# An Empirically Evaluated Checklist for Surveys in Software Engineering

Jefferson Seide Molléri[1], Kai Petersen[1], Emilia Mendes[1]

**1 BTH - Blekinge Tekniska Högskola**

**\* jefferson.molleri@bth.se**

## Abstract

**Context:** Over the past decade Software Engineering research has seen a steady increase in survey-based studies, and there are several guidelines providing support for those willing to carry out surveys. The need for auditing survey research has been raised in the literature. Checklists have been used to assess different types of empirical studies, such as experiments and case studies.

**Objective:** This paper proposes a checklist to support the design and assessment of survey-based research in software engineering grounded in existing guidelines for survey research. We further evaluated the checklist in the research practice context.

**Method:** To construct the checklist, we systematically aggregated knowledge from 14 methodological papers supporting survey-based research in software engineering. We identified the key stages of the survey process and its recommended practices through thematic analysis and vote counting. To improve our initially designed checklist we evaluated it using a mixed evaluation approach involving experienced researchers.

**Results:** The evaluation provided insights regarding limitations of the checklist in relation to its understanding and objectivity. In particular, 19 of the 38 checklist items were improved according to the feedback received from its evaluation. Finally, a discussion on how to use the checklist and what its implications are for research practice is also provided.

**Conclusion:** The proposed checklist is an instrument suitable for auditing survey reports as well as a support tool to guide ongoing research with regard to the survey design process.

## Introduction

A survey is a widely deployed research method in the area of Software Engineering (SE) and an increase in its usage has been highlighted by, e.g., Punter et al. [41]. Its purpose is to investigate a population, in order to construct explanatory models [3, 50] or to validate knowledge [25, 30]. Survey research is often employed when there is a need to study a large set of variables [50] or to perform a retrospective analysis [39]. It may be used to draw conclusions based on both quantitative and qualitative data [10].

Researchers have highlighted various challenges during the survey process. Common challenges are the formulation of questions [46], so to avoid shortcomings (e.g., introducing bias inside questions [50]), and the identification of invalid responses [51]. Other challenges are related to the recruitment of participants, such as how to obtain a sufficient number of responses and how to prevent high drop-out rates [1, 21].

The need for improving the completeness of reporting survey-based research, in particular with respect to the definition of the population and sampling strategies is evidenced by Stavru [45]. Furthermore, Stavru pointed out a lack of checklists for auditing surveys in SE which could be of help to both researchers conducting survey research as well as to those evaluating and reviewing the research.

Motivated by these needs, we employed an empirical approach to constructing and evaluating an assessment checklist[1]. Such approach comprises of two steps (see Sections 2 and 3):

First, we detail the process to **construct a checklist** to assess survey research in SE. The method used to derive the checklist was guided by two principles: (a) identify existing guidelines for survey research; in the context of SE, 14 methodological papers have been considered; (b) elicit the process stages, recommended practices, and related rationales. Those rationales support the cost-effectiveness analysis of employing a set of related practices.

The method for systematically deriving the checklist was based on thematic analysis [9]. Vote counting was applied to the themes identified in order to compute the frequency in which they occurred. Further, a co-occurrence was obtained through a relationship matrix relating different categories (e.g., practices versus rationales).

Later, we **evaluated the checklist** by using a mixed approach [25]. The evaluation process involved two distinct phases: (a) to apply the checklist on a set of published survey reports and register the assessment scores, and (b) to verify the results of this assessment with the corresponding authors of those survey reports.

The assessment produced a compliance coefficient for the selected studies in relation to each of the checklist items. We further investigated the authors' feedback in order to understand patterns we identified in the assessment scores. We also collected and addressed suggestions from the experts to improve the checklist instrument.

The remainder of the paper is structured as follows: Section 1 describes the background and related work. Section 2 details the systematic approach we used to construct the checklist. The evaluation of the checklist in research practice context is presented in Section 3. Section 4 discusses the findings and finally, and Section 5 concludes the paper.

# 1  Background & Related Work

We first present existing survey guidelines that are subject-independent or that have been proposed in other fields. Thereafter, an overview of SE specific guidelines is provided. After giving an overview of the guidelines we describe the literature on survey assessment. Finally, we looked into checklists proposed to support Empirical Software Engineering (ESE).

## 1.1  Survey process and guidelines

Survey as a research method has been established in social research for half a century. It has been employed in several academic fields, such as health care, politics, psychology, and sociology [8]. As a consequence, methodological knowledge on surveys has been published first in these fields.

As the survey research method matured, cross-field guidelines appeared (e.g., [2, 18, 19]). These publications aimed to provide methodological support independent from the subject of research. Nevertheless, it is not uncommon for their mentioned practices to focus on the social aspects of research.

---

[1]The resulting instrument is further detailed in Appendix A.2.

Those guidelines [2, 18, 19] describe a survey-research process comprising a set of stages, such as question design, sampling, data collection, instrument evaluation, measuring and data analysis [19]. Survey research is acknowledged for being flexible, although the process stages are often conducted sequentially.

It is worth mentioning that the survey process is a complete research method (i.e., including planning, execution, analysis, and reporting); the survey data collection instrument is called a questionnaire.

In addition to describing the process, the guidelines also recommend best practices based on desirable attributes for high-quality surveys. Such quality is based on evaluation of the produced evidence (e.g., precision, credibility), ethical issues (e.g., consent, privacy) and mitigating validity threats (e.g., sample error, non-responses) [19, 31].

## 1.2 Survey guidelines in Software Engineering

The need for specific and tailored guidelines to conduct empirical research in the context of SE has been pointed out, e.g. [36, 44, 49]. This demand is especially relevant to formal experiments and case studies, due to the popularity of such methods, but also applies to survey-based research. Methodological support for surveys in SE first appeared around the 1990s [35].

Three main guidelines [24, 30, 32] detail the survey research process in the SE field. They jointly provide a comprehensive structure for the research process, despite differing slightly from each other. Major differences are in relation to the breakdown structure of process stages and the recommended practices provided.

Besides these three main publications, a series of additional studies extend the guidance to particular stages of the survey process. For example, the challenges of identifying the target audience and establishing a sampling frame are discussed in [11–14]. The recommended practices in this set of papers are complementary, although some partially overlap.

Other studies provide lessons learned from carrying out the process in different contexts:

- Punter et al. [41] focus on self-administered online surveys and address issues such as monitoring real-time responses, identifying the reasons for dropouts and encouraging participants to complete a survey instrument;

- Ciolkowski et al. [5] addresses practical issues related to the process itself, such as managing resources and ensuring that deviations do not threaten the completion of the entire process; and

- Cater et al. [4] address replication challenges, such as updating a survey instrument and comparing the results.

Additional references for survey-based research are provided in our previous works [34, 35].

## 1.3 Survey assessment

When reviewing existing guidelines (see Section 1.2) we found out that several researchers highlighted the need for an instrument to audit survey research in SE context. This need is further stressed by the lack of reporting of the employed criteria to assess survey research [45].

Stavru's work [45] provides a critical review of surveys in the area of agile software development. In order to carry out the review, Stavru used 21 criteria by which the

thoroughness in reporting surveys was assessed. These criteria were extracted from different sources, cf. [5, 24, 30, 33, 40, 43]. Note that the method of eliciting the criteria was not detailed.

Stavru also highlighted that the different criteria were not equally important, and rated them on a scale from one to five. The most important criteria that ought to be documented were:

- Sampling frame, method, and size

- Response rate

- Assessment of a survey's trustworthiness

- Survey process

- Conceptual model comprising of the constructs investigated (e.g., variables and their relations)

- Target population

- Questionnaire design

## 1.4  Checklists in Software Engineering

Checklists have been proposed for various research methods with a specific focus on their usage in the SE context. Looking at the ways in which checklists were built, researchers most often based the construction of a new checklist upon existing ones (cf. [22, 48]).

As an example, Kitchenham et al. [27] combined two checklists [17, 28] to assess experiments and to evaluate whether researchers may use them objectively. Their findings indicate that a larger number of reviewers was needed (eight) to reliably assess studies using their checklist, which could be improved by having researchers conduct reviews in pairs (cf. [27]). Additional checklists proposed for assessing experiments are, e.g., [23, 27, 50].

Höst and Runeson [22] put forward a checklist for case study research, divided according to the research stages, including design, preparation for data & evidence gathering, data analysis, as well as reporting. To ease a reviewer's task, a condensed checklist abstracting the original checklist has been suggested to reduce the number of items to be checked. Additional checklists for case study research in SE are found in, e.g., [26, 37, 50].

Wieringa [47] observed that the individual checklists with the same focus differed, which may result in confusions for reviewers. The author highlights the need to find common checklist items across research types as they may share specific aspects. Thus, Wieringa et al. [48] used existing checklists (e.g., [22, 23, 42]) for experiments and case studies as a basis to synthesize an unified checklist. Later, the authors evaluated their checklist by having them used by PhD students and researchers in different research groups, as well as by conference participants.

Stavru's [45] filled a gap in the existing body of knowledge by complementing the set of available checklists with a set of criteria for assessing survey research. No other checklists to assess surveys were identified in our systematic literature search (cf. [34]).

Great emphasis was placed upon (a) basing the checklist on existing literature, and (b) following a systematic approach to eliciting checklist items [45, 47]. Thus, our work complements the above-mentioned by deriving and evaluating an assessment checklist grounded in existing guidelines for survey research.

# 2 Step 1. Construction of the checklist

The first step of our research approach entailed the systematic construction of the checklist. Three sub-contributions are made that ultimately lead to the checklist proposed:

C1 *Consolidation of survey processes and decision points:* We present a consolidated survey process based on existing guidelines. Key decisions points and implications of decision-making are highlighted. For example, a key decision in a survey process is the type of sampling used, which impacts participant recruitment and data analysis. Our checklist has to be adapted depending on the decisions taken.

C2 *Extraction of recommended practices and their mapping to the survey process:* We extracted the recommended practices to be carried out during a survey research process, which were later mapped to the research process identified in C1. Mapping the practices to the main stages aids researchers in the planning of surveys, as it indicates in which process step a practice is executed and where its impact needs to be considered.

C3 *Extraction of rationales for the recommended practices:* The reasons for considering existing survey research practices should be motivated by a rationale, thus making the value of adopting such practices explicit. This is particularly pressing because a survey's cost-effectiveness is an important consideration. Thus, understanding the rationales for the recommended practices supports the cost analysis of a practice and its effectiveness (i.e., the rationale regarding the value a given practice adds to the survey research).

## 2.1 Method

### 2.1.1 Research questions

We formulated three research questions corresponding respectively to each of the three contributions stated above, as follows:

RQ1 Which stages and key decisions are specified for the survey process (C1)?

RQ2 Which practices are suggested and how do they map to the stages of the research process (C2)?

RQ3 What is the rationale for conducting the respective recommended practices (C3)?

## 2.2 Study identification and selection

In order to select an appropriate set of primary studies, we used evidence from our previous studies [34, 35] identifying methodological papers for survey-based research in SE.

A set of eight papers provided the main guidelines covering all the stages of a survey research process [24, 30, 32]. We assume that these core papers were likely to hold all the information needed to derive our checklist, i.e., recommended practices and reasons for their adoption.

