



SEXTAMT: A systematic map to navigate the wide seas of factors affecting expert judgment software estimates[☆]

Patrícia Gomes Fernandes Matsubara^{a,b,*}, Bruno Freitas Gadelha^a, Igor Steinmacher^{c,d},
Tayana Uchôa Conte^a

^a Institute of Computing of the Federal University of Amazonas (UFAM), Brazil

^b Faculty of Computing of the Federal University of Mato Grosso do Sul (UFMS), Brazil

^c Federal University of Technology - Paraná (UTFPR), Brazil

^d Northern Arizona University (NAU), United States of America

ARTICLE INFO

Article history:

Received 30 April 2021

Received in revised form 11 November 2021

Accepted 12 November 2021

Available online 29 November 2021

Keywords:

Expert judgment

Software effort estimation

ABSTRACT

Context: Software projects involve technical and managerial activities, including software estimation. Inaccurate estimates are harmful and improving estimation methods is not enough: we need to understand more of the factors that impact estimates.

Objective: Our study aims to identify the existing evidence about the factors that affect estimates in software projects when using expert judgment.

Method: We executed a Systematic Literature Mapping (SLM) based on database and snowballing searches, selecting papers by first reading their titles and abstracts and later reading the full text.

Results: Researchers investigated a wide range of different factors employing mostly laboratory research strategies and relying primarily on differences of estimates and participants' perceptions to measure the factors' effects. Resulting from our analysis, we present the SEXTAMT (Software Estimates of eXperts: A Map of influencing facTors), a map of factors affecting estimates built on three dimensions: project/iteration phase, stakeholders, and type of effect.

Conclusion: Over the years, researchers have investigated a varied set of factors. Many of them were explored in different studies, employing diverse research strategies. Such studies provide compelling evidence on the elements that influence expert judgment estimates, which can be used to assess and improve everyday estimation in the software industry.

© 2021 Elsevier Inc. All rights reserved.

1. Introduction

An estimate is a quantitative assessment of a variable's likely outcome, such as project costs, resources, effort, or duration (IEEE, 2017b). Estimating tasks and projects is a critical part of developing and maintaining software, and researchers devoted a significant amount of effort to creating and assessing software estimation methods (Jørgensen and Shepperd, 2007). One such method is expert judgment: it is the preferred estimation method in the industry (Molokken and Jørgensen, 2003; Trendowicz et al., 2011). In agile software development, Planning Poker – based on expert judgment – is the most applied method (Usman et al., 2015). Expert judgment is also on the rise as a research topic in software effort estimation (Sehra et al., 2017).

Expert judgment differs from other estimation methods because the quantification step for generating the estimate is judgmental rather than mechanical (Halkjelsvik and Jørgensen, 2012). That is, experts use their human mind as a measurement instrument (Kahneman et al., 2021a). Therefore, the processes that we use for arriving at a prediction are largely unconscious (Halkjelsvik and Jørgensen, 2018b). Discovering and understanding the factors that affect expert judgment estimates is crucial for reducing errors and improving our accuracy when using such a method, and research on these factors is also a trend (Sehra et al., 2017). In addition, research and practice in other domains where evaluations and predictions rely on expert judgment have shown that countless triggers can drive variability in judgments, leading to bias, noise – and consequently, to error, unfairness, and losses (Kahneman et al., 2021b). For instance, in the seemingly exact science of forensic fingerprint analysis, where professionals have to decide whether fingerprints collected in crime scenes match exemplar fingerprints, researchers found that examiners can be misled by contextual information, such as eyewitness

[☆] Editor: Burak Turhan.

* Corresponding author at: Institute of Computing of the Federal University of Amazonas (UFAM), Brazil.

E-mail addresses: patriciagfm@icomp.ufam.edu.br (P.G.F. Matsubara), bruno@icomp.ufam.edu.br (B.F. Gadelha), igorfs@utfpr.edu.br (I. Steinmacher), tayana@icomp.ufam.edu.br (T.U. Conte).

recognition (Kahneman et al., 2021e). This led forensic laboratories to change their practices, sequencing information to which examiners are exposed before they analyze fingerprints.

Likewise, getting a comprehensive perspective of the factors researched in software estimation so far can guide researchers willing to build on the existing body of knowledge, to propose and assess new practices that minimize error and enhance the software estimation process. In addition, it can also help practitioners willing to identify the factors relevant to their context, to identify the good practices to adopt. In this article, we provide such perspective of factors through a Systematic Literature Mapping (SLM) using the guidelines of Kitchenham et al. (2015) and Petersen et al. (2015).

We found 131 relevant articles in our SLM, reporting 235 different factors – a myriad of diverse elements that somehow influence estimation results using expert judgment. Most of them (166 factors) was reported in one article and are provided as part of our supplementary material (Matsubara et al., 2021). Still, understanding the remaining 69 factors investigated in two or more articles is challenging. Therefore, we propose an instrument for researchers and practitioners to navigate the seas of factors affecting estimates: the SEXTAMT (Software Estimates of eXperts: A Map of influencing facTors).

Typically, a sextant is an instrument to aid overseas navigation by measuring the angle between the horizon and a celestial reference object like the sun, planets, or stars. The celestial object chosen as a reference depends on the period of the day the observer will take a sight. The observer can use the sun during the day or planets and stars during dawn or night. The measured angle serves as input for calculations that allow for identifying positions with the aid of nautical charts, thus supporting navigation overseas. The time the observer took the sight is also a necessary input (Hugh, 1911).

Likewise, the SEXTAMT uses reference points in the form of dimensions, which the interested reader can use to navigate these wide seas of factors. A temporal dimension alludes to the importance of time for calculating correct positions when using the physical sextant. In the SEXTAMT, it refers to a software project or iteration phases: initiating, planning, executing, monitoring and controlling, and closing – which we borrowed from the PMBOK (Project Management Body of Knowledge) group processes (Project Management Institute, 2017a). Most of the factors we found group at the planning and the executing phases. That is understandable because estimates emerge primarily at the planning phase, and the dynamics of project execution also affect our perceptions of accuracy and error of estimates.

Instead of finding a celestial object as a reference point, we included a stakeholder dimension to the SEXTAMT. The reader can define a stakeholder of interest to investigate only the factors associated with them, either because it relates to a task that the stakeholder is responsible for or because that stakeholder directly causes the factor. In some situations, the factor impacts the stakeholder somehow. Most factors are related to the estimator role, which is natural since stakeholders playing this role are responsible for estimating. However, we found factors associated with clients and users, higher management, project managers, requirement engineers, software developers, and testers. We also discovered factors that applied to the entire software team or no specific stakeholder at all.

The SEXTAMT also has a dimension regarding the type of effect of the factors. According to the direction of the effect, we had four types: positive direction for accuracy factors, negative direction for error factors, and neutral direction for value adjusting characteristics and empirical influence factors. If the reader wants to identify only the factors that increase accuracy when present, they can navigate the accuracy factors. Additionally, we grouped the factors in categories that represent the larger oceans and some smaller seas of our map.

2. Background

In this section, we present the relevant concepts for the context of our study (Section 2.1) and the related work Section 2.2, including two previous related reviews we found.

2.1. Software estimation

Software estimates are predictions about a variable, like the software project effort, cost, or duration (McConnell, 2006c). Given the importance of software estimation for industry, one critical concern is to devise improved methods to estimate software projects. More than 60% of research papers about estimation before 2007 proposed and evaluated estimation methods (Jørgensen and Shepperd, 2007). Boehm classifies these methods as algorithmic models, expert judgment, analogy-based, Parkinson, price-to-win, top-down, and bottom-up (Boehm, 1984). Our SLM focuses on expert judgment estimation, as it is the most used method in the industry (Trendowicz et al., 2011). To delineate what we mean by expert judgment-method, we used the guideline of Halkjelsvik and Jørgensen (2012): if the quantification step of the estimation method is judgmental, then the method is categorized as judgment-based. If this step is mechanical, then the method is categorized as model-based.

Another critical concern of software project estimation is the predicted variable, either size, effort, schedule, or cost of features (McConnell, 2006a). For instance, the functions of algorithmic models use size as their input (Jørgensen, 2007b). Then, considering software size estimates and productivity assumptions, estimators can generate effort estimates. From effort estimates and the project resources, estimators can generate estimates about cost, features, and duration (in calendar days), which project stakeholders use to establish the project commitments (McConnell, 2006a).

Nevertheless, many of the relationships among these software project variables are unstable and change from one context to another (Jørgensen, 2014), hampering the creation of a universal model of estimation. This instability may also explain why complex estimation models are not necessarily more accurate than simpler ones (Jørgensen, 2014). Despite this, many of the existing estimation methods can be applied to any software project variables (McConnell, 2006b). Therefore, in our SLM, we are not excluding studies based on the project variable.

2.2. Related work

Researchers have been investigating factors affecting estimates such as the anchoring bias (Aranda and Easterbrook, 2005), the impact of the development method (Molokken-Ostfold and Jørgensen, 2005), the influence of using checklists (Usman et al., 2018b), and others. In one of the related works, Halkjelsvik and Jørgensen (2012) present a review of studies about factors affecting judgment-based predictions of performance time, integrating results from the areas of psychology, engineering, and management science. Their review later inspired writing a more recent book about time predictions in general (Halkjelsvik and Jørgensen, 2018a). Given the multidisciplinary nature of their review, they opted to term performance time predictions as an equivalent for effort estimation. The authors described (i) the characteristics of estimates presented in the primary studies (ii) the details about the processes and strategies used in estimation, and (iii) the influence of task characteristics, estimators' characteristics, and contextual factors on estimates.

Halkjelsvik and Jørgensen (2012) included in their review studies correlational, quasi-experimental, and experimental designs. They excluded studies based on questionnaires and interviews describing respondents' opinions about reasons for estimation errors and biases because the authors affirm that they

do not have a suitable method to evaluate their validity. Also, they have included gray literature, like reports and unpublished manuscripts, bringing back the practice perspective and the practitioners' voice to their results that otherwise would be lost because of the exclusion of studies based on questionnaires and interviews.

In another related work, [Basten and Sunyaev \(2014\)](#) conducted an SLM focused on factors affecting software effort estimation accuracy. The authors presented four categories of factors affecting estimates: (i) factors related to the estimation process, (ii) factors related to the estimators' characteristics, (iii) features of the project to be estimated that may affect the estimates, and (iv) factors related to the external context, more specifically associated with the client. Although [Basten and Sunyaev \(2014\)](#) published their SLM in 2014, they only included papers written up to 2010. Also, their search strategy consisted of a manual search and snowballing procedures ([Basten and Sunyaev, 2014](#)). An automatic search may provide additional papers. Diverging from [Halkjelsvik and Jørgensen \(2012\)](#) and [Basten and Sunyaev \(2014\)](#) included papers reporting opinions from software experts, as they may indicate potentially influential factors.

Thus, we foresaw a need for an update and an expansion of such reviews. We executed our SLM on the scope of software engineering, including articles up to 2020, to satisfy this. On the one hand, our SLM differentiates from the review from [Halkjelsvik and Jørgensen \(2012\)](#) by applying a systematic mapping method and focusing on the software engineering domain alone. On the other hand, our SLM differentiates from the review of [Basten and Sunyaev \(2014\)](#) by extending the timeline of included papers, focusing on expert judgment only, and by including automated search instead of manual.

3. Research method

We started the SLM by defining a systematic mapping protocol, following the guidelines presented by [Kitchenham et al. \(2015\)](#) and [Petersen et al. \(2015\)](#), and by collectively inspecting it. The remaining of this section presents our research questions. It also presents our search, selection, extraction, and analysis procedures.

3.1. Research questions

Our primary research question is: **RQ 1 – How have researchers investigated the factors that affect expert judgment software estimation?** As we want to explore different aspects of the existing evidence about the factors, we further refined our primary research question in the following set of secondary research questions:

- **SQ 1.1 – What are the factors that affect expert judgment software estimation?**
- **SQ 1.2 – How was the impact of the factors over the expert judgment estimates measured?**
- **SQ 1.3 – What are the software project estimate variables investigated?**
- **SQ 1.4 – When and where are published the studies about factors affecting expert judgment software estimates? and**
- **SQ 1.5 – What research strategies and methods are used to investigate factors that affect expert judgment software estimation?**

Table 1

Second version of the search string.

("effort estimation" OR "effort estimate" OR "cost estimation" OR "cost estimate" OR "duration estimation" OR "duration estimate" OR "schedule estimation" OR "schedule estimate" OR "size estimation" OR "size estimate") AND (factor OR reason OR cause OR "anchor" OR "impact" OR "risk identification" OR "customer collaboration") AND (software OR system)

3.2. Search and selection

We started the search process by defining a known set of papers, which we used as an oracle to validate our search string's outcomes. Our oracle had 25 papers.¹ Our next step was defining the search string. The results of automated searches are highly dependent on the search string's quality ([Petersen et al., 2015](#); [Zhang et al., 2011](#)). We defined ours based on the extraction of the keywords of the titles and abstracts from the articles in our known set of papers, as [Petersen et al. \(2015\)](#) recommend.

