros-Justo, F. Pinciroli, S. Matalonga, N. Martínez-Araujo,
1149 What software reuse benefits have been transferred to the industry?
1150 a systematic mapping study, *Information and Software Technology* 103
1151 (2018) 1–21.
- 1152 [33] J. Kroll, I. Richardson, R. Prikladnicki, J. L. Audy, Empirical evidence
1153 in follow the sun software development: A systematic mapping study,
1154 *Information and Software Technology* 93 (2018) 30–44.
- 1155 [34] F. Wedyan, S. Abufakher, Impact of design patterns on software quality:
1156 a systematic literature review, *IET Software* 14 (1) (2020) 1–17.

- 1157 [35] A. E. Chacón-Luna, A. M. Gutiérrez, J. A. Galindo, D. Benavides, Em-
1158 pirical software product line engineering: a systematic literature review,
1159 Information and Software Technology 128 (2020) 106389.
- 1160 [36] Ö. Uludağ, P. Philipp, A. Putta, M. Paasivaara, C. Lassenius,
1161 F. Matthes, Revealing the state of the art of large-scale agile devel-
1162 opment research: A systematic mapping study, Journal of Systems and
1163 Software (2022) 111473.
- 1164 [37] A. R. Amna, G. Poels, Ambiguity in user stories: A systematic literature
1165 review, Information and Software Technology 145 (2022) 106824.
- 1166 [38] N. Dissanayake, A. Jayatilaka, M. Zahedi, M. A. Babar, Software se-
1167 curity patch management-a systematic literature review of challenges,
1168 approaches, tools and practices, Information and Software Technology
1169 144 (2022) 106771.
- 1170 [39] C. Robson, Small-scale evaluation: Principles and practice, Sage Publi-
1171 cations Ltd, 2017.
- 1172 [40] C. Wohlin, P. Runeson, M. Höst, M. C. Ohlsson, B. Regnell, A. Wesslén,
1173 Experimentation in software engineering, Springer Science & Business
1174 Media, 2012.
- 1175 [41] S. Easterbrook, J. Singer, M.-A. Storey, D. Damian, Selecting empiri-
1176 cal methods for software engineering research, in: F. Shull, J. Singer,
1177 D. I. Sjøberg (Eds.), Guide to advanced empirical software engineering,
1178 Springer, 2008, pp. 285–311.
- 1179 [42] Z. Li, P. Liang, P. Avgeriou, Application of knowledge-based approaches
1180 in software architecture: A systematic mapping study, Information and
1181 Software Technology 55 (5) (2013) 777–794.
- 1182 [43] F. Elberzhager, A. Rosbach, J. Münch, R. Eschbach, Reducing test
1183 effort: A systematic mapping study on existing approaches, Information
1184 and Software Technology 54 (10) (2012) 1092–1106.
- 1185 [44] J. Bailey, D. Budgen, M. Turner, B. Kitchenham, P. Brereton,
1186 S. Linkman, Evidence relating to object-oriented software design: A
1187 survey, in: First International Symposium on Empirical Software Engi-
1188 neering and Measurement, 2007, pp. 482–484.

- 1189 [45] R. Wieringa, N. Maiden, N. Mead, C. Rolland, Requirements engineer-
1190 ing paper classification and evaluation criteria: a proposal and a discus-
1191 sion, *Requirements Engineering* 11 (1) (2006) 102–107.
- 1192 [46] K. Petersen, R. Feldt, S. Mujtaba, M. Mattsson, Systematic mapping
1193 studies in software engineering, in: *12th International Conference on*
1194 *Evaluation and Assessment in Software Engineering*, 2008, pp. 1–10.
- 1195 [47] V. Berg, J. Birkeland, A. Nguyen-Duc, I. O. Pappas, L. Jaccheri, Soft-
1196 ware startup engineering: A systematic mapping study, *Journal of Sys-*
1197 *tems and Software* 144 (2018) 255–274.
- 1198 [48] P. Rodríguez, A. Haghightakhah, L. E. Lwakatare, S. Teppola, T. Suo-
1199 malainen, J. Eskeli, T. Karvonen, P. Kuvaja, J. M. Verner, M. Oivo,
1200 Continuous deployment of software intensive products and services: A
1201 systematic mapping study, *Journal of Systems and Software* 123 (2017)
1202 263–291.
- 1203 [49] M. Unterkalmsteiner, T. Gorschek, A. M. Islam, C. K. Cheng, R. B.
1204 Permadi, R. Feldt, Evaluation and measurement of software process im-
1205 provement—a systematic literature review, *IEEE Transactions on Soft-*
1206 *ware Engineering* 38 (2) (2011) 398–424.
- 1207 [50] L. Chen, M. A. Babar, A systematic review of evaluation of variabil-
1208 ity management approaches in software product lines, *Information and*
1209 *Software Technology* 53 (4) (2011) 344–362.
- 1210 [51] P. Brereton, B. Kitchenham, D. Budgen, Z. Li, Using a protocol tem-
1211 plate for case study planning, in: *12th International Conference on Eval-*
1212 *uation and Assessment in Software Engineering*, 2008, pp. 1–8.
- 1213 [52] M. Höst, P. Runeson, Checklists for software engineering case study
1214 research, in: *First International Symposium on Empirical Software En-*
1215 *gineering and Measurement*, IEEE, 2007, pp. 479–481.
- 1216 [53] R. K. Yin, *Case study research: Design and methods*, 5th Edition, Sage
1217 Publications, 2009.
- 1218 [54] R. L. Glass, I. Vessey, V. Ramesh, Research in software engineering: an
1219 analysis of the literature, *Information and Software Technology* 44 (8)
1220 (2002) 491–506.

- 1221 [55] C. Wohlin, A. Aurum, Towards a decision-making structure for selecting
1222 a research design in empirical software engineering, *Empirical Software*
1223 *Engineering* 20 (2015) 1427–1455.
- 1224 [56] D. S. Cruzes, T. Dybå, Research synthesis in software engineering: A
1225 tertiary study, *Information and Software Technology* 53 (5) (2011) 440–
1226 455.
- 1227 [57] C. Wohlin, A. Rainer, Challenges and recommendations to publishing
1228 and using credible evidence in software engineering, *Information and*
1229 *Software Technology* 134 (2021) 106555.
- 1230 [58] C. Wohlin, Writing for synthesis of evidence in empirical software en-
1231 gineering, in: *Proceedings of the 8th ACM/IEEE International Symposi-*
1232 *um on Empirical Software Engineering and Measurement*, 2014, pp.
1233 46:1–4.