Given that only a few papers covering the complete methodology may not cover recommended practices and reasons for their adoption sufficiently, we completed the core set with six additional additional supporting papers [4, 5, 11–13, 41] addressing specific stages of survey research such as sampling, instrument design, and validation, recruitment and response management.

$$c := n12 / (n1 + n2) - n12$$
$$c := 19 / (30 + 41) - 19$$
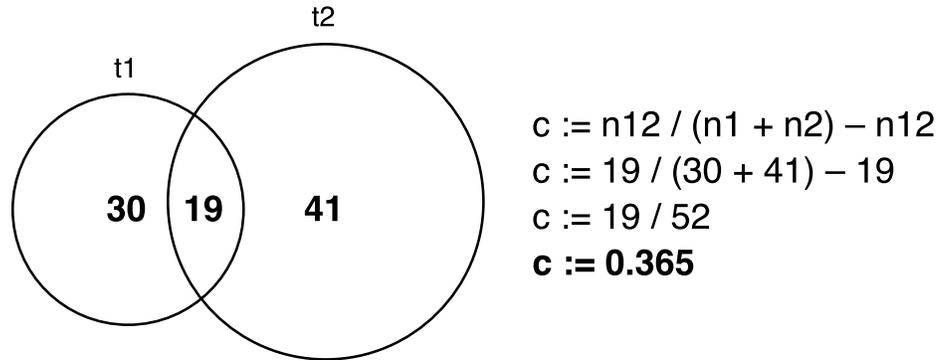$$c := 19 / 52$$
$$\mathbf{c := 0.365}$$

**Figure 1.** Example for computation of the co-occurrence coefficient, given that t1 occurs 30 times in the data set, t2 occurs 41 times, and they simultaneously occur 19 times.

### 2.2.1 Data extraction and analysis

We employed a systematic process based on thematic synthesis [9] to extract data from and to analyze the primary studies. First, we aggregated the papers in a common list using Atlas.ti [20] - a qualitative data analysis software - for text interpretation. Then, we incrementally read the papers, collected text segments and aggregated them into themes. The themes were classified into three major categories:

1. Process stages, e.g., data analysis.

2. Recommended practices, e.g., identify reasons for non-responses.

3. Rationale attributes, e.g., representativeness.

The terminology for initial codes is derived from the three guidelines analyzed first (i.e., [24, 30, 32]). Later, we iteratively improved and updated the initial coding according to the different views presented by the additional sources. The segments are characterized by a level of granularity of one paragraph, notwithstanding a single paragraph is often associated with more than one theme. Paragraphs containing no relevant information were associated with no theme.

Successive iterations of the process further refined the theme set, by combining or merging synonyms and aggregating related themes into families (e.g., representativeness is part of the external validity rationale). We also removed duplicates and combined successive occurrences of the same theme in larger segments, thus comprising several paragraphs.

Finally, we computed a co-occurrence coefficient [7] to analyze how frequently two related terms occurred alongside each other. This coefficient is calculated as follows: $c := n_{1,2}/(n_1 + n_2) - n_{1,2}$, whereas $n_1$ and $n_2$ are the vote-counting frequencies of two themes $t_1$ and $t_2$ respectively, and $n_{1,2}$ is the joint frequency, i.e., how many times the two themes co-occur. An example of the coefficient computation is given in Figure 1.

The resulting relationships are presented in a co-occurrence matrix [20], where the cells are filled in grey-tones according to the coefficient value (see Section 2.4). Darker cells represent a stronger co-occurrence between two themes. We opted for normalizing the themes in each matrix row since our analysis relates mainly to only comparing themes within the same category. Our normalized coefficient range from 0-100, whereas 100 relates to the maximum occurrence in the corresponding group, and 0 corresponds to no occurrence.

## 2.3 Threats to validity

**Construct validity.** To ensure a similar understanding and reduce research bias on the thematic analysis, we piloted the coding process between the three authors. The results imply a fair agreement, i.e., in average, 46.5% of the themes are similar, although worded differently. Further, based on the reflections of the pilot study, the first author coded the remaining papers. The co-occurrence coefficient used in the analysis takes into consideration the position and size of sentences but is prone to non-significant values when comparing themes that differ largely in size [20]. We partially address this potential bias by normalizing the values within the same row.

**Internal validity.** Internal validity relates to factors affecting the outcome of the study not accounted for by the researchers. One threat is the bias in interpreting the findings. Hence, at each stage of the research, the intermediate results were discussed among the researchers (observer triangulation).

Our data set consists of 14 papers gathered from a previous literature study [34]. We relied mostly on the original study design to the search and selection stages, using its reported evidence to collect our relevant set. We employed structured reading and coding to analyze the data set, producing themes and higher-level categories. The first author conducted the data extraction and analysis, further discussing the resulting themes with the other co-authors. We trust that this iterative process minimized the judgment bias of a solo researcher.

**Conclusion validity.** One threat to conclusion validity is whether the data based on which the survey was created is complete. The additional 6 papers in our selection (see Section 2.2) complemented the study results with 10 extra practices and 1 rationale. Those extra themes are related to particular challenges a researcher may face during the process, namely a large sample frame [11, 12], managing online surveys [41], and survey replication [4]. By adding guidelines very specific to the individual stages of survey research increased confidence in the results, despite being limited by the availability of the literature for each stage.

**External validity.** The resulting themes and frequencies were extracted from relevant methodological guidance for SE. However, we cannot assume that the practices and rationales identified are only important for this field. Moreover, there is the possibility of identifying a valid theme outside of our data set, e.g., a non-selected paper or the practical experiences of a researcher. We, therefore, conducted a evaluation of our proposed checklist with SE researchers, thus investigating its appropriateness and identifying potential improvements (see Section 3).

## 2.4 Results

The following section is structured according to the three contributions proposed in Section 2.1. Each section is, in turn, broken down into its units of analysis (i.e., decision-making points, process stages, and rationale aggregation).

**C1. Survey Research Process**

Figure 2 presents an aggregation of the survey research processes described by Kitchenham & Pfleeger [30], Kasunic [24], and Linåker *et al.* [32]. Although the main stages (and terminology) slightly differ among the guidelines, the processes are similar and follow a sequential flow. We adopted the view from [30] to describe the execution phase comprising two stages, one for recruiting participants and one for administering responses.

We also identified key decision-making points during the process, two of which should be addressed during the sampling stage (i.e., D1 and D2) and two others during
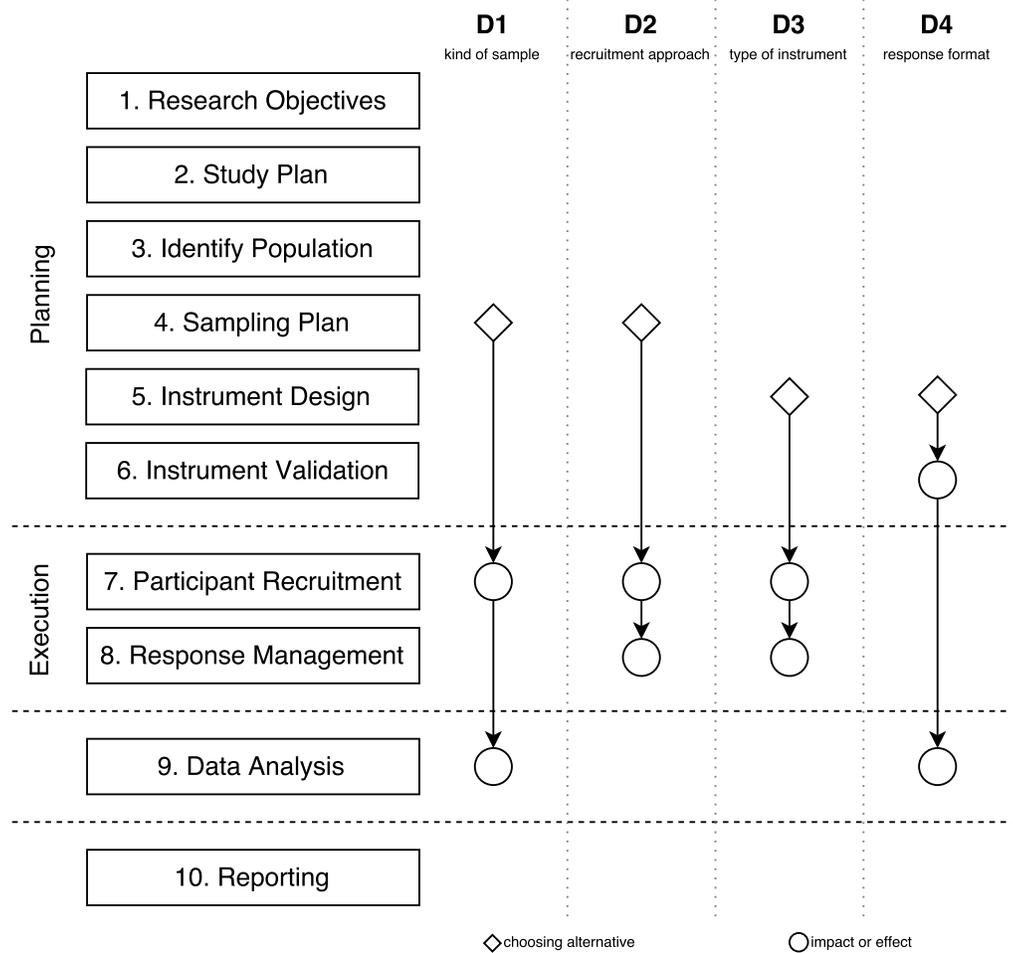
**Figure 2.** Model for the survey process, together with the decision points that could impact on the research. The diamond symbols (⋄) highlights the stages in which the conditional alternatives should be chosen, and the circles (○) mark the stages affected by those decisions.

the instrument design stage (i.e., D3 and D4). Those conditional nodes require researchers to make decisions regarding a survey's research design that can potentially impact the subsequent stages.

D1 *What kind of sample is selected?*
Depending on the strategy for selecting respondents, the researcher can choose a **P) probabilistic** (e.g., random selection) or **NP) non-probabilistic** (e.g., convenience, quota, snowballing) sample. This decision mainly affects the data analysis methods (as probabilistic samples are meant to be generalizable) and recruitment approaches (e.g., random selection) employed.

D2 *How are the participants recruited?*
On the one hand, **SS) self-selection** approaches allow for potential respondents to volunteer themselves which may introduce biases in the interpretation of the data.
On the other hand, **PS) personalized selection** (such as invitation letters and more rarely authorization codes) require specific actions for the recruitment and management of the responses. This decision-making point is often interdependent of the kind of sample (D1).

D3 *What type of survey instrument is designed?*
**SA) Self-administrated** surveys are mainly distributed in the form of Web pages or printed questionnaires, thus the respondents fill out the data themselves. **IA) Interviewer-administrated** surveys include face to face or phone interviews where respondents provide the information to a researcher, who records the data. This decision not only drives the instrument design but can also heavily impact the execution stages (i.e., recruitment and response management).

D4 *What response formats are collected?*
Question structure types could be **OE) open-ended** and **CE) close-ended**. Open-ended questions are less restrictive allowing for respondents to use their own words, whereas close-ended questions are represented by scales that can be easily quantified. This decision determines the data analysis methods employed, i.e., qualitative or quantitative approaches. Often survey instruments include a mix of both question types, thus requiring both analysis approaches.

## C2. Recommended Practices

We identified a list of recommended practices for the survey research process (see Figure 3). A stronger shade means that a higher number of citations were identified, however, the importance or the actual extent to which the practices and stages are related were not analyzed.