We executed the automated search restricting the search to title, abstract, and keywords whenever possible. Our sensitivity² goal for the automated search was 70%, as [Zhang et al. \(2011\)](#) recommended. After the first search round, we got a sensitivity of 60%—below our goal of 70%. We ran a trial search without restricting the search to title, abstract, and keywords, but the high number of results made this change prohibitive.³ We refined the search string, leading us to the second and final version, presented in [Table 1](#).

We carried out the automated search on ACM, IEEEExplore, Scopus, and El Compendex (Engineering Village), as illustrated in [Fig. 1](#) (Step 1), resulting in 5113 articles and a sensitivity of 84%, satisfying our goal of more than 70%. We did not include other publisher-specific databases, like SpringerLink and ScienceDirect, as they would probably yield a larger number of duplicates, according to [Dyba et al. \(2007\)](#).

After eliminating duplicates from the 5113 articles, we came to a total of 3654 articles ([Fig. 1](#), Step 2). Next, we executed the selection procedures, considering the following inclusion criteria: IC01 – The paper presents an empirical study that investigates factors that affect software project estimates related to expert judgment. We also selected the papers based on the exclusion criteria that we present in [Table 2](#). Additionally, [Table 2](#) presents the relationship between each exclusion criteria and the filter in which we applied it mostly: Filter 1 (title and abstract) and/or Filter 2 (full-text).

To reduce bias during the selection process, we independently selected a random sample of the articles retrieved by the search by reading their titles and abstracts. We calculated the researchers' level of inter-rater agreement on this sample of articles through the kappa coefficient ([Kitchenham et al., 2015](#)). We got a kappa level of 0.83, which is very good, according to [Kitchenham et al. \(2015\)](#). So, we considered the kappa level adequate, and we proceeded with the selection, getting to a total of 173 papers selected based on title and abstract ([Fig. 1](#), Step 3). After reading the full text of all the 173 articles, we selected 81 that satisfied the inclusion criteria and that we could not eliminate with our exclusion criteria ([Fig. 1](#), Step 4).

¹ The final list with the known set of papers is in the supplementary material, together with more details of the search and selection procedures ([Matsubara et al., 2021](#)).

² Number of relevant studies retrieved divided by the total number of relevant studies and then multiplied by 100 ([Zhang et al., 2011](#)). The number of papers in the known set is the number of relevant studies.

³ For ACM alone we had over 480,000 results.

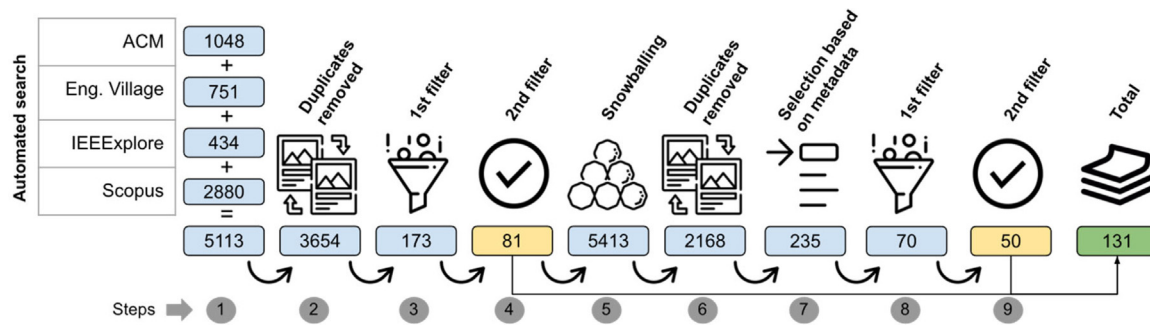


Fig. 1. Search and selection results.

Table 2

Exclusion criteria and their relationship with the selection filters.

ID	Exclusion criteria description	Filter
EC01	The paper presents a systematic mapping/review, lessons learned, or opinion paper, rather than an empirical study on factors that affect software project estimates related to expert judgment.	1, 2
EC02	The paper focus on factors affecting estimates related to estimation methods other than expert judgment.	1, 2
EC03	The paper presents non-peer-reviewed results.	1
EC04	The paper is not written in English.	1
EC05	The paper is not accessible in full-text online.	1
EC06	The study is published as a book or gray literature.	1
EC07	The paper is a duplicate or a previous version of another already selected paper.	2
EC08	The paper does not describe the factors to allow for categorization	2

The final set of papers selected from the database search formed the start-set for backward and forward snowballing (Wohlin, 2014). We aimed for a sensitivity of 100% after the snowballing step. We got to a total of 5413 articles through backward and forward snowballing (Fig. 1, Step 5), and to 2618 after removing duplicates (Fig. 1, Step 6). We selected a total of 234 of them based on their metadata – title, authors, and venue – and on their citation context on the original articles in the case of backward snowballing (Fig. 1, Step 7). We read their abstracts, reducing the number to a total of 70 articles (Fig. 1, Step 8). Following, we read their full text, leading to the inclusion of 50 articles (Fig. 1, Step 9). Therefore, the final list of articles included in our SLM contains 131 articles, and we satisfied our goal of 100% sensitivity of papers from our known set of papers.

3.3. Data extraction

We extracted the data using a form⁴ created and later refined after a pilot data extraction over the known set of papers. We extracted the following data:

- Title, authors and their affiliation, venue and year of publication;
- research strategy, according to the classifications of Stol and Fitzgerald (2018) and Storey et al. (2020), and research method;
- observations and context;
- factors and discussion about them;

- project variables that were the focus of estimation. These variables could be either size, effort, cost, productivity, or duration;
- how authors measured the impact of the factors over the estimates.

3.4. Data analysis

In Fig. 2, we provide an overview of our data analysis. After reading the full text of all selected articles and extracting text and data to our extraction form, we created codes to summarize the findings from the primary studies,⁵ supporting the aggregation of data into factors later during the analysis process.

Most of the codes we generated followed the structure we show in Fig. 2, with some variations. The **candidate factor** was the label that the original study authors provided. The **quantitative results** summarized whether the authors found significant results, sometimes informing p-values or other relevant information. It was optional, once only quantitative studies needed such data. The **brief description of effects** highlighted whether the candidate factor was a reason for accuracy, a reason for errors, an effort predictor, among others.

Next, we created mind maps aggregating similar candidate factors under a final factor label. We chose the final label to reflect the core of the candidate factors. In some situations, we had an intermediary factor label, reflecting essential variations of the core factor. We held regular meetings to review the mind maps with the categories, candidate factors, and codes. We analyzed the factors through the lenses of a few dimensions we considered relevant to interpret the results. The categories we used to organize the data relate to three dimensions, shown in Fig. 3.

The temporal dimension regards the phase of a software project/iteration that a factor is likely to happen or to cause an impact, based on the PMBOK project phases (Project Management Institute, 2017a). The stakeholder dimension informs one stakeholder or a group responsible for a task or process to which the factor is linked or that directly causes the factor. In some situations, the factor impacts the stakeholder. The type of effect dimension indicates the nature of the impact of the factor over the estimates, considering the results of the primary studies: (i) error factors are negative when present; (ii) accuracy factors lead to improvements in estimates' accuracy when present; (iii) value adjusting characteristics lead to a need for a higher or lower value of estimate and are inputs to estimation; and (iv) empirical influence indicate factors whose impact on the estimates are not definitely negative, positive, or leading to a need to a higher or lower value: it varies in direction and nature. Some of the factors under this label can lead to improvements in accuracy in

⁴ The form, as well as the complete extraction data are in the supplementary material (Matsubara et al., 2021).

⁵ All factors with their categories and codes are in the supplementary material (Matsubara et al., 2021).

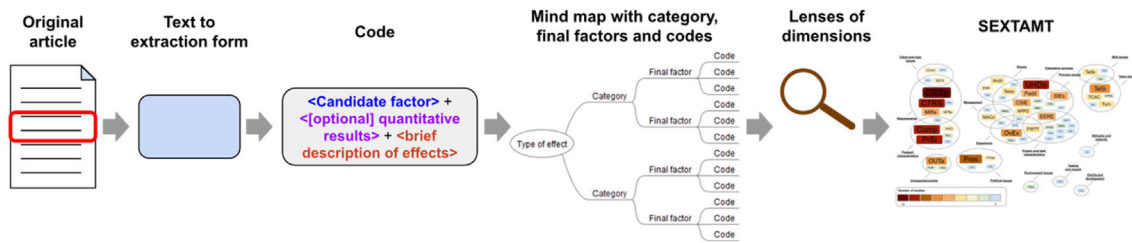


Fig. 2. Overview of the analysis.

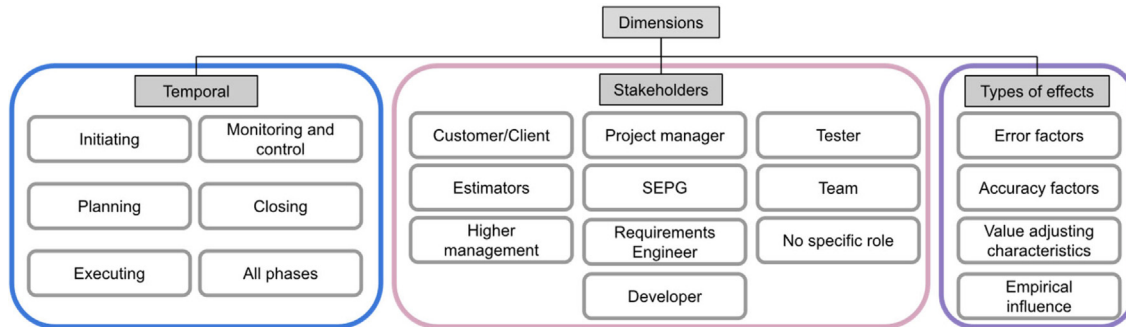


Fig. 3. SEXTAMT dimensions.

some circumstances, but to inaccuracies in others. For instance, the client’s expectation factor has an empirical influence over the estimates. If, by chance, such expectations are realistic, their impact are on the direction of making the estimate more accurate. Otherwise, they may lead to estimation error.

Finally, we created the SEXTAMT. We used the dimensions as the cornerstone for the navigation through the factors. However, we excluded from the SEXTAMT all the factors reported in only one article due to space restrictions, reporting them in our supplementary material. In the next section, we explore our results.

4. Results

In this SLM, we aim to answer the following primary research question: **RQ 1 - How have researchers investigated the factors that affect expert judgment software estimation?** In this section, we explore our results, considering each secondary research question presented in Section 3.1.

4.1. SQ 1.1 – What are the factors that affect expert judgment software estimation?

After analyzing all papers, we found 235 factors in total, from which we report the 69 that were explored in more than one research article. We present the 69 factors in Table 3, with an ID code in parenthesis, and the articles with the evidence about them.

In Section 5, we detail the factors, presenting them as part of the SEXTAMT. We also organized the factors considering the dimensions we described in Fig. 3.

4.2. SQ 1.2 – How was the impact of the factors over the expert judgment estimates measured?

This question’s motivation was to identify how researchers evaluate the impact of the factors over the estimates. Table 4 presents the associations between the strategy that researchers used for impact measurement with each article. Each article could have multiple different ways to measure impact.

Researchers’ most used strategy for investigating the impact of factors was participants’ perceptions: 45 articles adopted it, using either respondents or field research strategies. Some of these studies required participants to evaluate their companies or project accuracy subjectively. Another strategy widely used was assessing the difference of estimates between an experimental and a control group, with 44 occurrences. This is common in laboratory experiments, which was the most applied research strategy as discussed in Section 4.5. By analyzing the difference of estimates, researchers investigated the factors that could cause a shift from more realistic estimates to more optimistic ones – supposing that lower estimates lead to higher chances of error. Regarding more objective measures of accuracy, bias, and error, researchers used metrics like MRE (Magnitude of Relative Error), MREBias, BRE (Balanced Relative Error), and BREBias, as we show in Table 5.

Seven studies relied on less traditional metrics involving the estimated and actual values. While the critiques of MRE and MREBias focus on the use of actual values at the denominator of the formula – which is resolved in BRE and BREBias by using the minimum value between the estimated and actual values – seven studies use the estimated value at the denominator. We categorized these studies under the term “deviation”, since the researchers of such articles disagree about the best name for the metric, calling it effort deviation (Branco et al., 2015; Ohlsson et al., 1998), effort overrun (Jorgensen and Carelius, 2004), accuracy (Bergeron and St-Arnaud, 1992), effort variance (McGarry et al., 1998), overrun factor (Halstead et al., 2012),⁶ or project overrun (Lind and Sulek, 1998). Another three studies use the absolute error (estimated - actual value).

A total of six studies evaluates total effort. They are either based on regression analysis (He et al., 2010; Mendes et al., 2005) or correlations of effort with other variables (Arifin et al., 2017; Davis, 1989; Nugroho and Lange, 2008; Vijayakumar, 1997). Three studies relied on classifying projects according to ranges of

⁶ The original formula was actual duration = estimated value + estimated value*overrun factor for this study. Isolating the overrun factor, we get to the same formula as the other studies.

Table 3
List of factors.

Factor	Articles
Diligence (Dili)	Basten and Mellis (2011) and Lederer and Prasad (1995a)
Anchoring effect (Anch)	Shepperd et al. (2018), Aranda and Easterbrook (2005), Løhre and Jørgensen (2016), Jørgensen and Grimstad (2012) and Jørgensen and Gruschke (2009)
Effect of more and/or irrelevant information (EMII)	Jørgensen and Grimstad (2008), Usman et al. (2018a), Jørgensen and Grimstad (2011) and Grimstad and Jørgensen (2007)
Optimism (Opti)	Jørgensen et al. (2007) and Magazinius et al. (2012)
Sequence effects (Sequ)	Grimstad and Jørgensen (2009), Jørgensen (2016a), Jørgensen and Halkjelsvik (2020) and Jørgensen (2013b)
Time frame size (TFSi)	Jørgensen and Halkjelsvik (2010) and Halkjelsvik and Jørgensen (2011)
Unit effects (UnEf)	Jørgensen (2016a) and Jørgensen (2015)
Size (PrSi)	Conoscenti et al. (2019), Usman et al. (2015, 2018a), Lagerström et al. (2012), Layman et al. (2008), Silva-de-Souza and Travassos (2017), He et al. (2010), Vijayakumar (1997) and Usman et al. (2017)
Complexity (Comp)	Conoscenti et al. (2019), Usman et al. (2018b), Magazinius and Svensson (2014), Morgenshtern et al. (2007), Lee et al. (2011), Silva-de-Souza and Travassos (2017), Zapata and Chaudron (2013), Subramanian et al. (2006) and Altaieb et al. (2020a)
Integration and dependencies (InAD)	Britto et al. (2015), Magazinius and Svensson (2014), Usman et al. (2015) and Lee et al. (2011)
Platform (Plat)	Lagerström et al. (2012), Huang et al. (2015) and Altaieb et al. (2020a)
Programming language (Prog)	He et al. (2010) and Huang et al. (2015)
Collaboration and communication (CCAC)	Lederer and Prasad (1995a), Usman et al. (2017), Lederer and Mirani (1990) and Molokken-Ostfold and Furulund (2007)
Availability of knowledgeable/competent clients (AKCC)	Matos et al. (2013) and Grimstad et al. (2005)
Client's expectations (ClEx)	Jørgensen and Grimstad (2008) and Jørgensen and Sjøberg (2004)
Clarity of client's needs (ClCN)	Lederer and Prasad (1995a) and Matos et al. (2013)
Cultural differences (CuDi)	Britto et al. (2015) and Altaieb et al. (2020a)
Tool support and availability (TSAV)	Lee et al. (2011) and Jørgensen and Molokken-Ostfold (2004)
Use of historical data (UHDA)	Yang et al. (2008), Magazinic and Pernstål (2008), Lederer and Prasad (1995a), Conoscenti et al. (2019), Lee et al. (2011), Furulund and Molkken-stvold (2007), Jørgensen and Molokken-Ostfold (2004), Rahikkala et al. (2018), Shmueli et al. (2016) and Lederer and Mirani (1990)
Padding (Padd)	Magazinic and Pernstål (2008), Lederer and Prasad (1995a), Jørgensen and Molokken-Ostfold (2004), Lederer and Mirani (1990), Glass et al. (2008) and Lederer and Prasad (1991)
Anticipation of project' participants' skills (APPS)	Yang et al. (2008), Rahikkala et al. (2015a), Lederer and Prasad (1995a), Usman et al. (2018a) and Lederer and Mirani (1990)
Use of checklists (UsCh)	Usman et al. (2018b), Furulund and Molkken-stvold (2007) and Jørgensen and Molokken-Ostfold (2004)
Combination strategy of individual estimates (CSIE)	Jørgensen and Molokken (2002), Mahnič and Hovelja (2012), Gandomani et al. (2019), Haugen (2006), Molokken-Ostfold and Jørgensen (2004) and Molokken-Ostfold et al. (2008)
Involvement of technical staff (ITSt)	Yang et al. (2008), Lederer and Prasad (1995a) and Altaieb and Gravell (2019)
Informal basis for estimating (IBEs)	Usman et al. (2015), Lederer and Prasad (1998) and Keaveney and Conboy (0000)
Impact of early estimates (IEEs)	Jørgensen and Carelius (2004) and Jørgensen and Sjøberg (2001)
Reestimation and revision of estimates (REEs)	Usman et al. (2018a), Lagerström et al. (2012) and Layman et al. (2008)
Standards in estimation (StEs)	Yang et al. (2008), Magazinic and Pernstål (2008), Lederer and Prasad (1995a), Rahikkala et al. (2018), Altaieb and Gravell (2019) and Lederer and Mirani (1990)
Enough effort and resources spent on estimation (EERE)	Yang et al. (2008), Jørgensen and Carelius (2004), Lederer and Prasad (1995a), Jørgensen and Molokken-Ostfold (2004), Jørgensen and Gruschke (2009) and Rahikkala et al. (2018)
Overall experience (OvEx)	Britto et al. (2015), Usman et al. (2015), Morgenshtern et al. (2007), Magazinius et al. (2012), Matos et al. (2013), Jørgensen and Molokken-Ostfold (2004) and Karna and Gotovac (2014)
Technical experience (TeEx)	Conoscenti et al. (2019), Halstead et al. (2012) and Altaieb et al. (2020a)
Experience with similar/previous projects/tasks (ExSP)	Jørgensen and Molokken-Ostfold (2004) and Usman et al. (2017)
Familiarity with the product (FWTP)	Lee et al. (2011) and Davis (1989)
Estimation experience (EsEx)	Rahikkala et al. (2018) and Altaieb et al. (2020a)
Manager experience (MgEx)	Morgenshtern et al. (2007) and Altaieb et al. (2020b)
Monitoring and control (MACo)	Yang et al. (2008), Morgenshtern et al. (2007), Jørgensen and Molokken-Ostfold (2004), Grimstad et al. (2005) and Keaveney and Conboy (0000)
Risk assessment (RiAs)	Yang et al. (2008) and Morgenshtern et al. (2007)
Pressure (Press)	Yang et al. (2008), Lederer and Prasad (1995a), Magazinius et al. (2012), Zarour and Zein (2019), Keaveney and Conboy (0000) and Glass et al. (2008)
Price-to-win issues (PTWI)	Yang et al. (2008), Usman et al. (2015), Magazinius et al. (2012) and Jørgensen and Molokken-Ostfold (2004)

(continued on next page)

Table 3 (continued).

Factor	Articles
Goals and targets (GATa)	Magazinovic and Pernstål (2008) and Magazinius et al. (2012)
Negotiations games in estimates (NGIE)	Magazinius et al. (2012) and Glass et al. (2008)
Use of flexible/agile development model (UFAM)	Molokken-Ostvold and Jorgensen (2005), Koch and Turk (2011) and Brown et al. (2013)
Resources dependencies (ReDe)	Conoscenti et al. (2019), Usman et al. (2018a) and Layman et al. (2008)
Simplicity (Simp)	Jorgensen and Molokken-Ostvold (2004) and Jørgensen and Gruschke (2009)
Project flexibility (PrFI)	Jorgensen and Molokken-Ostvold (2004) and Grimstad et al. (2005)
Similarity with previous tasks/projects (SWPP)	Jorgensen and Molokken-Ostvold (2004) and Jørgensen and Gruschke (2009)
Task size (TaSi)	Usman et al. (2018b) and Hill et al. (2000)
Business area (BuAr)	He et al. (2010) and Huang et al. (2015)
Type of project (TyPr)	He et al. (2010) and Altaieb et al. (2020a)
Longer projects (LoPr)	Lagerström et al. (2012) and He et al. (2010)
Familiarity with the technology (FWTT)	Basten and Mellis (2011), Lee et al. (2011), Furulund and Molken-stvold (2007), Jørgensen and Gruschke (2009) and Keaveney and Conboy (0000)
Clear requirements specification (CRSp)	Yang et al. (2008), Lederer and Prasad (1995a), Conoscenti et al. (2019), Usman et al. (2015), Furulund and Molken-stvold (2007), Jorgensen and Molokken-Ostvold (2004), Grimstad et al. (2005), Altaieb et al. (2020b), Zarour and Zein (2019) and Arnuphaptrairong (2018)
Changes to requirements or scope (CTRS)	Yang et al. (2008), Usman et al. (2015, 2018a), Layman et al. (2008), Grimstad et al. (2005), Zapata and Chaudron (2013), Arnuphaptrairong (2018), Keaveney and Conboy (0000), Lederer and Mirani (1990), Lederer and Prasad (1995a), Matos et al. (2013) and Jorgensen and Molokken-Ostvold (2004)
Misunderstanding of requirements (MiRe)	Conoscenti et al. (2019), Magazinius et al. (2012), Matos et al. (2013), Jorgensen and Molokken-Ostvold (2004), Jørgensen and Gruschke (2009) and Keaveney and Conboy (0000)
Non-functional requirements (NFRe)	Usman et al. (2015), Lee et al. (2011), Silva-de-Souza and Travassos (2017) and Usman et al. (2017)
Familiar problem or requirements (FPre)	Layman et al. (2008) and Jørgensen and Gruschke (2009)
Dependencies between user stories/backlog items (DUBI)	Conoscenti et al. (2019) and Altaieb et al. (2020a)
Technical skill (TeSk)	Jørgensen et al. (2007), Furulund and Molken-stvold (2007), Jorgensen and Molokken-Ostvold (2004), Jorgensen et al. (2020) and Keaveney and Conboy (0000)
Estimation skills (EsSk)	Magazinovic and Pernstål (2008) and Keaveney and Conboy (0000)
Training in Estimation (TrEs)	Yang et al. (2008) and Rahikkala et al. (2018)
Team Size (TeSi)	Conoscenti et al. (2019), Lagerström et al. (2012), Silva-de-Souza and Travassos (2017), He et al. (2010), Huang et al. (2015), Altaieb et al. (2020a) and Hill et al. (2000)
Team Collaboration and communication (TCAC)	Yang et al. (2008), Britto et al. (2015), Usman et al. (2018a), Matos et al. (2013) and Altaieb et al. (2020a)
Turnover (Turn)	Lederer and Prasad (1995a), Lenarduzzi (2015), Magazinius and Svensson (2014), Usman et al. (2015) and Lind and Sulek (1998)
New team members (NTMe)	Yang et al. (2008), Conoscenti et al. (2019) and Keaveney and Conboy (0000)
Team Stability (Stab)	Usman et al. (2015) and Silva-de-Souza and Travassos (2017)
Team Skill (Skil)	Britto et al. (2015), Usman et al. (2015) and Usman et al. (2017)
Overlooked and unplanned tasks (OUTa)	Lederer and Prasad (1995a), Conoscenti et al. (2019), Furulund and Molken-stvold (2007), Magazinius et al. (2012), Jorgensen and Molokken-Ostvold (2004), Jørgensen and Gruschke (2009) and Lederer and Mirani (1990)
Incorrect assumptions (InAs)	Conoscenti et al. (2019), Furulund and Molken-stvold (2007) and Jørgensen and Gruschke (2009)
Occurrence of unforeseen problems (OUPr)	Yang et al. (2008), Conoscenti et al. (2019) and Jørgensen and Gruschke (2009)

over/underestimation or over/underruns. Two of them were respondent studies, and therefore the classification depended on respondents' memories (Britto et al., 2015; Lederer and Prasad, 1995a). The other study was a data one (Benschop et al., 2020). Also, two studies used pred(x) (Jørgensen, 2007a).

4.3. SQ 1.3 – What are the software project estimate variables investigated?

Regarding the project variables investigated in the primary studies, we extracted the metrics that authors reported as within their studies' scope. Fig. 4 shows the results we obtained, making evident that most of the studies focus on effort estimation.

Most of the studies focused on effort estimation (96 in total). Twenty-five studies claimed to investigate factors related to cost, while 13 focused on duration. Eight studies explored prediction intervals – mostly of effort – and we classified them separately

to emphasize the importance of avoiding single values when estimating. Three studies reported factors associated with productivity. Only two studies claim to investigate factors associated with size, probably because the focus is on other metrics when using expert judgment.

4.4. SQ 1.4 – When and where are published the studies about factors affecting expert judgment software estimates?

Our sample includes articles published between 1989 and 2020. The past two decades have been very fruitful regarding research about factors affecting estimates, as shown in Fig. 5, revealing an increasing interest in them. We also show a trend-line reporting the moving average (past five years), revealing a relative degree of stability of the number of papers published regarding factors affecting expert judgment estimates since 2016.

Table 4
Impact measurement strategy by article.