Some practices are likely more relevant to a stage other than pointed by the cell shade. As an example, the recommended practice "keep a diary/log book (P1)" is strongly related to stage S10. Reporting, although this practice is often initiated in stage S1. Research Objectives.

The practices are mostly focused on specific stages of the survey process. However, they can also influence preceding and succeeding stages. The effect on the preceding stages is usually in the form of planning the strategy to be carried out. The impact on succeeding stages usually entails follow-up actions and the consequences of a given decision (see Section 2.4).

1. *Research objectives:* Survey-based research is motivated by a specific goal. Thus it is important to state the research questions that correspond to such goal. The two
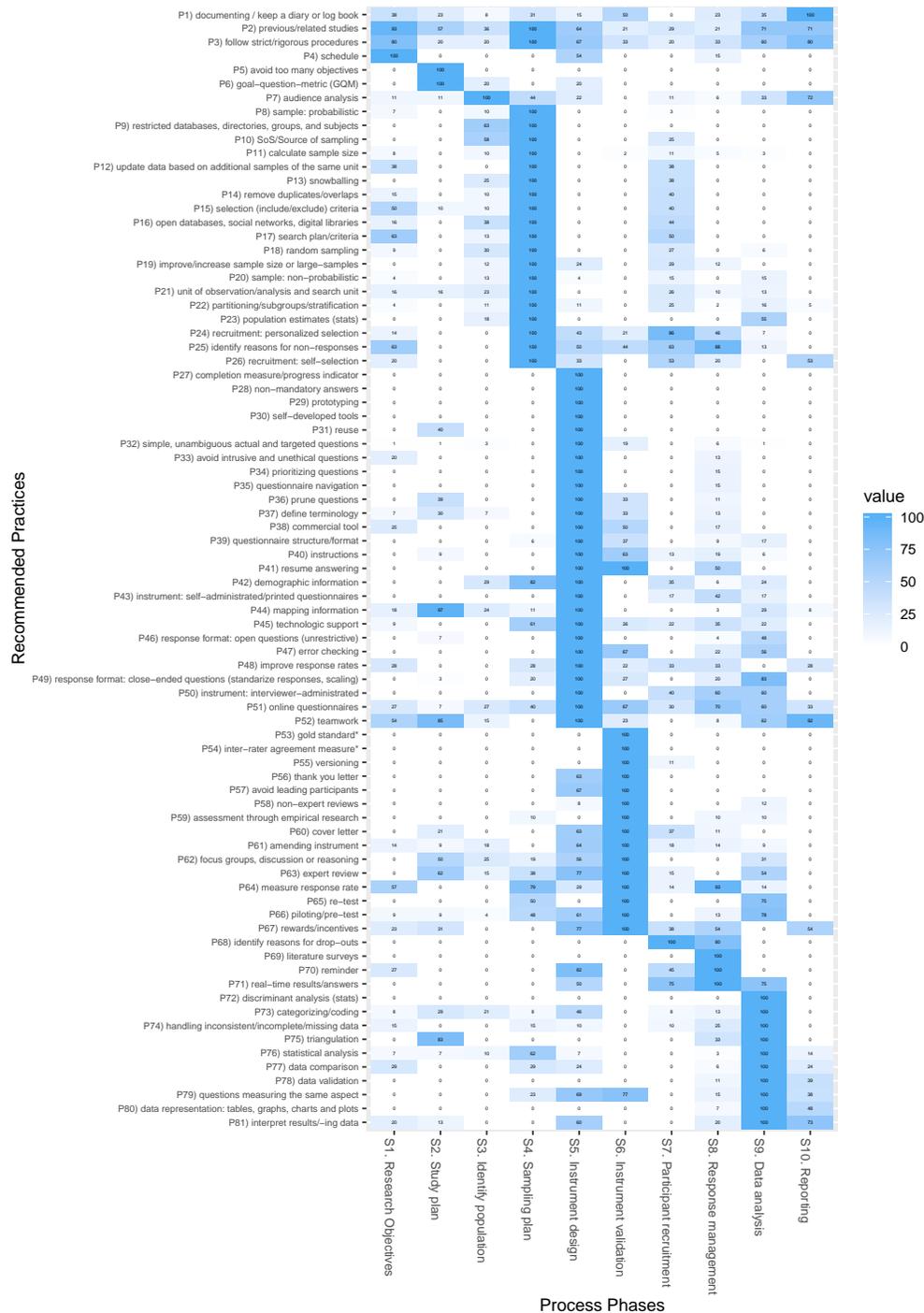
**Figure 3.** Co-occurrence matrix of recommended practices according to the process stages. Each cell contains a normalized coefficient, i.e., the highest co-occurrence value in each row is assigned a value 100 whereas the lowest value is 0. Cell shading illustrates the strength of the normalized co-occurrence coefficient, ranging from white (0) to blue (100).

main recommendations related to this stage are P1) to limit the scope, as this could impact upon the survey's complexity, and P2) to apply the goal-question-metric (GQM) approach to define its objectives. Moreover, the questionnaire items and collected data should be mapped to the research questions (P44).

2. *Study plan:* The need for designing the survey research is set at the beginning of the process, often along with the research questions [24]. The main suggested practices for this stage are to: P3) investigate related work; P4) define a set of procedures to guide the process; P5) develop a schedule plan for the stakeholders; and P6) start a diary or log book. The study plan should then be iteratively revised during the process, and the updates recorded in the log book (P1). This information is specially required for the reporting stage, at the end of the process.

3. *Identify and characterize the population:* Audience analysis (P7) is often employed to identify and select the characteristics of the population addressed by the research. This task has a strong effect on the sampling stage, in which the sources of sampling (P10) should be defined. Surveys often target potential participants at open databases (P16), but could employ restricted databases (P10) as an alternative or complementary source of sampling. Restricted databases should be investigated prior to the sampling stage.

4. *Sampling plan:* It is often employed in order to sample the population representatively. A sample plan should contain the sources of sampling (P10), units of observation and search unit (P21). The type of sample (P8 and P20) should potentially lead the decision for the data analysis methods employed. Other essential aspects to be considered are the P11) size of the sample and P19) how to manage large samples.

   Additional practices for this stage include to P14) remove the redundant units; P15) apply criteria for selecting the units of observation; P17) plan the retrieval of search units; and P22) partition the population according to the chosen characteristics. Strategies for recruitment (e.g., P8, P18, and P20) are likely to impact the participant selection stage.

5. *Instrument design:* A questionnaire or similar instrument is designed to gather data from the sample representative of the target population. Depending on the choice of distribution, the instruments can be P43) self-administered, e.g., online forms (P51), or P50) interview-administered e.g. interview or phone survey. They can be P29) prototyped, P30) implemented from the sketch, or acquired through P38) commercial tools or P31) reuse.

   Several recommendations to design and present an instrument are provided in the literature, e.g., avoid P33) intrusive and unethical questions, and P57) to lead the respondents; provide P27) a progress indicator, P35) questionnaire navigation, P40) instructions of use, P41) option to resume answering, and mainly P32) ask simple, unambiguous, actual and targeted questions. Responses can assume P46) open-ended or P49) close-ended formats.

6. *Instrument validation:* After design, the ability of the instrument to measure what is intended should be assessed. The most frequently cited approaches for the assessment are P66) piloting, P65) retest, P62) focus groups, and P63) expert or P58) non-expert reviews. Additionally, user-related metrics (e.g., usability, readability, time to respond) can result in improvements to the instrument design (P61). Ancillary documents supporting the recruitment stage should also be

reviewed, e.g., cover letter (P60) and thank you letter (P56), likely providing incentives to the respondents (P67).

7. *Participant recruitment:* The strategies to select potential participants are previously defined in the sampling plan stage, such as P24) invitations and authorization codes, P26) self-recruitment, and P13) snowballing. By adopting proper actions and technology support, researchers can even investigate the potential threats to the process related to drop-outs (P68).

8. *Response management:* After distribution of the instrument to the selected participants, it is important to observe the response rate (P64) in order to identify the reasons for non-responses (P25). To ensure that the expected number of responses is achieved, researchers are likely to send reminders (P70) or to provide rewards for participation (P67).

9. *Data analysis:* Prior to the synthesis, the collected data should be validated (P78) in order to handle incomplete and missing values (P74). Furthermore, qualitative (e.g., P73) or quantitative (e.g., P76) analysis methods can be employed according to the survey's sample and response format. The results should then be P80) presented, P81) interpreted and likely P77) compared to particular subsets of the population. An additional suggestion to ensure their reliability is P79) to have more than one item measuring the same variable.

10. *Reporting:* The main practice related to the reporting phase is to produce an output of the information contained in the process documentation (P6). Ideally, the documentation is to be updated during the survey process, including the data analysis and results' interpretation. Both the related work (P3) and the adopted guidelines (P4) are used as additional information sources for this stage. Finally, it is important to consider the report's intended audience (P7).

The frequency with which the practices occur in the segments of the text does not denote its importance for the process. The reasons to adopt a particular practice over another depends on the researchers' conscious decision supported by the guidelines employed.

## C3. Rationales and Outcomes

We define a rationale as the motivation to choose a particular practice. They are often described as desirable process attributes (e.g., cost-effectiveness, generalizability) or outcomes of such actions (e.g., minimize or introduce bias). As an example, to achieve generalizability, a researcher should utilize probabilistic sampling (P8) and estimates of the population size (P23).

Moreover, several of the rationales can be related to validity threats, i.e., concerns about the methodology in order to achieve valid conclusions [38]. The additional rationales are not fully linked to validity, but can still positively or negatively affect it (e.g., willingness is likely to influence data quality, and thus validity). The coding resulted in a set of nine rationale categories, as shown in Figure 4.

R1 *Conclusion validity:* the actual extent to which conclusions about the investigated relationship are true or correct. Survey-based research employing quantitative analysis methods are prone to significance, effect size, and magnitude factors. Moreover, the reliability or confidence of the results is inherent to the conclusion validity.
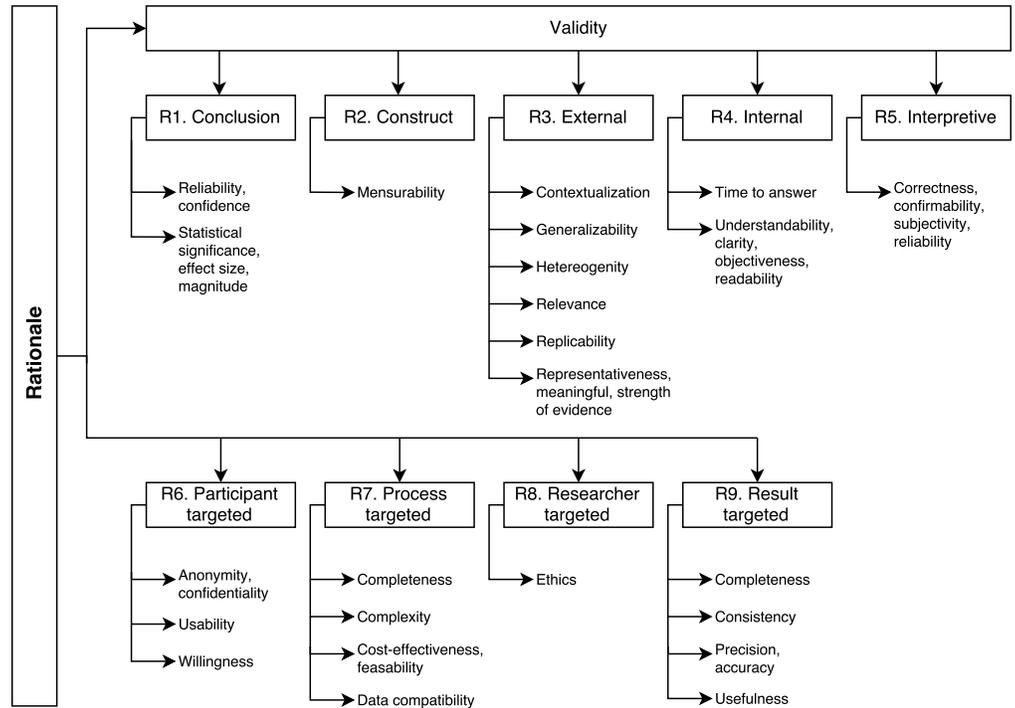
**Figure 4.** Rationale Aggregation Tree. Boxes represent the major categories and the related rationales are listed below. The topmost five categories are related to the validity of the research, whereas the bottom four are combined according to targeted aspects, i.e., participants, process, researcher and results.