Impact measurement strategy	Article
Difference of estimates	Grimstad and Jorgensen (2009), Jorgensen and Carelius (2004), Passing and Shepperd (2003), Shepperd et al. (2018), Aranda and Easterbrook (2005), Jørgensen and Grimstad (2008), Jørgensen et al. (2007), Valerdi (2007), Jørgensen (2011), Moløkken and Jørgensen (2005), Jørgensen and Løhre (2012), Jørgensen (2010a), Lagerström et al. (2012), Jørgensen and Sjøberg (2001), Gren et al. (2017), Løhre and Jørgensen (2016), Atas et al. (2018), Jørgensen (2010b), Jorgensen and Grimstad (2012), Tripathi et al. (2017), Jørgensen and Halkjelsvik (2010), Jørgensen and Sjøberg (2004), Jørgensen and Grimstad (2011), Grimstad and Jorgensen (2007), McDonald (2005), Jørgensen (2016a), Brown et al. (2013), Subramanian et al. (2006), Huang et al. (2015), Subramanian et al. (2017), Shmueli et al. (2016), Jørgensen (2007c), Jorgensen et al. (2020), Jørgensen and Halkjelsvik (2020), Jørgensen (2015, 2014c), Halkjelsvik and Jorgensen (2011), Tamrakar and Jørgensen (2012), Jørgensen and Moløkken (2004), Hill et al. (2000), Moløkken-Østvold and Jørgensen (2004), Jørgensen (2013b), Jørgensen and Halkjelsvik (2010) and Moløkken-Østvold et al. (2008)
Participants' perception	Yang et al. (2008), Rahikkala et al. (2015a), Britto et al. (2015), Jorgensen and Carelius (2004), Passing and Shepperd (2003), Magazinovic and Pernstål (2008), Jørgensen et al. (2004), Lederer and Prasad (1995a), Conoscenti et al. (2019), Lenarduzzi (2015), Usman et al. (2018b), Magazinius and Svensson (2014), Usman et al. (2015), Tanveer et al. (2017), Usman et al. (2018a), Henry et al. (2007), Suliman and Kadoda (2017), Lee et al. (2011), Jørgensen and Sjøberg (2001), Magazinius et al. (2012), Mendes et al. (2005), Silva-de-Souza and Travassos (2017), Lederer and Prasad (1995b), Matos et al. (2013), Jorgensen and Molokken-Ostvold (2004), Grimstad et al. (2005), Jørgensen and Gruschke (2009), Rahikkala et al. (2018, 2015b), Rozalina and Mansor (2018), Altaieb et al. (2020b), Usman et al. (2017), Altaieb and Gravell (2019), Zapata and Chaudron (2013), Altaieb et al. (2020a), Magazinius and Feldt (2011), Ramessur and Nagowah (2020), Zarour and Zein (2019), Arnuphaptrairong (2018), Jørgensen (2016b), Lind and Sulek (1998), Lederer and Prasad (1998), Lederer and Mirani (1990), Glass et al. (2008) and Lederer and Prasad (1991)
MRE	Molokken-Ostvold and Jorgensen (2005), Basten and Mellis (2011), Jørgensen et al. (2007), Grimstad and Jørgensen (2007), Mendes et al. (2005), Jorgensen and Molokken-Ostvold (2004), Vicinanza et al. (1991), Jørgensen and Gruschke (2009), Karna and Gotovac (2014), Gandomani et al. (2019), Gruschke and Jørgensen (2008), Haugen (2006), Host and Wohlin (1998) and Cao (2008)
MREBias	Molokken-Ostvold and Jorgensen (2005), Jorgensen and Molokken-Ostvold (2004), Jørgensen and Gruschke (2009), Jørgensen (2013a) and Haugen (2006)
BRE	Molokken-Ostvold and Jorgensen (2005), Usman et al. (2018a), Mahnič and Hovelja (2012), Grapenthin et al. (2016), Molokken-Ostvold and Furulund (2007) and Moløkken-Østvold et al. (2008)
BREBias	Molokken-Ostvold and Jorgensen (2005), Basten and Mellis (2011), Usman et al. (2018b,a), Furulund and Molokken-Ostvold (2007), Mahnič and Hovelja (2012), Grapenthin et al. (2016) and Moløkken-Østvold et al. (2008)
Deviation	Branco et al. (2015), Jorgensen and Carelius (2004), Bergeron and St-Arnaud (1992), McGarry et al. (1998), Halstead et al. (2012), Lind and Sulek (1998) and Ohlsson et al. (1998)
Absolute error	Passing and Shepperd (2003), Gandomani et al. (2019) and Jorgensen et al. (2020)
Total effort	Mendes et al. (2005), He et al. (2010), Vijayakumar (1997), Arifin et al. (2017), Nugroho and Lange (2008) and Davis (1989)
Interval of over/underrun (over/underestimation)	Britto et al. (2015), Lederer and Prasad (1995a) and Benschop et al. (2020)
Pred(X)	Basten and Mellis (2011) and Gray et al. (1999)
Confidence related	Jørgensen et al. (2004), Jørgensen and Moløkken (2002), Løhre and Jørgensen (2016), Jorgensen (2004), Grimstad and Jorgensen (2007), Jørgensen and Gruschke (2009), Jørgensen (2016a, 2018), Gruschke and Jørgensen (2008) and Jørgensen and Teigen (2002)
Not informed/not defined	Bhatt et al. (2006), Boetticher and Lokhandwala (2007), Bukhari and Malik (2012), Morgenshtern et al. (2007), Koch and Turk (2011), Keaveney and Conboy (0000), Javed et al. (2013), Taff et al. (1991) and Bratthall et al. (2001)
Other	Passing and Shepperd (2003), Jørgensen (2014b), Layman et al. (2008), Jorgensen and Grimstad (2005), Little (2006), Vicinanza et al. (1991) and Jørgensen (2014a)

Table 6 shows all the venues concentrating three or more studies about factors affecting estimates. In total, we represent 65 articles in Table 6.

There is a balance between publishing in conferences (63 occurrences) and journals (68 occurrences). The Journal of Systems and Software, IEEE Transactions on Software Engineering and Information and Software Technology, concentrated the highest number of articles.

4.5. SQ 1.5 – What research strategies and methods are used to investigate factors that affect expert judgment software estimates?

To answer SQ 1.5, we classified the studies considering the taxonomies proposed by Storey et al. (2020), which is focused on

Table 5
Objective metrics of accuracy, bias, and error.

Accuracy metrics	#	Bias metrics	#
MRE	13	MREBias	5
BRE	6	BREBias	8

human factors of software engineering, identifying four research strategies: respondents, lab, field, and data, as we show in Table 7. Each article can report more than one study and, accordingly, could be associated with more than one research strategy.

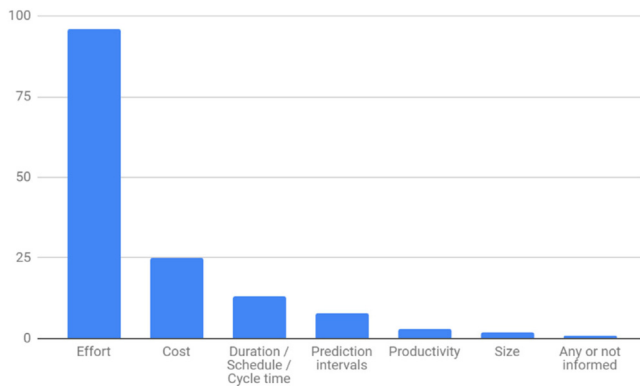


Fig. 4. Variables investigated in primary studies.

Table 6

Top venues.

Venue	# citations
Journal of Systems and Software	15
IEEE Transactions on Software Engineering	10
Information and Software Technology	8
Euromicro Conference on Software Engineering and Advanced Applications	5
International Conference on Evaluation and Assessment in Software Engineering	5
IEEE Software	4
Empirical Software Engineering	4
International Symposium on Empirical Software Engineering and Measurement	4
International Conference on Product Focused Software Process Improvement	4
International Journal of Project Management	3
International Software Metrics Symposium	3

Table 7

Research strategies distribution.

Research strategy	Number of studies
Data	31
Field	31
Lab	51
Respondents	33

In general, the different available research strategies had been used in a balanced way, except for lab strategies, which detach from the others as the most used one. That is, most of the studies in our sample evaluate one factor in a controlled setting through hypothesis testing (Storey et al., 2020). Studies investigating or reporting more than one factor generally employ respondent or field strategies, each one having 33 and 31 occurrences, respectively, in our data. In Fig. 6 we show the use of the research strategies throughout the years.

Research about factors affecting estimates became prolific after the year 2005. Since then, the distribution of studies using different strategies has been relatively uniform. However, it seems that laboratory strategies are outperforming the others in the past decade.

5. The SEXTAMT

As we informed in Section 4.1, we found a total of 235 factors, of which 69 were reported in two or more articles. We gathered these 69 factors in one instrument: the SEXTAMT. It has three dimensions to allow the navigation through the seas of factors:

1. The temporal dimension provides a view of the factors relevant for different software project or iteration phases.

2. The stakeholders' dimension focuses on the factors associated with different roles in the software process.
3. The type of effect' dimension, based on the direction of the effect of the factor.

In Fig. 7, we present the overall map of factors affecting estimates — a bird's eye view of the SEXTAMT. We represent the factors as rounded rectangles, labeled with the factors' codes we indicated in Table 3. We marked some of them with symbols related to their stakeholders' dimension. The size and color of each factor represent the number of articles investigating it. We also grouped them by major categories represented in the form of ellipses. We also provided an expanded view of Fig. 7 as part of our supplementary material (Matsubara et al., 2021), in which we added the studies that investigated each factor.

Fig. 7 shows two larger oceans, formed by categories that share common factors. The larger one contains the categories: estimation process, biases, management, experience, skill issues, team issues and project and task characteristics. It also concentrates many of the top investigated factors: the use of historical data, padding, the combination strategy of individual estimates, standards in estimation, enough effort and resources spent on estimation, overall experience, and team size.

Client/customer issues, requirements, and product' characteristics are categories that also share factors, forming another larger ocean with some of the factors that stand out: changes to requirements or scope, clear requirement specifications, misunderstanding of requirements, complexity, and product size. The map also has some categories representing smaller seas, of which political issues and unexpected events are the larger ones. Pressure and overlooked and unplanned tasks are the most investigated factors, respectively.

The remaining of this section describes the factors composing the SEXTAMT in more detail, from the perspective of dimensions we presented in Fig. 3. In each of the following subsections, we show the factors for each different class of stakeholders, organizing them per project phase. Therefore, the reader may easily navigate through the factors by stakeholder and by phase. We also present the type of effect for each factor.

5.1. Customer/client

Fig. 8 shows all the factors related to customers/clients, each one represented by a blue box. We wrote the factors using positive statements representing the presence of a factor, like in the clarity of the client's needs, representing such presence through green circles in Fig. 8. However, the existing evidence may refer to the absence of such an aspect, like the lack of clarity of the client's needs, represented in Fig. 8 by a red circle inside the factor box. Fig. 8 also presents the timeline of the typical project or iteration phases when a factor may happen or cause an impact over the estimates: the temporal dimension of the SEXTAMT. We also mapped each factor to their type of effect at the right of the figure. Some factors are organizational or overarching, and we represent them at the left of the image. We did not present their types of effects on the figure to keep it simple: we discuss it in the text only. In addition, the gray hexagons associated with each factor represent the articles that published results regarding them. The numbering of each hexagon indicates the article ID in the extraction forms (part of our supplementary material).

At the planning phase, four factors stand out. Two studies report findings related to the lack of clarity of client's needs as an error factor. Lederer and Prasad (1995a) present a survey where the users' lack of understanding of their requirements is a reason for inaccuracy. Matos et al. (2013) report a qualitative study where clients who do not know what they want hinder software estimation and accuracy in the context of web effort

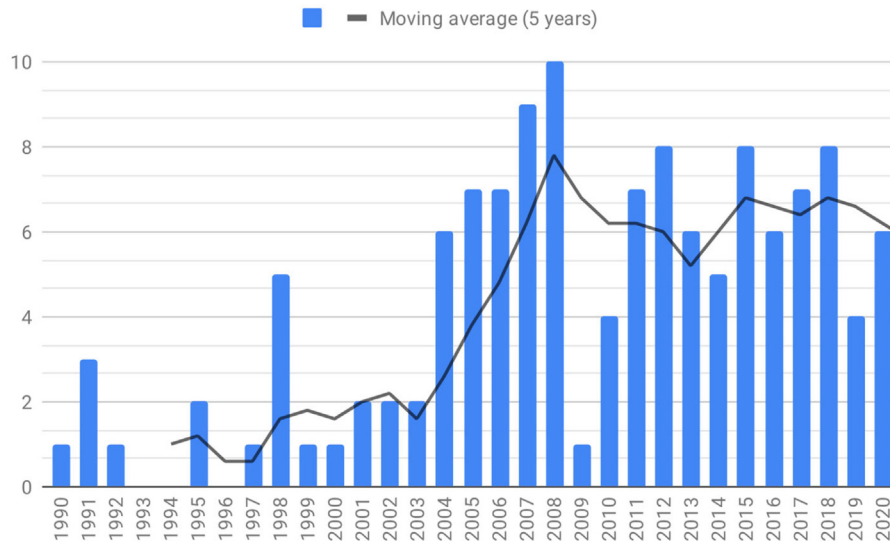


Fig. 5. Research about factors affecting the estimates over the years.

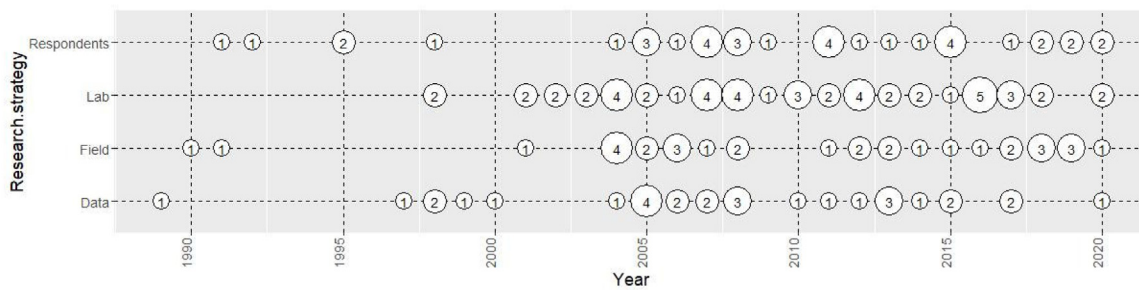


Fig. 6. Research strategies throughout the years.