R2 *Construct validity:* refers to the interaction between the underlying theory and measurement constructs, i.e., if the variables are actually measuring what they mean to. The main rationale in this category is mensurability, mainly addressed by the instrument validation stage.

R3 *External validity:* the degree to which the results of the survey can be applicable to other scenarios, such as different contexts and strata of the population. Surveys are largely impacted by external validity factors, such as generalizability, replicability, and relevance to practice. Some major factors in this category are whether the sample is representative and heterogeneous to the overall target population.

R4 *Internal validity:* represents an estimate of the degree to which conclusions about the investigated relationships can be drawn based on the measures and the research process. In survey-based research, the rationales in this category are mostly related to the sampling and instrumentation stages, e.g., understandability and time to respond.

R5 *Interpretive validity:* related to the inference of the participants' opinions from the collected responses. Unlike the conclusion validity, the interpretive is more focused on the analysis of the qualitative data. Thus, factors such as correctness, confirmability, and subjectiveness play an important role in interpreting the data.

R6 *Participant-targeted:* additional factors related to the respondents include concerns about anonymity and confidentiality, usability and willingness. Those

factors are likely to impact the data quality, as they can positively or negatively influence the participants while answering the survey.

R7 *Process-targeted:* the improvement of the process itself is a target of several rationales. Researchers carrying out surveys should pay special attention to cost-effectiveness, as their decisions are likely to require extra resources. Other practical considerations include the complexity of the instruments and techniques, completeness of the sampling sources, and compatibility of produced data.

R8 *Researcher-targeted:* by providing their opinions, survey respondents trust that the gathered data will be processed responsibly. Thus, it is essential that the researchers are aware of potential ethical issues and their responsibilities regarding the survey process.

R9 *Result-targeted:* one would expect a properly conducted survey process to produce useful results. Data validation tasks are meant to assess data consistency and completeness. Moreover, the precision of the results can be achieved by properly addressing a representative sample of the target population.

Rationales are potential standards for quality of survey research. Moreover, by relating them to the practices can potentially support design decisions, e.g., the open-ended question format has implications on the interpretive validity from a researcher's perspective (R5); whilst standardized close-ended are leaning to conclusion validity (R1). Therefore, it is important to reflect upon the importance of different rationales while conducting the research process.

### C4. Checklist Instrument

The checklist had originally 38 items, distributed according to the survey process stages (see pre-evaluation checklist in Appendix A.1). At the end of each question, we also identify the related practices and rationales, which are intended to support reviewers while using the checklist. The rightmost column identifies the key decision points and related conditions that could impact the research (see Section 2.4, e.g., D1:P means kind of sample: probabilistic as D1:NP is related to the non-probabilistic alternative).

Later, we complemented the reporting section of the checklist with more generic questions selected from the checklist proposed by Dybå and Dingsøyr [16]. Those questions are meant to assess the quality of the evidence produced by empirical studies, regardless of the research method employed.

One can notice that several practices are presented within a stage not in accordance to the co-occurrence table (Section 2.4). Although some of the practices require early planning, they could eventually be carried out in a later stage. Therefore, the checklist is organized in accordance with the stages in which those actions are more likely to take place.

Several process stages and their recommended practices are subject to decision-making (e.g., the kind of sample, instrument type). The choice should be guided by the motivations (i.e., rationales) and desired outcomes of the process. Due to this decision-making aspect, not all checklist items can be achieved to the same degree at once. We, therefore, rely on the researchers to prioritize the checklist items and hence made trade-offs according to their research goals.

## 3 Step 2. Evaluation of the checklist

In the second step of our work, we conducted a evaluation of the checklist in a research practice context, i.e., with researchers that published survey research papers. This

evaluation intended to assess the appropriateness of the checklist, i.e., how well it addresses the needs of the research community to assess survey-based research. In particular, we were interested in evaluating the completeness, relevance and fairness of the checklist. To address these goals, we formulated three evaluation questions:

EQ1 **Completeness:** Is the checklist missing any important aspect?

EQ2 **Relevance:** Does the checklist contains items that are not relevant for SE research? and

EQ3 **Fairness:** Is the assessment using the checklist too lenient or stringent?

## 3.1 Method

The goal of our evaluation was to verify whether researchers agreed with our independent assessment of their work, and furthermore whether they though the checklist was complete and fair to assess survey research reports. To address such goal, we employed a mixed quantitative-qualitative approach, using our assessment through the checklist as the subject of study. Similar strategies for evaluating methods and tools in the SE context are described in [25].

A set of conditions favor the decision to adopt such a evaluation approach within the research community, such as:

- The subject of study is widely available, i.e., survey-based articles published in ESE journals and conferences. We can directly measure the compliance to the survey practices by assessing these objects with the support of our checklist.

- Researchers experienced with conducting surveys (i.e., corresponding authors of the above-mentioned papers) are likely to be interested in the checklist's potential use. Their expert judgment is also essential to identifying limitations and to gather improvement suggestions.

- Some of the benefits from using the checklist to assess research (i.e., rigour and fairness) are difficult to quantify. Thus, open-ended questions are more likely to provide a deep understanding of the opinions and reasoning of the experts regarding such aspects.

Our research evaluation process is based on a series of steps derived from a literature search [29] and survey-based recruitment and data collection [24]. Later, we analyzed the data according to both quantitative and qualitative synthesis procedures [9].

### 3.1.1 Selection and recruitment

**Search strategy.** At first, we identified survey-based articles that can be assessed using our checklist. We searched for potential candidates in nine venues (four journals and five conferences) publishing empirical research studies in SE, namely:

**TSE** Transactions on Software Engineering

**IST** Information and Sofware Technology

**ESEJ** Empirical Software Engineering Journal

**JSERD** Journal of Software Engineering Research and Development

**ICSE** International Conference on Software Engineering

**SEAA** Euromicro Conference on Software Engineering and Advanced Applications

**IWSM-Mensura** International Workshop on Software Measurement

**EASE** International Conference on Evaluation and Assessment in Software Engineering

**ESEM** International Symposium on Empirical Software Engineering and Measurement

Despite well-established guidelines [24, 30], our checklist also incorporate practices mentioned in recent guidelines (e.g. [15, 15, 32]). Thus, we opted for candidate papers that were published in the 5 most recent years (i.e., from 2012 to 2017), as they are more likely to incorporate such practices. From this database, we identified 3429 potential publications matching these characteristics.

**Selection process.** We further filtered the papers that mentioned the term "survey" in the title or abstract, thus narrowing the original list down to 177 candidates. We gathered these papers and selected them according to an inclusion criterion: *Does the paper clearly reports survey-based research?* This resulted in 62 included papers.

**Recruitment.** Later, we invited by e-mail the corresponding authors of the selected papers to participate in our evaluation. Two of the corresponding authors have more than one paper in our candidate list, thus we sent 60 invitations related to 62 resulting papers. The invitation letter presented the goal and the context of the research and also described the assessment procedure (see evaluation procedure, below).

**Responses.** Three invitation e-mails could not be received with the given e-mail address. Out of the 57 authors who received an e-mail, 22 agreed to participate. One of them consented in assessing two of the papers we asked for and also provided an extra paper which was not part of our dataset. The additional paper was selected and aggregated to our list.

### 3.1.2 Evaluation process

The process to evaluate our checklist consisted of:

1. collecting the referred paper and applying the checklist to assessing it;

2. providing the corresponding authors with the filled out checklist so that they could review our assessment; and finally

3. asking the corresponding authors to provide feedback regarding the checklist instrument and the resulting assessment.

The resulting scores from our assessment using the checklist (see item 1, above) were aggregated into a dataset. This dataset was further analyzed in order to explore how many of the papers addressed each checklist item. Identifying patterns such as checklist items poorly addressed by most of the papers is essential for the next steps of our study. After receiving the participant's feedback, we compared theirs review to the patterns we identified.

Out of the 22 corresponding authors who agreed to participate, 12 provided us with a feedback, which consisted of:

- a **review of our assessment**, in which the corresponding author can point out disagreements with our assessment, and and refine the assessment scores;

- **response to three opinion questions** regarding the checklist, as follows:

  1. Do you consider the checklist complete? If not, what should be included?

  2. Is there anything you would like to remove, or do you think it is irrelevant?

  3. Do you think our assessment by means of this checklist is fair? That is, was our assessment of the paper too rigid or too lenient?

### 3.1.3 Data analysis

Our analysis considered the feedback provided by the participants in the form of (1) a review of our assessment, and (2) answers to a set of opinion questions.

In order to analyse the **review of our assessment**, we gathered the notes and comments provided by the participants regarding each of the checklist items. These notes were used to assess the *completeness and relevance* (RQ1 and RQ2, respectively) of our checklist. In particular, we look for suggestions to improve the checklist, whether by removing, adding, or rephrasing. We responded to each comment and highlighted any action we took to improve the checklist based on the participants' feedback (see Appendix A.3).

Furthermore, we assessed the *fairness of our assessment* (RQ3) by computing the inter-rater agreement between the scores in our assessment and the ones reviewed by the corresponding authors. The inter-rater agreement is expressed in accordance with Cohen's kappa coefficient.

We also aggregated the participants' **answers to three opinion questions** into a common list[2]. These open-ended answers comprise the respondents own phrasing and reasoning regarding the three topics of our evaluation (i.e. *completeness, relevance, and fairness*).

We read each of the answers and assigned a value in a scale of yes/no/partial, representing their agreement with the question. We used both information types (i.e., assigned value and open-ended text) to answer our evaluation questions. Ultimately, we compared the participants' opinions with the findings from the review of our evaluation in order to identify recurrent themes.

## 3.2 Threats to validity

The checklist results from a systematic process to elicit recommended practices from survey guidelines in SE research. This process is not biased-free, however. In order to assess how well the checklist addresses the needs of the research community, we ought to evaluate its completeness and relevance, and fairness of the produced assessment.

**Construct validity.** A major threat to validity concerns the ability to assess the constructs with qualitative questions. We asked the participants to provide their own opinions regarding the checklist and our assessment. In particular, one of the participants questioned whether completeness could be assessed based on opinions.

Conventionally, open questions are associated to subjective responses, which is likely to constrain the analysis of data. To reduce such effect, we assigned an agreement value (using a 3-point scale) to the participant's opinion. Besides the participants' opinions, we support our findings with the scores resulting from our assessment. These scores are used to identify the practices often not reported or addressed.

**Interpretative validity.** Another potential threat to validity relates to the interpretation of the findings. In particular, we formulated three open-ended questions to collect participants' opinions. The questions themselves are not bias-free, as they are formulated to extract a positive/negative response. As an example, "Do you consider the checklist complete?" received more positive than negative answers. To decrease this threat, the data analysis and interpretation of our evaluation study were conducted by the first author and discussed with the other co-authors.

**Reliability.** Our great involvement in constructing the checklist is likely to introduce personal biases on our assessment scores. We aimed to mitigate these by building a traceable chain of evidence. First, we assessed the selected papers and recorded notes to support the given scores. We later asked the corresponding authors to review our scores and notes, and to refine any disagreement they identified. We further

---

[2]Available at `https://goo.gl/XE7wQF`.

computed the inter-rater agreement between ours and the participants' scores, resulting in a very strong agreement (k = 0.91, according to weighted Cohen's Kappa [6]).

**External validity.** Our selection process aimed to identify a diverse set of survey-based articles, i.e. surveys in different areas and/or surveys of different quality. The sample of papers collected covers a wide range of SE topics, e.g., testing, modeling, and industry practice. These papers were peer-reviewed, so we assume they present a rigorous and sound description from the survey process. This assumption is supported by the results of our assessment, in which the selected papers comply with 65% of the items in our checklist. Thus, our sample is not diverse with regard to the methodological quality of the papers. Besides that, the participation of experienced researchers supports the generalization of our findings by expertise.

## 3.3 Results

### 3.3.1 Our assessment using the checklist

We applied our checklist to assess 24 papers reporting survey research. Each of the checklist items was ranked as fully addressed (F), partially addressed (P), not addressed (N), or not applicable (NA). A summary of our assessment is presented in Table 1.

**Table 1.** Summary of the combined scores obtained by the papers in our sample. Each row represents a checklist item, and the relative amount of papers (out of 24) ranked as fully addressed (F), partially addressed (P), not addressed (N), or not applicable (NA). The last column computes a compliance score based on how many papers address the related item.

| # | | N | P | F | NA | Compl. |
|---|---|---|---|---|---|---|
| **1. Research Objectives** | | | | | | |
| 1A | Are the research question(s)... | 1 | 0 | 23 | 0 | 95.8% |
| 1B | Is the research context defined?... | 1 | 0 | 23 | 0 | 95.8% |
| 1C | Are the needs for the survey... | 1 | 0 | 23 | 0 | 95.8% |
| **2. Study plan** | | | | | | |
| 2A | Is the survey process supported by guidelines?... | 13 | 2 | 9 | 0 | 41.7% |
| 2B | Is there a reflection on the need to update the research plan?... | 19 | 0 | 5 | 0 | 20.8% |
| 2C | Are the roles and responsibilities... | 20 | 2 | 2 | 0 | 12.5% |
| **3. Identify population** | | | | | | |
| 3A | Is the population characterized...? | 13 | 0 | 11 | 0 | 45.8% |
| 3B | Is the size of the population... | 19 | 1 | 4 | 0 | 18.7% |
| **4. Sampling plan** | | | | | | |
| 4A | Is the kind of sample... defined? | 8 | 5 | 11 | 0 | 56.2% |
| 4B | Is the sample size calculated... | 7 | 1 | 16 | 0 | 68.7% |
| 4C | Are the sources of sampling... | 1 | 2 | 21 | 0 | 91.7% |
| 4D | Are the strategies and criteria to select units... | 10 | 0 | 14 | 0 | 58.3% |
| **5. Instrument design** | | | | | | |
| 5A | Is the type of instrument... defined? | 1 | 1 | 22 | 0 | 93.7% |
| 5B | Is the instrument design process... | 5 | 2 | 17 | 0 | 75% |
| 5C | Are the demographic questions ... | 2 | 2 | 20 | 0 | 87.5% |
| 5D | Does particular care is taken to make the questions understandable...? | 10 | 2 | 12 | 0 | 54.2% |
| 5E | Is the number and order of the questions taken in consideration? | 16 | 1 | 7 | 0 | 31.2% |
| 5F | Is there a reflection on the type of responses... for the questions? | 4 | 1 | 19 | 0 | 81.2% |
| 5G | If employing close-ended questions, are the standardized response... | 1 | 3 | 20 | 0 | 89.5% |
| 5H. | Is there a reflection on the adoption of additional sources... | 18 | 2 | 4 | 0 | 20.8% |
| **6. Instrument validation** | | | | | | |
| 6A. | Is the validation process of the survey instrument detailed?... | 6 | 0 | 18 | 0 | 75.0% |
| 6B. | Is the instrument measuring what is intended?... | 6 | 5 | 13 | 0 | 64.6% |
| 6C. | In case of an electronic or online questionnaire, is the usability ... | 21 | 2 | 1 | 0 | 8.3% |
| 6D. | Are the results of the instrument validation discussed?... | 10 | 1 | 12 | 1 | 54.3% |
| **7. Participant recruitment** | | | | | | |
| 7A. | Are the strategies to select participants... | 0 | 1 | 23 | 0 | 97.9% |
| 7B. | Are the ancillary documents... | 13 | 4 | 7 | 0 | 37.5% |
| 7C. | If rewards or incentives to respondents are provided... | 0 | 0 | 2 | 22 | 100% |
| **8. Response management** | | | | | | |
| 8A. | Are the responses monitored?... | 4 | 2 | 18 | 0 | 79.2% |
| 8B. | Is there any action to be taken in case of non-responses...? | 16 | 0 | 5 | 3 | 23.8% |
| **9. Data analysis** | | | | | | |
| 9A. | Is the data validated... | 16 | 1 | 7 | 0 | 31.2% |
| 9B. | Is the method for data analysis... | 2 | 2 | 20 | 0 | 87.5% |
| 9C. | If statistical analysis is employed, is the hypothesis testing process... | 0 | 1 | 15 | 8 | 96.9% |
| 9D. | If using qualitative synthesis... | 0 | 0 | 10 | 14 | 100% |

*Continued on next page*

Table 1 – *Continued from previous page*

| # | | N | P | F | NA | Compl. |
|---|---|---|---|---|---|---|
| 9E. | If a stratified sample is defined... | 0 | 0 | 3 | 21 | 100% |
| **10. Reporting** | | | | | | |
| 10A. | Are the instrument and ancillary documents accessible... | 5 | 1 | 18 | 0 | 77% |
| 10B. | Has a discussion of both positive and negative findings... | 0 | 1 | 23 | 0 | 97.9% |
| 10C. | ...Are limitations of the study (e.g. threats to validity) discussed? | 0 | 3 | 21 | 0 | 93.7% |
| 10D. | Are the conclusions justified... | 0 | 1 | 23 | 0 | 97.9% |
| | **Mean** | 11.2 | 2.04 | 21.88 | 2.88 | 65% |

For each assessed item, we also added notes for possible improvements in the study's documentation, e.g., due to missing information. As an example, in relation to the checklist item 2A, which assess the detailed procedures when designing a survey, we provided the following note to one of the participants: "*The paper cited guidelines to survey research to characterize the sample and recruitment. It is not clear if the method provided in the guidelines are followed thought all the research process.*" In order to preserve the anonymity, we do not report the complete notes here. They were however shared with the corresponding authors.

The last column of the table presents a compliance score, i.e., the relative amount of papers that addresses the related checklist item. A score of 100% means that all papers were rated "F". The NA ratings are not computed, and each P counts as half of a full score. The compliance score does not take into consideration the importance of different items for the research. Thus it should be interpreted merely as an account of possible improvements to be taken into consideration.

**EQ1. Completeness:** Overall, our sample of papers complies with 65% of the items in the checklist. Some of the checklist items and groups presented better compliance, such as the items related to the research objectives (1A to 1C) and three out of four items related to reporting (10B to 10D). One expects any research work, regardless of research method employed, to meet these requirements.

**EQ2. Relevance:** Three checklist items are fully addressed by all the papers assessed. They cover practices such as incentives to responses (7C), qualitative synthesis (9D) and stratified data analysis (9E). These items are optional, and thus the assessment is rated not applicable (NA) for all the papers that do not employ such strategies. These results imply that researchers applying such strategies are likely to report them explicitly.

Among the checklist groups that are more scarcely addressed are:

2) study plan;

3) identify the population;

6) instrument validation; and

8) response management.

The low compliance scores show that the same kind of information is missing in several assessed studies. This implies that some of the recommended practices proposed by the guidelines are not followed. If we consider that this sample of papers is a good representation of the overall survey-based research in SE, the low-compliance items point out to gaps that should be part of wider discussion so to see if they are relevant in survey research.

**EQ3. Fairness:** We later compared the participants' scores to ours via inter-rater agreement. The resulting weighted Cohen's Kappa coefficient k = 0.91 [6], suggesting a

very strong level of agreement. We assume that two reviewers using the checklist independently are not likely to achieve such stronger agreement. However, the results showed that the corresponding authors judged the assessment as mostly fair. This is reinforced by the authors' answers the opinion questions (see Section 3.3.2).

### 3.3.2 Evaluation by the corresponding authors

After assessing all the papers, we provided the corresponding authors with the filled out checklist and our notes. We then asked them to provide feedback based on our assessment. Out of the 22 corresponding authors contacted, 12 replied to our request, providing feedback regarding the completeness of the checklist, irrelevant checklist items, and fairness of our assessment.

We addressed the participants' comments individually, responding to each issue in need of due attention and detailing the actions we took to improve the checklist. A subset of our responses are provided in Appendix A.3, and the complete set is available online in `https://goo.gl/YDj1XA`.

**EQ1. Completeness:** Most of the participants (7 out of 12) agreed that the checklist was complete and included all the main aspects of survey-based research. Two participants thought that the checklist was partially complete, and it could be improved by clarifying a few items.

One participant highlighted their confidence that our method of creating the checklist was grounded in methodological publications, such as [30]. This information was not provided beforehand, so we assumed that the participant is familiar with such work, thus relating our checklist items to the recommended practice described in Kitchenham's guidelines [30].

The three remaining participants who did not agree with the checklist completeness, raised issues such as:

- internal and external validity are not completely addressed in relation to the sampling plan and the instrument validation;

- more details are needed for novice researchers using the checklist; and

- validating completeness is not possible as an opinion.

These aspects are addressed individually in our feedback document (see Appendix A.3), as mentioned above.

**EQ2. Relevance:** Three participants mentioned irrelevant checklist items they believed should be removed:

2C) the checklist item addressing research roles and responsibilities was considered irrelevant for the report, but it could be part of the research plan (2B);

6A/6C) these two items should be combined, as they both address the instrument validation;

5H) using additional sources for data collection is optional, therefore if not mentioned in the paper it should be rated NA; and

7B) to provide ancillary documents (e.g., cover letter, invitation letter) is irrelevant to the research report.

The only issue raised by more than one participant is related to unifying 6A and 6C. The results of our assessment point out that most of our sample studies are in compliance with 6A (75%), but just a few (8.3%) actually address item 6C. We think that it is important to keep these two aspects separated, thus making explicit the needs

for validating the usability of the questionnaire (see e.g., recommended practices P27, P28, P35, P39, P40, P41, P45). All the issues abovementioned are discussed in our feedback document (see Appendix A.3).

**EQ3. Fairness:** Most participants (9 out of 12) considered our assessment fair. Two of them also mentioned that despite rigid, the assessment was fair. Another one highlighted the need for instruments that promote rigorous assessment of the research methods. None of the participants described our assessment as completely unfair, although three of them pointed out that items we missed in our assessment were limitations to fairness.