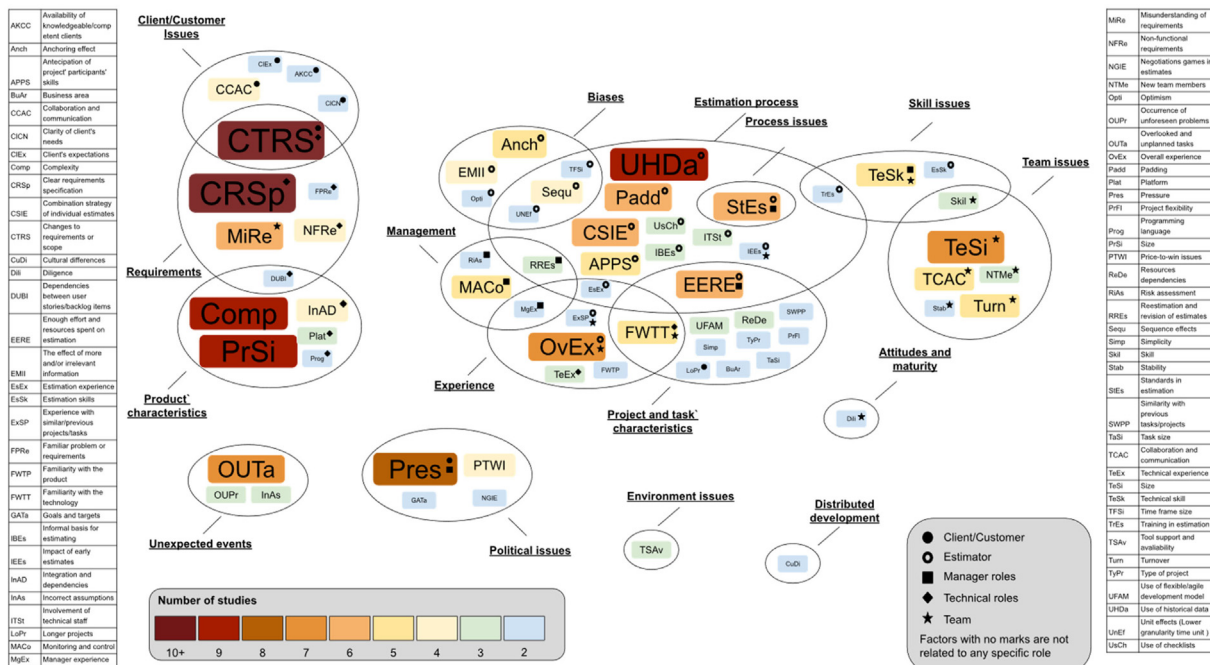


Fig. 7. The SEXTAMT.. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

estimation. Other two studies report that *longer projects* relate to higher costs (Lagerström et al., 2012) and that increasing calendar

time will increase total effort (He et al., 2010). Therefore, it is a value adjusting characteristic.

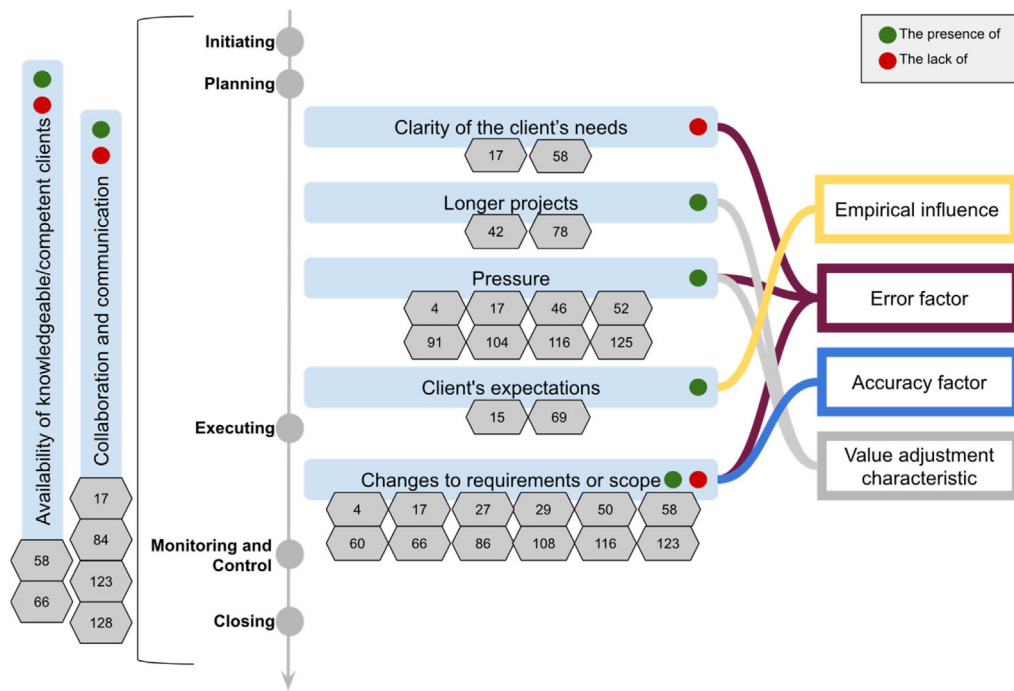


Fig. 8. Factors related to Customer/Client.. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

Eight studies declare that *pressure* impacts estimating, either as an error factor or as a value adjusting characteristic. Nevertheless, the articles describe pressure in varying levels and originating from different sources. It can, for example, be an overall pressure, directed by management or related to the schedule alone. Therefore, we created intermediary factors for pressure, and in this section, we explore only the customer pressure, which appears in two studies. *Yang et al. (2008)* point out that pressure from senior managers and clients to set or change the estimation results is a reason for inaccurate estimates. *Keaveney and Conboy (2000)* report that pressures from customers or managers result in lower estimates than would be realistically expected.

The final factor at the planning phase is the *client's expectations*, which have an empirical influence over the estimates. Estimators were impacted by the effort informed by the client at the specification of one experiment (*Jørgensen and Sjøberg, 2004*). This result repeated even when estimators are told to disregard such information (*Jørgensen and Grimstad, 2008*).

At the **executing** phase, *changes to requirements or scope* emerge as an error factor when present, with twelve studies discussing it. Some studies report that requirement changes are a reason for inaccuracies (*Layman et al., 2008; Usman et al., 2015; Zapata and Chaudron, 2013*), and two studies indicate that frequent changes are the problem (*Arnuphaptrairong, 2018; Lederer and Prasad, 1995a*). Others emphasize that requirement changes contribute to overruns (*Grimstad et al., 2005; Halstead et al., 2012*), are a challenge (*Usman et al., 2018a*), or a potential problem for estimation (*Keaveney and Conboy, 2000; Lederer and Mirani, 1990*). Finally, some researchers identify changes in scope (*Layman et al., 2008*) and scope creep (*Jørgensen and Molokken-Ostfold, 2004; Usman et al., 2015*) as reasons for inaccuracies. When the client's needs are stable, it facilitates software estimation and raises accuracy (*Matos et al., 2013*), so the absence of changes to requirement or scope is an accuracy factor.

Some factors intersect **all phases**. For instance, the availability of clients who understand the project's business rules facilitates software estimation and accuracy (*Matos et al., 2013*) - therefore, *the availability of knowledgeable/competent clients* is an accuracy

factor. Moreover, the lack of it leads to errors (*Matos et al., 2013*) and is a reason contributing to overruns (*Grimstad et al., 2005*). *Collaboration and communication* with the customer and users is an additional factor trespassing all phases. Researchers report that good collaboration with customers, facilitated by frequent communication, was associated with projects that experienced a lesser magnitude of effort overruns (*Molokken-Ostfold and Furulund, 2007*). Also, researchers found that insufficient user-analyst communication and understanding was a potential cause of estimating problems in a case study (*Lederer and Mirani, 1990*), confirming it is a reason for inaccuracy later on in a survey (*Lederer and Prasad, 1995a*). Additionally, in the agile context, customer communication is an effort predictor (*Usman et al., 2017*). Thus, *collaboration and communication* with the customer and users is an accuracy factor and a value adjusting characteristic. When absent, it is also an error factor.

5.2. Estimator

Fig. 9 presents all the factors related to anyone assuming the role of an estimator. Only one factor is related to the **initiating** phase: *early estimates* – two studies indicate that they impact estimates in later phases (*Jørgensen and Carelius, 2004; Jørgensen and Sjøberg, 2001*). In one of them, project leaders believed that pre-planning estimates impacted detailed estimates, an effect confirmed in a laboratory experiment (*Jørgensen and Sjøberg, 2001*). In a field experiment about project bidding, companies providing early price indications based on limited and uncertain information gave higher estimates in the next bidding round. Such findings surprised the researchers, who expected the early estimates to act as anchors, leading to lower bids. Next, they carried out a laboratory experiment to explore further this finding, concluding that early estimates act as anchors to final estimates only when estimators have nothing to lose (*Jørgensen and Carelius, 2004*).

All the other factors mapped to estimators concentrate on the **planning** phase. Many of them are biases, such as the *anchoring effect*, which is our tendency to be influenced by values presented to us before the estimation activity (*Løhre and Jørgensen,*

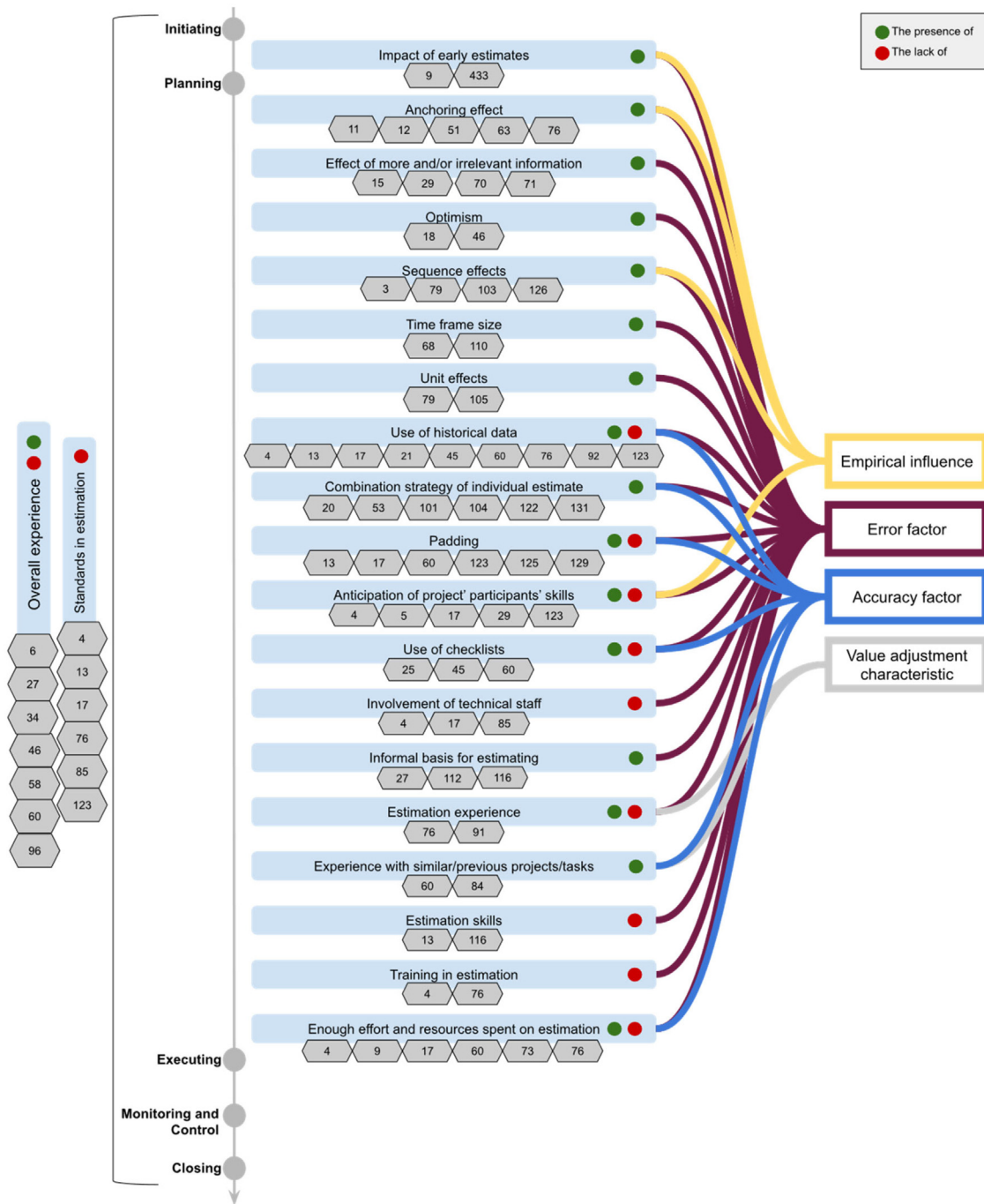


Fig. 9. Factors related to estimators.