We noted that two participants mentioned the lack of information due to size limitations of the publication. This issue is further highlighted in the comments of other participants (see Appendix A.3). We sympathize with the participants' concern regarding a fair assessment due to the size limitation. However, we stress the importance to provide all the details needed to properly assess the research based on its report. As a recommended practice, researchers are encouraged to make additional information (e.g., research diary, questionnaire instrument, ancillary documents) accessible to the target audience.

# 4 Discussion

## 4.1 Checklist Usage

In order to assess survey-based research, reviewers can employ the proposed checklist. Prior to assessment, we suggest verifying the availability of research process information (i.e., research report, survey instrument, and ancillary documents). Thereafter, each checklist item should be carefully read and then evaluated with respect to whether the question can be answered and was reflected on in the research report.

Several checklist items comprise two or more nested questions. Those items are intertwined and should not be assessed separately. Moreover, the checklist items can be addressed as partial coverage, due to the higher level of abstraction where answer is likely to be subjective. In such cases, we rely on the reviewers' best judgment regarding the adoption of partial scores (i.e., 0.5).

It is possible to derive a scoring measure based on the checklist marks (e.g., 23 out of 38). However, we do not encourage the simple aggregation of scores in such a way, as it is likely to lead to a loss of assessment information. We suggest reviewers report the reasoning to score each question, thus highlighting the strengths and weaknesses of the assessed survey.

## 4.2 Implications to Research

The objective of our checklist is twofold: first, it is intended to audit reported survey-based research; and second, to support researchers in making research design decisions and reporting them. Ideally, both the researchers employing the checklist to plan and report their studies and the reviewers assessing the same research should obtain similar scores.

Alternatively, this reflexive checklist can be used to improve the survey process; researchers are encouraged to think and reflect upon the questions they are aiming to use. In particular, trade-offs have to be made. The completeness of the survey as well as the ability to obtain a large and representative sample are desired, but also costly. Thus, as highlighted in the survey guidelines, the research process decisions have to be reflected upon with respect to cost-effectiveness. This is not to say that researchers

should aim at minimizing the cost, but rather reflect on what is needed to fulfil the research goals.

Reviewers using our checklist are strongly encouraged to report the resulting scores along with their reflections about the checklist itself. We also foster independent evaluations to verify the appropriateness of the checklist to assess survey-based research by the research community.

Finally, the proposed checklist is intended to assess survey-based research in SE, but it has the potential to address different domains' studies (e.g., social sciences). It is important to identify the differences of the survey-based process employed in SE and in other fields, thus evaluating the checklist in a cross-domain study.

### 4.2.1 Research Practice

During the checklist evaluation, we assessed 24 papers reporting survey research. The results of our assessment (Section 3.3.1) point to a list of recommended practices (see Section 2.4) that are scarcely addressed. We believe that by communicating these insights to the community, we can encourage researchers to consider the recommended practices in their research. The scarcely addressed practices are:

2. *Study plan:* A research plan or log book (see recommend practice P6) is important to guide the research efforts. This protocol should detail the responsibilities of each stakeholder (P52) and a timetable (P5). The document should be updated as new information becomes known, and ultimately make it accessible by the end of the research.

3. *Identify the population:* Very often the demographics of the participants are described, but scarce information is provided regarding the target population. An audience analysis (P7) is likely to identify and supply these characteristics, and the census of practitioners and institutions could provide estimates of the population size (P23).

5. *Instrument design:* When designing the data collection instrument, one should carefully consider the order (P34) and amount (P36) of questions. Additional sources of data (P69), such as work repositories can provide means to cross-check the results.

6. *Instrument validation:* In general, the reports stated that some kind of validation of the questionnaire (e.g., a pilot) was performed. When employing online surveys, we should ideally check the usability of the questionnaire instrument. A series of practices (see e.g. P27, P28, P35, P39, P40, P41, and P45) can be used as a guideline or checklist to this validation. Furthermore, it is also important to report the improvements made as resulting of the instrument validation (P61).

7. *Participant recruitment:* Besides the need for making the research plan available, it is suggested to provide access to the standardized communication with participants. These documents include the invitation and cover letters (P60) as well as follow-ups and thank you letters (P56).

8. *Response management:* Besides the response rates, it is valuable to report any strategy used to improve responses, such as reminders (P70) or searching for additional databases (P16). These strategies are likely to affect the sample size, thus we should ideally discuss the implications of them for the study validity.

9. *Data analysis:* Before the analysis, researchers are encouraged to validate the data (P78) and check for inconsistent, incomplete or missing information (P74). There

are several strategies to deal with missing data, from discarding to the imputation based on statistical models. Ideally, the implications due to employing such approaches should be described.

## 4.3 Comparison with Related Work

We previously identified related studies providing checklists for assessing empirical research in SE (see Section 1.4). The existing checklists in SE target mostly experiments [23, 27] and case study research [48]. Our proposed checklist is intended to address the issues of survey-based research and its specific stages, e.g., sampling, instrument design, and recruitment.

Our resulting checklist can be comparable to Stavru's set of criteria to assess the thoroughness of surveys [45]. Similar to their work, we also systematically derived our instrument from the literature. We detailed herein the empirical processes to construct our checklist.

Our checklist differs from Stavru's criteria, as it covers a larger set of practices and provide instructions based on the literature. The recommended practices often taken into account further activities of the survey research process. As an example, in relation to the research objectives, our checklist item A1: *"Are the research objective expressed in measurable terms? E.g. as research questions, or using the goal-question-metric approach."* In this particular case, the objectives should describe what is intended to be measured, thus favoring data analysis, as proposed by [5, 24, 30, 32].

The items in our checklist are not weighted, as the relative importance of each practice depends on the survey process, and researchers are foster to reflect on the key decision points. Besides that, we evaluated our checklist with experienced researchers. The evaluation identifies room for improvement, and we provide an updated post-evaluation version in Appendix A.2.

More generic checklists [16, 47] share some similarities with our proposed checklist. Similar questions are mostly related to the research objectives and reporting stages. This is expected since some actions (e.g., formulate research questions, define the scope, discuss the results and limitations of research) are inherent to all empirical research.

We considered the question formulation of the generic checklists when phrasing our proposed instrument. Researchers who used those similar checklists are likely to recognize the overlaps. This could be beneficial, as they can employ identical reasoning when assessing the common items. Similar questions provide the opportunity to compare the checklists or the studies assessed through them.

# 5 Conclusions

In this paper, we described a process to derive an instrument to assess survey-based research and its further evaluation. The motivation for such is grounded in the increased usage of survey-based research in SE. In a previous study, we identified several guidelines supporting the survey processes, but no instrument provided a checklist to assess their quality.

A set of 14 methodological papers provided qualitative data that was collected and analyzed through thematic analysis. The resulting themes resulted in sets of practices and rationales for carrying them out. We built the proposed checklist based on those extracted themes, reporting it through 38 questions organized by the survey process stages.

Later, we employed an empirical evaluation approach to collect experts' opinions of our checklist. We provided the experts with an assessment produced by applying the checklist to their own reported surveys. Overall, our instrument was evaluated as

complete and the assessment as rigorous and fair. Issues regarding understandability and subjectivity of the checklist items were collected and based on this feedback, we update our proposed checklist.

We believe that the empirical software engineering community can benefit from our checklist for survey research. It can be a valuable asset for both researchers conducting and reporting survey-based studies, and for reviewers auditing survey reports.

As future work, we plan to investigate the potential benefits of using of the checklist by independent reviewers. We also intend to compare our checklist with other assessment tools (e.g. [23, 27, 48]) with respect to quality standards for empirical research.

# References

1. R. Akbar, M. F. Hassan, and A. Abdullah. A framework of software process tailoring for small and medium size IT companies. In *Computer & Information Science (ICCIS), 2012 International Conference on*, pages 914–918. IEEE, 2012.

2. P. L. Alreck and R. B. Settle. *The survey research handbook*. McGraw-Hill, 1994.

3. E. R. Babbie. *Survey research methods*. Wadsworth, 1973.

4. A. Cater-Steel, M. Toleman, and T. Rout. Addressing the challenges of replications of surveys in software engineering research. In *Empirical Software Engineering, 2005. 2005 International Symposium on*, pages 10–pp. IEEE, 2005.

5. M. Ciolkowski, O. Laitenberger, S. Vegas, and S. Biffl. *Practical experiences in the design and conduct of surveys in empirical software engineering*. Springer, 2003.

6. J. Cohen. Weighted kappa: Nominal scale agreement provision for scaled disagreement or partial credit. *Psychological bulletin*, 70(4):213, 1968.

7. R. B. Contreras. Examining the context in qualitative analysis: The role of the co-occurrence tool in atlas. ti. *Newsletter*, 2011:2, 2011.

8. J. M. Converse. *Survey research in the United States: Roots and emergence 1890-1960*. Routledge, 2017.

9. D. S. Cruzes and T. Dyba. Recommended steps for thematic synthesis in software engineering. In *Empirical Software Engineering and Measurement (ESEM), 2011 International Symposium on*, pages 275–284. IEEE, 2011.

10. C. W. Dawson. *Projects in computing and information systems: a student's guide*. Pearson Education, 2005.

11. R. M. de Mello, P. C. da Silva, P. Runeson, and G. H. Travassos. Towards a framework to support large scale sampling in software engineering surveys. In *Proceedings of the 8th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement*, page 48. ACM, 2014.

12. R. M. de Mello, P. C. da Silva, and G. H. Travassos. Sampling improvement in software engineering surveys. In *Proceedings of the 8th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement*, page 13. ACM, 2014.

13. R. M. de Mello, P. C. Da Silva, and G. H. Travassos. Investigating probabilistic sampling approaches for large-scale surveys in software engineering. *Journal of Software Engineering Research and Development*, 3(1):1–26, 2015.

14. R. M. de Mello and G. H. Travassos. Would sociable software engineers observe better? In *Empirical Software Engineering and Measurement, 2013 ACM/IEEE International Symposium on*, pages 279–282. IEEE, 2013.

15. R. M. de Mello and G. H. Travassos. Characterizing sampling frames in software engineering surveys. In *Proc. 12th Workshop on Experimental Software Engineering (ESELAW)*, 2015.

16. T. Dybå and T. Dingsøyr. Empirical studies of agile software development: A systematic review. *Information and software technology*, 50(9):833–859, 2008.

17. T. Dybå and T. Dingsøyr. Strength of evidence in systematic reviews in software engineering. In *Proceedings of the Second International Symposium on Empirical Software Engineering and Measurement, ESEM 2008, October 9-10, 2008, Kaiserslautern, Germany*, pages 178–187, 2008.

18. A. Fink. *The survey handbook*, volume 1. Sage, 2003.

19. F. J. Fowler Jr. *Survey research methods*. Sage publications, 2013.

20. S. Friese. Atlas. ti 7 user manual. *Berlin: ATLAS. ti Scientific Software Development GmbH*, 2012.

21. M. Galster and D. Tofan. Exploring web advertising to attract industry professionals for software engineering surveys. In *Proceedings of the 2nd International Workshop on Conducting Empirical Studies in Industry*, pages 5–8. ACM, 2014.

22. M. Höst and P. Runeson. Checklists for software engineering case study research. In *Proceedings of the First International Symposium on Empirical Software Engineering and Measurement, ESEM 2007, September 20-21, 2007, Madrid, Spain*, pages 479–481, 2007.

23. A. Jedlitschka and D. Pfahl. Reporting guidelines for controlled experiments in software engineering. In *2005 International Symposium on Empirical Software Engineering (ISESE 2005), 17-18 November 2005, Noosa Heads, Australia*, pages 95–104, 2005.

24. M. Kasunic. Designing an effective survey. Technical report, CMU/SEI-2005-HB-004. Software Engineering Institute, Carnegie Mellon University., 2005.

25. B. Kitchenham, S. Linkman, and D. Law. Desmet: a methodology for evaluating software engineering methods and tools. *Computing & Control Engineering Journal*, 8(3):120–126, 1997.

26. B. Kitchenham, L. Pickard, and S. L. Pfleeger. Case studies for method and tool evaluation. *IEEE software*, 12(4):52–62, 1995.

27. B. Kitchenham, D. I. K. Sjøberg, P. Brereton, D. Budgen, T. Dybå, M. Höst, D. Pfahl, and P. Runeson. Can we evaluate the quality of software engineering experiments? In *Proceedings of the International Symposium on Empirical Software Engineering and Measurement, ESEM 2010, 16-17 September 2010, Bolzano/Bozen, Italy*, 2010.

28. B. A. Kitchenham, P. Brereton, M. Turner, M. Niazi, S. G. Linkman, R. Pretorius, and D. Budgen. Refining the systematic literature review process - two participant-observer case studies. *Empirical Software Engineering*, 15(6):618–653, 2010.

29. B. A. Kitchenham, D. Budgen, and P. Brereton. *Evidence-Based Software Engineering and Systematic Reviews*, volume 4. CRC Press, 2015.

30. B. A. Kitchenham and S. L. Pfleeger. Principles of survey research: parts 1 – 6. *ACM SIGSOFT Software Engineering Notes*, 26–28, 2001–2003.

31. P. J. Lavrakas. *Encyclopedia of survey research methods*. Sage Publications, 2008.

32. J. Linåker, S. M. Sulaman, R. Maiani de Mello, and M. Höst. Guidelines for conducting surveys in software engineering. 2015.

33. M. K. Malhotra and V. Grover. An assessment of survey research in pom: from constructs to theory. *Journal of operations management*, 16(4):407–425, 1998.

34. J. S. Molléri, K. Petersen, and E. Mendes. Survey guidelines in software engineering: An annotated review. In *Proceedings of the 10th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement, ESEM 2016, Ciudad Real, Spain, September 8-9, 2016*, pages 58:1–58:6, 2016.

35. J. S. Molléri, K. Petersen, and E. Mendes. Cerse-catalog for empirical research in software engineering: A systematic mapping study. *Information and Software Technology*, 2018.

36. D. E. Perry, A. A. Porter, and L. G. Votta. Empirical studies of software engineering: a roadmap. In *Proceedings of the conference on The future of Software engineering*, pages 345–355. ACM, 2000.

37. D. E. Perry, S. E. Sim, and S. M. Easterbrook. Case studies for software engineers. In *Software Engineering, 2004. ICSE 2004. Proceedings. 26th International Conference on*, pages 736–738. IEEE, 2004.

38. K. Petersen and C. Gencel. Worldviews, research methods, and their relationship to validity in empirical software engineering research. In *Proceedings of the 2013 Joint Conference of the 23Nd International Workshop on Software Measurement (IWSM) and the 8th International Conference on Software Process and Product Measurement*, IWSM-MENSURA '13, pages 81–89, Washington, DC, USA, 2013. IEEE Computer Society.

39. S. L. Pfleeger. Experimental design and analysis in software engineering. *Annals of Software Engineering*, 1(1):219–253, 1995.

40. A. Pinsonneault and K. Kraemer. Survey research methodology in management information systems: an assessment. *Journal of management information systems*, 10(2):75–105, 1993.

41. T. Punter, M. Ciolkowski, B. Freimut, and I. John. Conducting on-line surveys in software engineering. In *Empirical Software Engineering, 2003. ISESE 2003. Proceedings. 2003 International Symposium on*, pages 80–88. IEEE, 2003.

42. K. F. Schulz, D. G. Altman, and D. Moher. Consort 2010 statement: updated guidelines for reporting parallel group randomised trials. *BMC medicine*, 8(1):18, 2010.

43. A. K. Shenton. Strategies for ensuring trustworthiness in qualitative research projects. *Education for information*, 22(2):63–75, 2004.

44. D. I. Sjoberg, T. Dyba, and M. Jorgensen. The future of empirical methods in software engineering research. In *2007 Future of Software Engineering*, pages 358–378. IEEE Computer Society, 2007.

45. S. Stavru. A critical examination of recent industrial surveys on agile method usage. *Journal of Systems and Software*, 94:87–97, 2014.

46. M. Torchiano and F. Ricca. Six reasons for rejecting an industrial survey paper. In *Conducting Empirical Studies in Industry (CESI), 2013 1st International Workshop on*, pages 21–26. IEEE, 2013.

47. R. Wieringa. Towards a unified checklist for empirical research in software engineering: first proposal. In *16th International Conference on Evaluation & Assessment in Software Engineering, EASE 2012, Ciudad Real, Spain, May 14-15, 2012. Proceedings*, pages 161–165, 2012.

48. R. Wieringa, N. Condori-Fernández, M. Daneva, B. Mutschler, and O. Pastor. Lessons learned from evaluating a checklist for reporting experimental and observational research. In *2012 ACM-IEEE International Symposium on Empirical Software Engineering and Measurement, ESEM '12, Lund, Sweden - September 19 - 20, 2012*, pages 157–160, 2012.

49. C. Wohlin. Is there a future for empirical software engineering? In *Proceedings of the 10th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement*, page 1. ACM, 2016.

50. C. Wohlin, P. Runeson, M. Höst, M. C. Ohlsson, and B. Regnell. *Experimentation in Software Engineering.* Springer, 2012.

51. C. Yang, P. Liang, and P. Avgeriou. A survey on software architectural assumptions. *Journal of Systems and Software*, 113:362–380, 2016.

# A    Appendices

## A.1    Survey assessment checklist (pre-evaluation)

**1. Research objectives**

1A  Are the research question(s) specified?

1B  Is the context of research defined? Does it consider a reasonable set of objectives? I.e. too many objectives require that particular considerations to the size and complexity of the questionnaire instrument are discussed. [P1, R7]

1C  Are the needs for the survey motivated? E.g. grounded on background and related studies. [P3, R5]

**2. Study plan**

2A  Is the survey process supported by guidelines? Does the researchers describe how the guidelines has been implemented? [P4, R7]

2B  Is there a reflection on the need to update the research plan? E.g. through keeping a research diary or log book. [P6, R3]

2C  Are the roles and responsibilities of researchers and other stakeholders defined? E.g. through creating a schedule or timetable. [P5, P52, R7]

**3. Identify population**

3A  Is the population clearly characterized (e.g. through audience analysis)? [P7, R3]

3B  Is the size of the population stated? If it is not possible to gather this data, are statistic estimates of the population drawn? [P23, R1, R9]

**4. Sampling plan**

4A  Is the kind of sample (i.e. probabilistic, non-probabilistic) defined? [P8, P20, R3, R7]  — D1:P / D1:NP

4B  Is the sample size calculated and presented? Are the actions needed to obtain the sample described? [P11, P19, R1, R9]

4C  Are the sources of sampling (e.g. particular databases or directories, open or restricted) defined? E.g. through a search plan. [P9, P12, P16, P17, P21, R3, R7]

4D  Are the strategies and criteria to select units (of observation, of analysis and search unit) stated? E.g. through a sampling frame. [P13, P14, P15, P18, P21, R7, R9]  — D2:SS / D2:PS

**5. Instrument design**

5A  Is the type of instrument (i.e. self- or interviewer-administrated) defined? [P43, P50, R7]  — D3:SA / D3:IA

5B  Is the instrument design process (acquisition, development, prototyping, versioning, reuse) described in the report? [P29, P30, P31, P38, P55, R4, R7]

5C  Are the demographic questions formulated according to the audience? If a stratification of the sample is planned, are the demographics adequate to characterize subsets the participants? [P22, P42, R3, R7]

5D  Does particular care is taken to make the questions understandable and to ensure that the participant can provide an unbiased answer? [P32, P33, P37, P57, R4, R6, R8]

5E  Is the number and order of the questions taken in consideration? [P34, P36, R4]

5F  Is there a reflection on the type of responses (i.e. open-ended, close-ended or a mix of both) required for the questions? [P46, P49, R2, R4]  — D4:OE / D4:CE

5G  If employing close-ended questions, are the standardized response formats (i.e. nominal, ordinal, interval or ratio) stated? [P44, P49, R1]  — D4:CE

5H  Is there a reflection on the adoption of additional sources for data collection? E.g. through the participant's profile or supporting literature. [P69, R7]

**6. Instrument validation**

6A  Is the validation process of the survey instrument detailed? E.g. through piloting, pre-test, retest, focus groups, experiments, expert or non expert reviews. [P53, P54, P58, P62, P63, P65, P66, P51, R2, R4]

⋮   ⋮

⋮   ⋮

6B  Is the instrument measuring what is intended? Are the questionnaire items mapped to the research question(s)? [P44, R2]

6C  In case of an electronic or online questionnaire, is the usability evaluated? E.g. questionnaire navigation, instructions of use, option to resume answering, progress indicator, required/non- inputs, aesthetics and layout. [P27, P28, P35, P39, P40, P41, P45, P51, R4, R6]  — D3:SA

6D  Are the results of the instrument validation discussed? After main problems been identified, were the instrument updated/amended according to the validation results? [P61, R4]

**7. Participant recruitment**

7A  Are the strategies to select participants (stage 4. Sampling plan) implemented? E.g. through invitations, authorization codes, self-recruitment, or snowballing [P13, P24, P26, R3, R6]  — D1:P / D2:SS / D2:PS

7B  Are the ancillary documents (e.g. invitation, cover and thank you letter) provided? If they were not produced, are the reasons for that discussed and convincing? [P56, P60, R6]  — D3:SA / D3:IA

7C  If rewards or incentives to respondents are provided, are the reasons and implications (e.g. ethical concerns, biases) discussed? [P67, R6, R8]

**8. Response management**

8A  Are the responses monitored? E.g. response rate, non-responsiveness and drop-out questions [P25, P64, P68, P71, R1, R4]  — D2:SS

8B  Is there any action to be taken in case of non-responses (e.g. reminders)? [P70, R6]  — D2:PS

**9. Data analysis**

9A  Is the data validated prior to analysis? E.g. through checking inconsistent, incomplete and missing values [P74, P78, R1, R5, R9]

9B  Is the method for data analysis specified? Are the steps of the analysis process described? Are they appropriate for the response formats collected? [P46, P49, R1, R5]

9C  If statistical analysis is employed, is the hypothesis testing process clearly documented? Are the standardized responses clearly presented? E.g. through tables, graphs, charts and plots [P49, P49, P72, P76, P80, R1]  — D4:CE

9D  If using qualitative synthesis (e.g. meta-ethnography, thematic or content analysis), is it clear how the categories/themes were derived from the data? [P46, P73, R5]  — D4:OE

9E  If a stratified sample is defined (see 5C), are the data analysed according to demographics? Are there meaningful comparisons drawn from them? [P22, P77, R2, R3]  — D1:P / D1:NP

**10. Reporting**

10A  Are the instrument and ancillary documents accessible (e.g. url link, external reference, appendix) to readers? If not, are the reasons for that discussed and convincing? If data resulting from the survey were disclosure, were anonymity and confidentiality of data discussed? [P56, P60, R4, R7]

10B  Has a discussion of both positive and negative findings been demonstrated? Are the discussion addressing the research question(s) or hypothesis? Does the discussion take in consideration the generalization of the findings? [P81, R1, R3]

10C  Are the results of the assessment checklist reported? Are limitations of the study (e.g. threats to validity) discussed? [R9]

10D  Are the conclusions justified by the results? Furthermore, are the implications and potential use of the results discussed? [R1]

⋮   ⋮

**Figure 5.** Survey assessment checklist proposed. This pre-evaluated version is later improved and updated (see Appendix A.2).