2016). In a field study it is reported to hinder the creation of a meaningful estimate (Rahikkala et al., 2018) and, thus, is an error factor. Many laboratory experiments also report that the anchoring effect impacts software estimation (Aranda and Easterbrook, 2005; Jorgensen and Grimstad, 2012; Løhre and Jørgensen, 2016; Shepperd et al., 2018) – therefore providing evidence of its empirical influence over the estimates. Aranda and Easterbrook (2005) found a statistically significant impact of numerical anchors on time estimates. Jorgensen and Grimstad (2012) also found a significant impact of numerical anchors over estimates, reporting a medium to large effect size. They also found a small to medium effect size when using a textual anchor: putting the same requirements specification as a “minor extension” work led to lower estimates than putting it as “new functionality” work.

Løhre and Jørgensen (2016) found a slight tendency for a larger anchoring effect with interval anchors compared to single value anchors when dealing with numerical anchors. Additionally, they expected the expertise – defined as the length of experience – of the anchor's source would act as a moderator for the anchoring effect. Surprisingly, they found that the receiver's expertise that acted as such. Beyond investigating anchoring itself, Shepperd et al. (2018) discovered that raising awareness about anchoring reduces the impact of high anchors on productivity estimations but does not eliminate the effect.

Another relevant factor for estimators is the effect of more and/or irrelevant information over the estimates. Usman et al. (2018a) found that the availability of more detailed information may increase underestimation bias. Grimstad and Jorgensen

(2007) report that specifications with irrelevant information lead to higher estimates in laboratory experiments. Jørgensen and Grimstad (2008) explored different aspects of irrelevant and misleading information that have an effect over the estimates: (i) the client's cost expectations, (ii) the wording of the specification (words associated with small and simple tasks lead to underestimation, while words associated with complex and large tasks lead to overestimation), (iii) the suggestion of future opportunities for work contingent on performance in current projects (lead to underestimation), and (iv) the amount of information, even when they are irrelevant (more information leads to overestimation). Asking people to highlight relevant information or strike irrelevant ones is not enough to eliminate the observed impact (Jørgensen and Grimstad, 2008). Additionally, in a field experiment, Jørgensen and Grimstad concluded that informing that the customer required development in a short period with start-up several months ahead also led to lower estimates, though supposedly this information is irrelevant to estimation (Jørgensen and Grimstad, 2011).

Optimism is an additional error factor, leading to estimates' unintentional distortions, for instance Magazinius et al. (2012). Jørgensen et al. (2007) measured general optimism in varying ways in an experiment. They discovered that explanatory style, life orientation, and higher self-assessed level of optimism are all weakly connected with optimistic predictions. Also, merely asking estimators whether they assess themselves to be more or less optimistic seems to be enough as an indicator of optimistic predictions – instead of using more complex measures of optimism as the scales for explanatory style or life orientation (Jørgensen et al., 2007).

Estimators should also be aware of *sequence effects* relative to the order of estimation of tasks and projects with different sizes.⁷ Overall, starting the estimation processes with small tasks or projects leads to lower estimates for the subsequent project or task. The opposite happens when starting with larger ones (Grimstad and Jørgensen, 2009; Jørgensen, 2013b). Also, when estimating projects or tasks of similar sizes in a sequence, estimators tend to estimate the target project as more extensive compared to the reference project (the first one) (Jørgensen, 2013b; Jørgensen and Halkjelsvik, 2020).

Two articles address the *time frame size*: shorter time frames tend to lead to more optimistic estimates than larger ones (Halkjelsvik and Jørgensen, 2011; Jørgensen and Halkjelsvik, 2010). Another two articles investigate *unit effects*: asking for estimates using a lower granularity time unit led to lower estimates compared with using a higher granularity one (Jørgensen, 2015, 2016a). Therefore, both time frame size and unit effects are error factors.

A comprehensive set of factors affecting estimates relates to the estimation process's particularities, such as the *use of historical data*. A field study connected it with a lesser magnitude of effort overruns (Furulund and Molkken-stvold, 2007). A relevant number of studies also reported that the lack or no use of historical data is related to errors and problems in estimating – with evidence coming from respondent studies (Lederer and Prasad, 1995a; Yang et al., 2008), laboratory studies (Shmueli et al., 2016), and field studies (Conoscenti et al., 2019; Jørgensen and Molokken-Ostfold, 2004; Lederer and Mirani, 1990; Magazinius and Pernstål, 2008; Rahikkala et al., 2018).

The *combination strategy of individual estimates* rose as a factor in our SLM, either for combining single values or interval estimates – with minimum and maximum values. We found evidence for three strategies regarding single values: statistical combination, unstructured group estimates, or Planning Poker. Three

articles report evidence in favor of estimating in groups over averaging: unstructured group estimates (Molokken-Ostfold and Jørgensen, 2004) and Planning Poker (a structured approach) (Gandomani et al., 2019; Molokken-Ostfold et al., 2008) led to less optimistic estimates compared with the average of individual estimates. When combining interval estimates, the results also favor group discussion over averaging (Jørgensen and Molokken, 2002). Mahnic and Hovelja (Mahnič and Hovelja, 2012) found the same result for Planning Poker estimates compared with the statistical combination, but only when the participants in their experiments were software professionals. They found the opposite effect when students were estimating. In another study, the results suggested that planning poker is more accurate when the team has previous experience from similar tasks compared to unstructured group estimation sessions (Haugen, 2006). In summary, there is evidence in favor of estimating in groups over averaging estimates in general and in favor of Planning Poker more specifically.

Padding also impacts estimates' accuracy. The inclusion of a large buffer to deal with unexpected events and/or changes in the specification is a reason for accurate estimates (Jørgensen and Molokken-Ostfold, 2004). The greater the preference for projects completed within the estimates, the greater the padding frequency (Lederer and Prasad, 1991). More evidence about it comes from the fact that the removal of padding by management is related to estimating problems (Lederer and Mirani, 1990; Magazinius and Pernstål, 2008) and is a reason for inaccuracies (Lederer and Prasad, 1995a). Nevertheless, it is reported as an intentional increase in estimates aimed at the holding back reserves, which gives it a negative denotation (Glass et al., 2008).

The *anticipation of project' participants skills* emerged as a relevant factor for estimators. The inability to anticipate the team members' skills, abilities, or characteristics is a problem for estimating (Lederer and Mirani, 1990) and a reason for inaccuracies (Lederer and Prasad, 1995a; Yang et al., 2008). The knowledge about who will execute testing allows for the definition of testing effort (Rahikkala et al., 2015a). However, one study suggests that the team's knowledge of who will work on the project may increase underestimation bias (Usman et al., 2018a). It might be the case that anticipating the project participants' skills may not work for all contexts.

Another essential aid is the *use of checklists*, leading to a lesser magnitude of effort overruns (Furulund and Molkken-stvold, 2007). Personalized checklists reduces the underestimation bias (Usman et al., 2018b). Such evidence indicates that the use of checklists is an accuracy factor. Also, the lack of checklists is a reason for estimation error (Jørgensen and Molokken-Ostfold, 2004).

The lack of *involvement of technical staff* during estimating is a reason for inaccuracies in three respondent studies (Altaieb and Gravell, 2019; Lederer and Prasad, 1995a; Yang et al., 2008). Other three studies (Keaveney and Conboy, 0000; Lederer and Prasad, 1998; Usman et al., 2015) also reported that an *informal basis for estimating* is an error factor. Lederer and Prasad (1998) considered informal bases for estimating, comparing similar, past projects based on personal memory, guessing, and intuition as reasons for inaccuracies. The other two studies emphasized the lack of formality of the estimation process as a reason for inaccurate estimates (Keaveney and Conboy, 0000; Usman et al., 2015).

Four factors associated with estimators regard their experience and skills. The first one is the *estimation experience*: an effort predictor in the context of mobile development (Altaieb et al., 2020a), whose absence hinders the creation of a meaningful estimate (Rahikkala et al., 2018). The second is *experience*

⁷ The use of the word size here is for simplicity. A task or project is larger in the sense that it requires more effort to be executed/implemented compared to others.

with similar/previous projects/tasks, which is also an effort predictor (Usman et al., 2017) and a reason for accurate estimates (Jorgensen and Molokken-Ostfold, 2004). The third factor is the lack of *estimation skills*, an estimation inhibitor (Magazinovic and Pernstål, 2008) that can cause estimation problems (Keaveney and Conboy, 0000). The fourth is the lack of *training in estimation*, which hinders creating a meaningful estimate (Rahikkala et al., 2018) and is a reason for inaccurate estimates (Yang et al., 2008).

The final factor related to estimators at the planning phase is *enough effort and resources spent on estimation*, which is an accuracy factor and, when lacking, an error factor. On the one hand, a respondent study reports that spending enough time on estimating is a reason for accurate estimates (Jorgensen and Molokken-Ostfold, 2004). On the other hand, making quick, rough estimates is not motivating and hinders creating a meaningful estimate (Rahikkala et al., 2018). Also, insufficient time, effort, or resources for estimating is a reason for inaccurate estimates (Jorgensen and Molokken-Ostfold, 2004; Jørgensen and Gruschke, 2009; Magazinovic and Pernstål, 2008; Yang et al., 2008).

Two factors intersect **all the phases**. One of them is the *overall experience* of the estimator. In one study, experts' experience (including total experience, company experience, project experience, and the number of projects expert has participated) predicted estimation performance, leading to less estimation error (Karna and Gotovac, 2014). Therefore, the presence of overall experience improved accuracy. Additionally, other studies indicate that the lack of overall experience is an error factor, leading to unintentional distortions of software estimates in varying directions – reducing or increasing them (Magazinovic et al., 2012), hindering software estimation and accuracy (Matos et al., 2013), being a reason for estimation error (Jorgensen and Molokken-Ostfold, 2004).

The other factor affecting all phases is *standards in estimation*. All evidence about it is related to its shortage, and all results point to it as an error factor. It has many facets, in any case. For instance, in one case study, participants revealed that the lack of methodology or guidelines and the lack of setting and review of standards is a potential cause of estimating problems (Lederer and Mirani, 1990). A follow-up survey confirms that these are reasons for inaccuracies (Lederer and Prasad, 1995a). Also, no development of estimation standards and no record-keeping of estimates and actual results make it difficult to capitalize on lessons learned (Magazinovic and Pernstål, 2008), and no documented estimation procedure hinders the creations of a meaningful estimate (Rahikkala et al., 2018). Researchers also report that the lack of appropriate software cost estimation methods and processes (Yang et al., 2008) and the lack of clear guidance for estimating (Altaieb and Gravell, 2019) are reasons for the inaccuracy of estimates.

5.3. Management roles

We present the factors regarding management roles – including higher management, project managers, and the Software Engineering Process Group (SEPG) – in Fig. 10. We explored some of them thoroughly in previous sections: longer projects (Section 5.1), enough effort and resources spent on estimation (Section 5.2), and standards in estimation (Section 5.2). We explore all the others in the current section.

At the **planning phase**, *pressure* came up as an error factor. Yang et al. (2008) report that the company's survival pressure and that the senior manager's pressure to set or change the estimation results is a reason for inaccurate estimates – a finding that echoes in other studies (Lederer and Prasad, 1995a; Zarour and Zein, 2019). It leads people to change their estimates intentionally (Magazinovic et al., 2012), to cave in to people with

more power (Glass et al., 2008), resulting in lower estimates than would be realistically expected (Keaveney and Conboy, 0000). Another facet of pressure is work pressure, which Altaieb, Altherwi, and Gravell report as an effort predictor (Altaieb et al., 2020a). A final facet is schedule pressure, which leads to more effort in test tasks (Silva-de-Souza and Travassos, 2017) – and thus is a value adjusting characteristic.

Risk assessment is another factor in the planning phase. Systematic risk assessment related to lower error in duration estimates (Morgenshtern et al., 2007), and the lack of it is a reason for inaccurate estimates (Yang et al., 2008). Surprisingly, some laboratory experiments' results indicate that identifying more risks immediately before software estimation leads to increased overconfidence (Jørgensen, 2010a). Nevertheless, the authors stress that they have not investigated a complete risk management process – only the impact of simple risk identification.

Low *technical skills* also are among factors related to managers. One study report that project managers not skilled in planning multi-disciplinary projects are reasons for estimation error (Jorgensen and Molokken-Ostfold, 2004). Other studies also report technical skill issues but concerning the team, and we describe them further in Section 5.5.

At the **executing phase**, the only factor is the *reestimation and revision of estimates*. In a large company with two estimation points in its process, the reestimation at the analysis stage improves the accuracy of the effort estimates (Usman et al., 2018a). In a data study, more budget revisions were related to higher costs (Lagerström et al., 2012) – and therefore, we considered it a value adjusting characteristic. Nevertheless, in another data study, more estimation updates were connected with larger errors in effort estimates (Layman et al., 2008). Regarding the last result, the authors explain that more extensive features had more frequent estimation updates. Another possible explanation is that projects already in trouble may undergo more estimation updates.