## A.2 Survey assessment checklist (post-evaluation)

**1. Research objectives**

1A Are the research objective expressed in measurable terms? E.g. as research questions, or using the goal-question-metric approach.

1B Is the research context defined? Does it consider a reasonable set of objectives? Obs. too many objectives requires that particular aspects relating to a questionnaire's size and complexity be discussed.

1C Is the need for a survey research motivated (i.e. grounded on background and related studies)?

**2. Study plan**

2A Is the survey process conducted based upon detailed procedures? Ideally, the survey process should also be based upon methodological guidelines.

2B Is there a reflection on the need to update the research plan? E.g. through keeping a research diary or log book.

2C Are the roles and responsibilities of researchers and other stakeholders defined? This information can be detailed in the research plan.

**3. Identify population**

3A Is the population or the survey's target audience characterized (e.g. through audience analysis)?

3B Is the size of the population stated? If it is not possible to gather this data, are statistic estimates of the population drawn?

**4. Sampling plan**

4A Is the kind of sample (i.e. probabilistic, non-probabilistic) defined? Obs. impact for data analysis, its representativeness and/or generalization should be discussed.

4B Is the sampling process described, and the resulting sample size presented?

4C Are the sources of sampling (e.g. particular databases or directories, open or restricted) defined? E.g. through a search plan.

4D Are the strategies and criteria to select units (of observation, of analysis and search unit) stated? E.g. through a sampling frame.

**5. Instrument design**

5A Is the type of instrument (i.e. self- or interviewer-administered) defined? Obs. impact for participant recruitment and manage responses should be discussed.

5B Is the instrument design process (acquisition, development, prototyping, versioning, reuse) described in the report?

5C Are the demographic questions formulated according to the audience? If a stratification of the sample is planned, are the demographics adequate to characterize subsets the participants?

5D Has special care been taken to make the questions understandable by the respondents? E.g. through using a terminology familiar to the target population, or by providing a thesaurus.

5E Has special care been taken to avoid intrusive and unethical questions? E.g. such biases may include questions that lead the respondent to a particular answer, or to expose personal data or behavior.

5F Is the number and order of the questions taken in consideration? In case of a potential bias related to the order of questions is identified, different versions of the instrument can be distributed.

5G Is there a reflection on the type of responses (i.e. open-ended, close-ended or a mix of both) required for the questions? Ideally, it should be possible to assess the type of each question, but the report could present the overall reasoning for the choices and provide a way to access the instrument.

5H If employing close-ended questions, are the standardized response formats (i.e. nominal, ordinal, interval or ratio) stated? Appropriate scales should be attributed to the questions according to the mapped variables.

5I Is there a reflection on the adoption of additional sources for data collection? E.g. through the participant's profile or supporting literature? Such additional sources may provide means for characterizing strata of participants or for validating data through cross-verification and triangulation.

⋮ ⋮

⋮ ⋮

**6. Instrument validation**

6A Is the validation process of the survey instrument detailed? E.g. through piloting, pre-test, retest, focus groups, experiments, expert or non expert reviews.

6B Is the instrument measuring what is intended? Are the questionnaire items mapped to the research question(s)?

6C In the case of an electronic or online questionnaire, is the usability evaluated? E.g. questionnaire navigation, instructions of use, option to resume answering, progress indicator, required/non-inputs, aesthetics, and layout.

6D Are the results of the instrument validation discussed? After the main problems been identified, were the instrument updated/amended according to the validation results?

**7. Participant recruitment**

7A Are the strategies to select participants (stage 4. Sampling plan) implemented? E.g. through invitations, authorization codes, self-recruitment, or snowballing

7B Are the ancillary documents (e.g. invitation, cover and thank you letter) provided? If they were not produced, are the reasons for that discussed and convincing?

7C If rewards or incentives to respondents are provided, are the reasons and implications (e.g. ethical concerns, biases) discussed? Those actions are likely to impact the participant's willing to respond and the research's ethical concerns, thus introducing validity bias.

**8. Response management**

8A Are the responses monitored? E.g. response rate, non-responsiveness, and drop-out questions. In case of inadequate response rate, the reasons for non-responses and drop-out items were investigated?

8B Is there any action to be taken in case of non-responses (e.g. reminders)? If reminders are employed, is the process for selecting and inviting new participants described? Moreover, are the implications of reminders discussed? I.e. changes in the sample size are likely to impact the heterogeneity and generalizability of data.

**9. Data analysis**

9A Is the data validated prior to analysis? E.g. through checking inconsistent, incomplete and missing values

9B Is the method for data analysis specified? Are the steps of the analysis process described? Are they suitable for the response formats collected?

9C If statistical analysis is employed, is the hypothesis testing process documented and the standardized responses presented? E.g. through tables, graphs, charts and plots

9D If using qualitative synthesis (e.g. meta-ethnography, thematic or content analysis), is it clear how the categories/themes were derived from the data?

9E If a stratified sample is defined, are the data analysed according to demographics? Are there meaningful comparisons drawn from them?

**10. Reporting**

10A Are the instrument and ancillary documents accessible (e.g. URL link, external reference, appendix) to readers? If not, are the reasons for that discussed and convincing? If data resulting from the survey were disclosure, were anonymity and confidentiality of data discussed?

10B Has a discussion of both positive and negative findings been demonstrated? Are the discussion addressing the research question(s) or hypothesis? Does the discussion take into consideration the generalization of the findings?

10C Are the results of the assessment checklist reported? Are limitations of the study (e.g. threats to validity) discussed?

10D Are the conclusions justified by the results? Furthermore, are the implications and potential use of the results discussed?

**Figure 6.** Survey assessment checklist after evaluation (see Section 3). A digital version of the checklist is available at `https://tinyurl.com/se-survey-checklist`.

## A.3  Suggestions to improve the checklist (excerpt)

Here we present a sample of the feedback provided by the participants (i.e., corresponding authors) of our evaluation in the professional context (Section 3). The comments are listed according to the checklist item they are related to. For each comment, we present our responses and actions we took to address the mentioned topics. A complete list detailing all the comments is provided at `https://goo.gl/jNXx7U`.

1B) Two comments regarding the understandability of this checklist item:

   (a) *What type of context and limitations should be described? Many questions would benefit from having a more detailed guide along with the checklist.*; and

   (b) *The term "limitations of scope" may be misleading. Reading quickly I first thought that you referred to whether the study scope has some limitations to be able to answer RQs (related to study validity).*

   **Response:** We agree that term "limitations of scope" can leave room for interpretations. It could also be misleading, as limitations are often described as threats to the validity of a study.
   **Action:** To improve the checklist understanding, we rephrased item 1B, as follows: "Is the research context defined? Does it consider a reasonable set of research objectives? Obs. too many objectives requires that particular aspects relating to a instrument's questionnaire's size and complexity be discussed.".

2) Two comments regarding the description of this checklist item:

   (a) *The sub-questions for me do not address the main question of whether a survey is appropriate. 2A-C are more about what is reported, rather than whether survey is the right method. That for me is more about whether other approaches were considered etc. Roles and responsibilities I would generally not note in a paper.*; and

   (b) *(...) I am not convinced that the three questions that are included in this category would be enough to assess if a survey study research design is appropriate to address its research aims (as it is stated in the question). The fact that guidelines are followed, a research diary is kept and responsibilities are defined does not guarantee that the research design is appropriate to answer the research questions. (...)*

   **Response:** We agree with the authors that the description of checklist item 2 does not match what is assessed in sub-items 2A-2C. These sub-items assess whether the study plan is accessible and complete, instead of "appropriate".
   In order to assess whether the survey design is appropriate to address its research objectives, we designed a specific question (see 6B), that assess if the questionnaire items are related to the research questions described in the study plan.
   The recommended practices for the checklist item 2 are: (i) to provide a survey plan document; (ii) to report the guidelines used; and (iii) to detail the responsibilities of each researcher. These aspects are important to allow for the study to be reviewed and replicated.
   **Action:** We updated phase 2's description to match what is assessed by items 2A-2C: "Study plan: Is the survey design accessible and complete?".

2C) Two comments regarding the relevance of this checklist item:

   (a) *Is this information really relevant for the report? I think this information would fits better into a protocol or research plan than in the report itself.*; and

(b) *It sounds irrelevant (e.g., how relevant it is for assessing the survey itself to know timetables)?*

**Response:** We agree with the authors that a schedule or timetable is not relevant for assessing the quality of the survey. However, these artifacts can provide means to assess the roles and responsibilities of researchers in the survey process [24]. As suggested by one of the participants, the information regarding the roles and responsibilities of each researcher could be provided in a survey plan document. We highlight here the need to make this document accessible to reviewers [4].
**Action:** We removed the references to schedule and timetables in the item 2C, instead stating of that "This information can be detailed in the research plan (see item 2B)".

5H) Three comments regarding the understandability of this checklist item:

(a) *I didn't understand it.*;

(b) *It makes no sense for me. If man don't mention another source of information it means there is not. Why do you assume that there is another unmentioned source?*; and

(c) *I don't fully follow this question. Do you mean data triangulation so findings from the survey are triangulated with other data?*

**Response:** The authors pointed out a very important understanding issue. Our proposed checklist accounts for the need to discuss if additional sources are required, e.g., to characterize stratas of the participants, or to cross validate the data related to the investigated phenomenon from multiple sources [24, 32]. As an example, after the survey, the findings can be compared to other sources, such as personal profile information or related work.
**Action:** We added a note on 5H making explicit the reasons to adopt additional sources of data collection, as follows: "Such additional sources may provide means for characterizing strata of participants or for validating data through cross-verification and triangulation".

10A) *Is data disclosure / open data also a criterion? I think it should be as people should be pushed in the general direction of open science to foster reproducibility.*
**Response:** The author suggests that, besides the ancillary documents we already mentioned in the checklist, the resulting data from the survey is also make available. We see the value on that, but acknowledge potential implications due to anonymity and confidentiality that should be taken into consideration. Therefore, we rely on the judgment of researchers conducting the survey to provide a discussion on such aspect.
**Action:** We added a note on 10A making explicit that the data disclosure can be provided and thus should be discussed. The new item is "If data resulting from the survey are disclosure, does a consideration about the anonymity and confidentiality of data is discussed?".

10A) Two additional comments regarding applicability of this checklist item:

(a) *I can't imagine including things like invitations, thank you notes, etc. My ethics approval process requires so much documentation, I would need a separate 10 pages just for all of the information I provide to participants.* ; and

(b) *not all materials can be added to an appendix especially with paper length limitation.*

**Response:** We agree with the author that limitations due the publication size are likely to constrain the amount of information provided in a paper. As an alternative, the additional documents mentioned in 10A could compose a "survey research package", available as a web reference provided in the paper [4].