The only factor at the **monitoring and control** phase is its homonym and is an accuracy factor. One field study reports that good cost control is a reason for accurate estimates (Jorgensen and Molokken-Ostfold, 2004). One a respondent study reports that adequate project administration is a reason for the prevention of overrun (Grimstad et al., 2005).

The factor that intersects **all phases** is the *manager's experience*. For instance, the number of projects previously managed correlates with duration error – more projects managed leads to lower error (Morgenshtern et al., 2007). It is, therefore, an accuracy factor. Also, when the estimates used for the project contract are based on the project manager's previous experience only, it requires the developers to work over their capacity, which is a reason for low accuracy (Altaieb et al., 2020b).

5.4. Technical roles

We found factors related to technical roles: requirement engineers, software designers, developers, and testers. Fig. 11 brings such factors to the surface. None of them apply to all phases. We explained two of them in Section 5.1: changes to requirements or scope and pressure – including the factor associated with the **tester** role.

We found four factors related to requirements at the **planning phase**, which we associated with the **requirements engineer** role. One of them is a clear *requirements specification*. Some studies present results in more general terms, indicating that poor, unclear, or ill-defined requirements are one reason for inaccuracies (Grimstad et al., 2005; Jorgensen and Molokken-Ostfold, 2004; Lederer and Prasad, 1995a; Usman et al., 2015; Yang et al., 2008; Zarour and Zein, 2019). Other studies emphasize specific

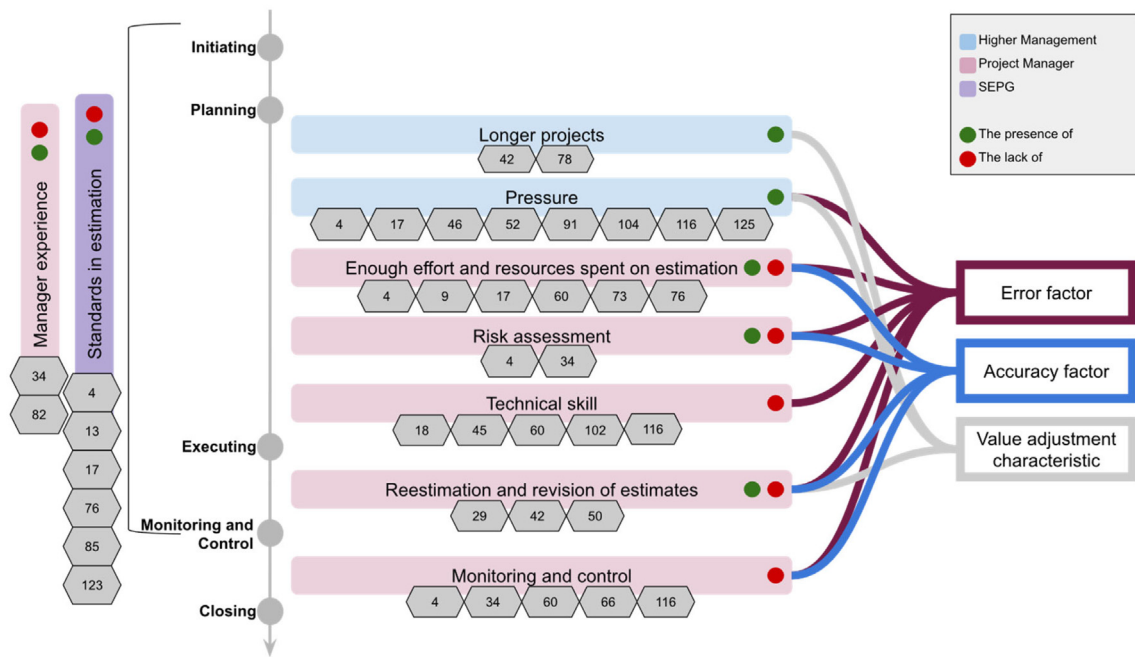


Fig. 10. Factors related to managers.

facets that make requirements poor, like the redundancy of user stories (Conoscenti et al., 2019), missing requirements (Usman et al., 2015), weak or ambiguous requirements (Furulund and Molkken-stvold, 2007), incomplete requirements (Jorgensen and Molokken-Ostvold, 2004), and the user's lack of understanding of their requirements (Arnuphaptrairong, 2018). All this evidence indicates that the lack of clear requirements specifications is an error factor.

Familiar problems or requirements was also classified as an error factor when they are absent. Layman et al. (2008) report that unfamiliar feature requirements are a reason for estimation inaccuracy. Jørgensen and Gruschke (2009) report that too little knowledge about the problem is a reason for estimation inaccuracy.

The third factor associated with the requirements engineer is dependencies between user stories/backlog items. Conoscenti et al. (2019) found that links to other stories serve as indicators for a possible inaccurate estimation. Altaieb et al. (2020a) found that dependency between backlog items is an effort predictor in the mobile development context.

The fourth factor we found regards studies reporting that non-functional requirements are an effort predictor or a cost driver (Usman et al., 2017, 2015). We also found studies reporting that specific non-functional requirement types are associated with higher effort, like the high legal or regulatory impact of the code (Lee et al., 2011), the required level of performance, and the required security level (Silva-de-Souza and Travassos, 2017). So, we classified it as a value adjustment characteristic.

Still in the **planning phase**, three factors emerge for the **developer** role. One of them is *integration and dependencies*. One study report that technical dependencies are an effort predictor in agile global development (Britto et al., 2015). Another one considers that integration issues are a cost driver, also in the context of agile development (Usman et al., 2015). In the context of corrective maintenance of object-oriented systems, a high level of code/system dependencies leads to higher effort (Lee et al., 2011). Therefore, the integration and dependencies factor is a value adjustment characteristic. Another study informs that integration complexity is an estimation challenge (Magazinius and Svensson, 2014), suggesting it is also an error factor.

The other factor regarding developers is the *platform*. In the context of mobile development, the supported Platform type (IOS/Android./Win./etc.) and the supported device (phone, tablet, smartwatch) are both effort predictors (Altaieb et al., 2020a). Other two studies report that the type of platform impacts software costs (Lagerström et al., 2012) and that the interaction of team size and development platform has a significant impact on productivity (Huang et al., 2015).

Finally, the developer's *familiarity with the product* is a value adjustment characteristic. When low, it leads to more need for effort (Lee et al., 2011). In another study, the programmer's familiarity in the number of months of experience with the system was a significant predictor of debugging effort (more experience leads to less effort) (Davis, 1989).

Two data studies inform the *programming language's* importance as an empirical influence over the estimates related to the **developer** role at the **executing phase**. It has a significant impact on the effort needed (He et al., 2010) and on time-to-market (Huang et al., 2015). Huang et al. (2015) also report that team size and language type interaction significantly impact productivity.

The *technical experience* related to the **developer** role is an additional factor we found. Altaieb et al. (2020a) evidence that developer implementation experience is an effort predictor. Also, developers' lack of experience leads to estimation inaccuracy (Conoscenti et al., 2019), and the lack of technology experience leads to a higher probability of effort overrun (Halstead et al., 2012).

5.5. Team

Some of the factors we found regarded the whole software team. We show them in Fig. 12. We thoroughly discussed two of these factors in previous sections: *involvement of technical staff* in estimating and *experience with similar/previous projects/tasks* – both at Section 5.2.

At the **planning phase**, *familiarity with the technology* is a value adjustment characteristic because when it is low, it leads to a higher need for effort (Lee et al., 2011). Other studies

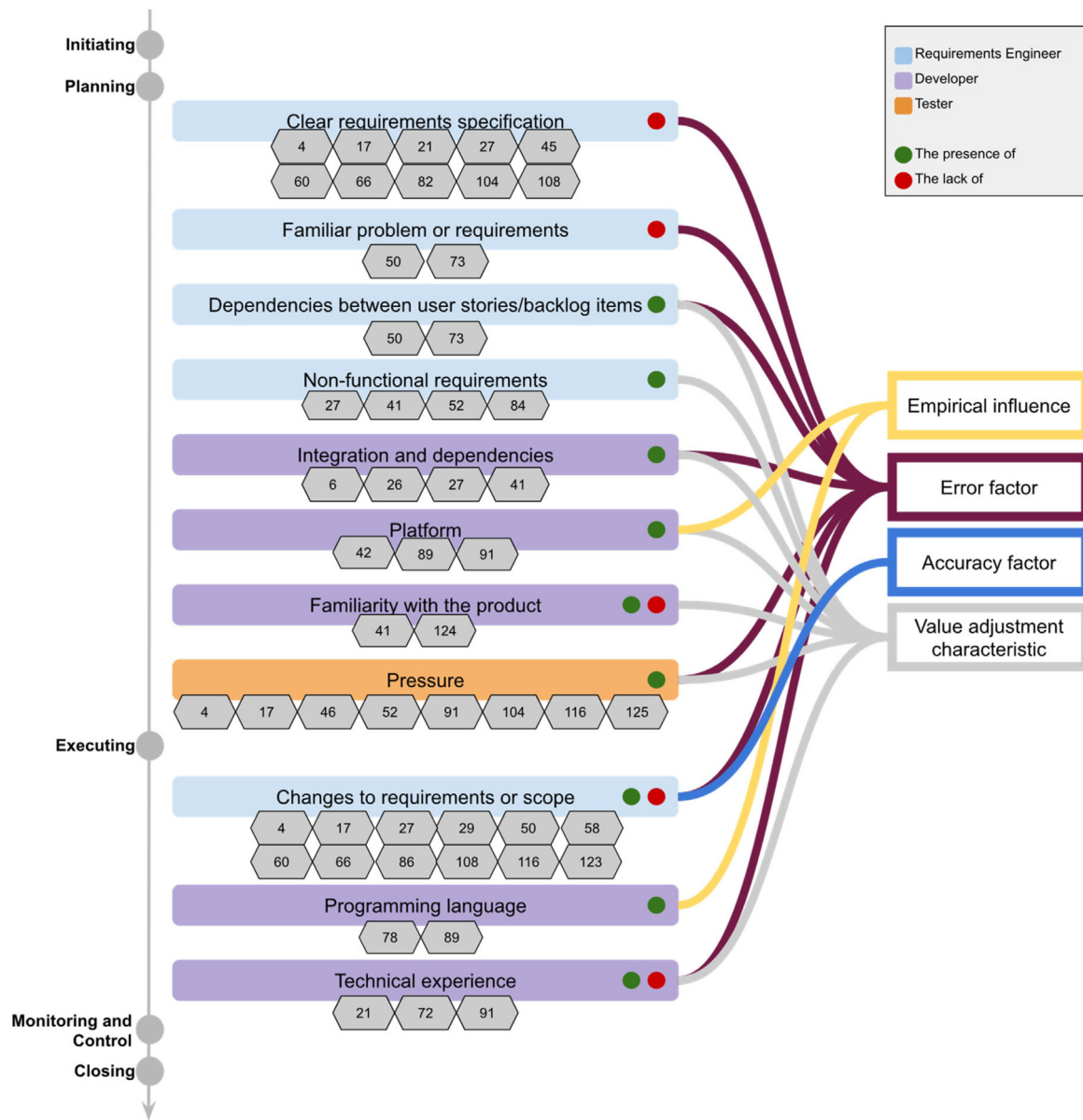


Fig. 11. Factors related to people in technical roles.

also indicate that the use of new or little-known technology is a reason for estimation inaccuracies (Basten and Mellis, 2011; Furulund and Molkenstovold, 2007; Jørgensen and Gruschke, 2009) and a significant threat to estimates (Keaveney and Conboy, 0000). Also, many studies report results regarding how the *misunderstanding of requirements* leads to estimation inaccuracy and errors (Conoscenti et al., 2019; Jørgensen and Molokken-Ostovold, 2004; Jørgensen and Gruschke, 2009; Keaveney and Conboy, 0000; Matos et al., 2013). It also causes unintentional distortions of software estimates in different directions: either as increases or decreases of estimates (Magazinius et al., 2012).

The team's *skill* is a value adjustment characteristic at the **executing phase** once three studies present it as either an effort predictor or a cost driver (Britto et al., 2015; Usman et al., 2017, 2015). Another more specific factor is the *technical skill*, which we partially addressed in Section 5.3. The presence of unskilled members in the team is a reason for inaccurate estimates (Usman et al., 2015). Lack of technical skills (Furulund and Molkenstovold, 2007) and technical expertise in a particular

area (Keaveney and Conboy, 0000) lead to estimation inaccuracies. Less software development skill is weakly connected with optimistic predictions too (Jørgensen et al., 2007). More specifically, Jørgensen et al. (2020) reported that lower programming skills connect with higher over-optimism in larger tasks and higher over-pessimism in smaller tasks.

Two respondent studies report how *diligence* issues may impact estimates negatively. Lack of diligence by systems analysts and programmers is a reason for inaccuracy (Lederer and Prasad, 1995a). Also, the delay of decisions concerning requirements due to team members' lack of responsibility and motivation is a reason for a higher need for effort than estimated (Basten and Mellis, 2011). So, lack of diligence is an error factor.

Many studies report findings regarding a range of issues related to team's size and stability issues. The team's size is an effort predictor (Altaleb et al., 2020a; Huang et al., 2015), and larger teams connect with higher effort and costs (He et al., 2010; Lagerström et al., 2012; Silva-de-Souza and Travassos, 2017). It is, therefore, a value adjustment characteristic. The interaction of team size and language type and the interaction of team size and

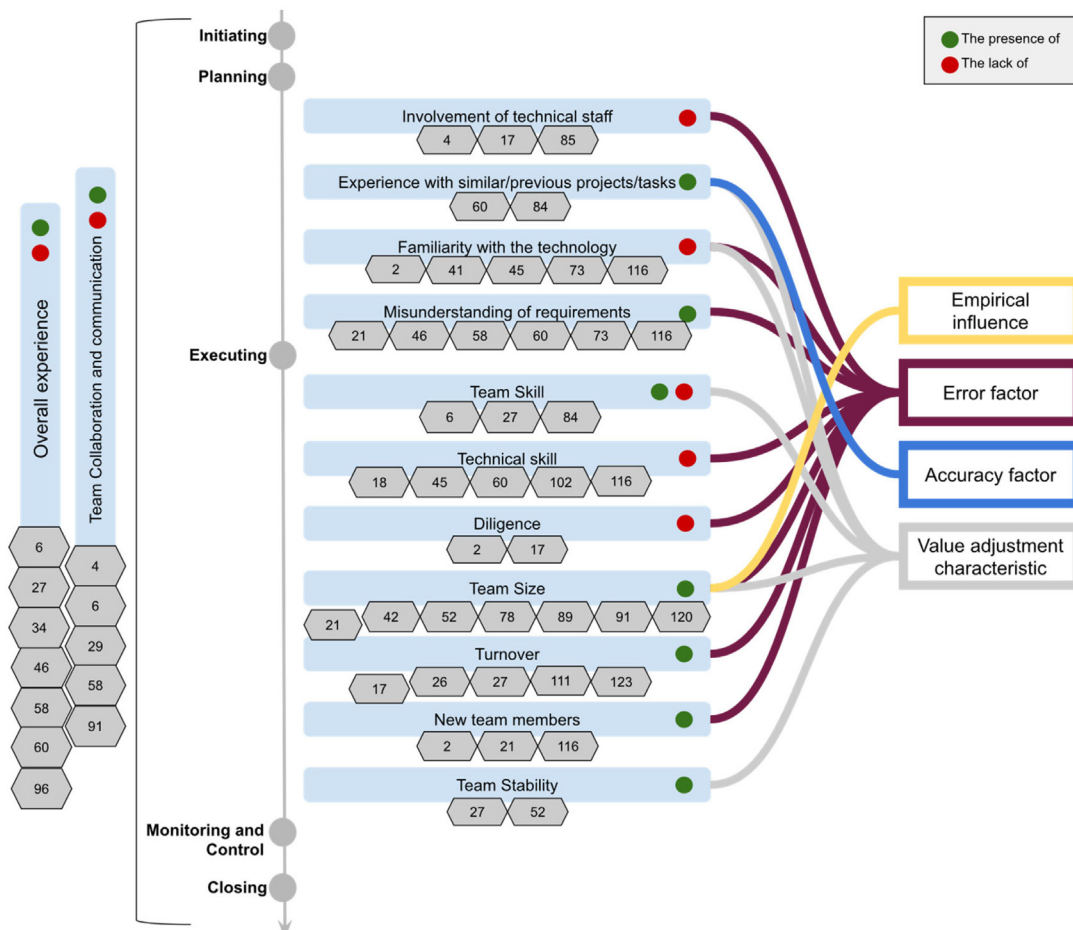


Fig. 12. Factors related to the team.

development platform significantly impact productivity (Altaieb et al., 2020a; Huang et al., 2015). Interestingly, two studies suggest that multiple developers' involvement in a story or a task may lead to over or underestimations (Conoscenti et al., 2019; Hill et al., 2000). So, larger team size also is an error factor.

The last three factors of the executing phase are intricately connected. *Turnover* is a reason for inaccuracies in estimates (Lederer and Prasad, 1995a; Lind and Sulek, 1998; Usman et al., 2015) and estimating problems (Lederer and Mirani, 1990). The loss of organizational knowledge due to high turnover is an estimation challenge (Magazinius and Svensson, 2014). The existence of *new team members* leads to estimation inaccuracies (Conoscenti et al., 2019) and a higher need for effort than estimated (Basten and Mellis, 2011). Another study reports that the introduction of new people is a major threat to accurate estimates (Magazinius and Svensson, 2014) – and therefore, we classified it as an error factor. Finally, regarding *team stability*, one study reports it as a cost driver (Usman et al., 2015), while another one stresses that team continuity leads to less effort in the context of testing tasks (Silva-de-Souza and Travassos, 2017). Therefore, team stability is a value adjustment characteristic that estimators should account for when estimating.

Two factors impact **all phases**. The *team's overall experience* is one of them – and we explored some of its facets in Section 5.2. Three studies report it as more specifically connected with the team. Two respondent studies put the team's overall experience as an effort predictor or a cost driver (Britto et al., 2015; Usman et al., 2015). Another respondent study indicates that low team experience correlates with duration error (Morgenshtern et al., 2007).

The other factor related to all phases is *collaboration and communication*. The communication process and the communication model are effort predictors (Altaieb et al., 2020a; Britto et al., 2015). On the one hand, team collaboration facilitates software estimation and accuracy (Matos et al., 2013). On the other hand, the lack of stakeholder collaboration is a reason for inaccurate estimates (Yang et al., 2008). Also, inherent difficulties related to communication and coordination present in multi-site arrangements lead to higher effort overruns (Usman et al., 2018a).

5.6. No specific role

In Fig. 13, we present a whole set of factors we found that is not specifically connected with any roles. They may impact or be caused by any or all of them.

During the **planning phase**, *price-to-win issues* play a role in estimation when present. Price-to-win is described as an estimate defined by the price or schedule needed to win a job (Boehm, 1984). An estimate strongly impacted by price-to-win is a reason for estimation error (Jorgensen and Molokken-Ostfold, 2004). Allowing the project bidding requirements to predefine the project cost (Yang et al., 2008) or purposefully underestimating the effort to obtain a contract (Usman et al., 2015) are reasons for inaccurate estimates. Magazinius et al. (2012) also report intentional distortions of software estimates in varying directions because estimates are budget determined. Somewhat related is the *goals and targets* factor. In field studies, the authors report that personal goals affect the estimates (Magazinovic and Pernstål, 2008), and

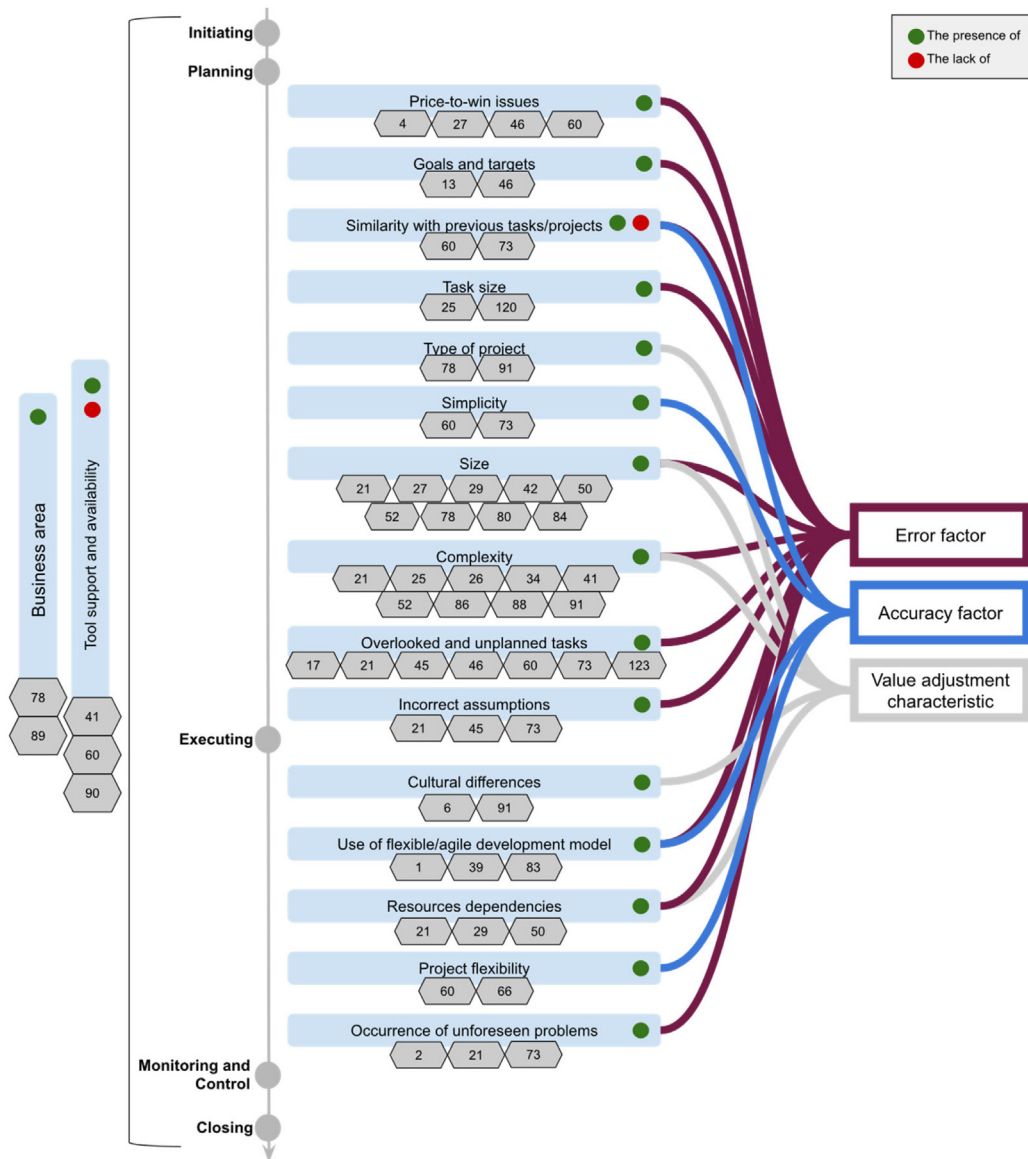


Fig. 13. Factors unrelated to any specific role.

that personal or organizational agendas lead to intentional distortions of software estimates in varying directions (Magazinius et al., 2012).

We identified that some of the project and task characteristics also are relevant factors for estimation, such as the *similarity with previous tasks/projects*. Task similarity is a reason for improving estimation accuracy (Jørgensen and Gruschke, 2009). However, projects frequently different from previous ones are a reason for estimation error (Jørgensen and Molokken-Ostfold, 2004). The *task size* is also an error factor: larger tasks are more prone to effort overruns (Usman et al., 2018b), and tasks with more subtasks were underestimated compared to tasks involving fewer ones (Hill et al., 2000). *Project type* also matters: whether it is related to a new or enhanced application in mobile development (Altaleb et al., 2020a). He et al. (2010) also report that re-development needs less effort than enhancement, and new development consumes even less than re-development. Therefore, the project type is a value adjustment characteristic. Finally, two studies inform that task or project *simplicity* is a reason for accuracy (Jørgensen and Molokken-Ostfold, 2004; Jørgensen and Gruschke, 2009).

A subset of the planning phase factors regards the product characteristics: the product size and *complexity*. *Size* is a value adjustment characteristic since many studies report it as a cost driver, effort predictor, or as correlated to effort (He et al., 2010; Usman et al., 2017, 2015; Vijayakumar, 1997) – with larger project sizes leading to more effort (Lagerström et al., 2012). *Size* is also an error factor. For instance, Conoscenti et al. (2019) report that user story size serves as an indicator for a possible inaccurate estimation. In a data study, more extensive features correlated to larger errors in effort estimates (Layman et al., 2008). Finally, a field study indicates that smaller product customizations tend to be overestimated, while larger ones tend to be underestimated (Usman et al., 2018a).

Complexity is a factor with many facets. Requirements complexity (Silva-de-Souza and Travassos, 2017) and high technical complexity (Lee et al., 2011; Silva-de-Souza and Travassos, 2017; Subramanian et al., 2006) leads to more effort. In mobile development, application form complexity is an effort predictor (Altaleb et al., 2020a). Therefore, complexity is a value adjustment characteristic. Some studies report technical complexity (Magazinius and Svensson, 2014; Morgenshtern et al., 2007; Usman

et al., 2018b; Zapata and Chaudron, 2013) and feature complexity (Conoscenti et al., 2019; Magazinius and Svensson, 2014) as estimation challenges or as related to inaccuracies, delays, and under or overestimations.

Overlooked and unplanned tasks is another impacting error factor: it is a challenge for estimation (Lederer and Mirani, 1990) and a source of inaccuracies and errors (Conoscenti et al., 2019; Furulund and Molkken-stvold, 2007; Jorgensen and Molokken-Ostfold, 2004; Jørgensen and Gruschke, 2009; Lederer and Prasad, 1995a; Magazinius et al., 2012). Unplanned tasks or re-work also is a reason for estimation error (Jorgensen and Molokken-Ostfold, 2004). Closely related, *incorrect assumptions* when estimating is also an error factor that may be related to the code (Jørgensen and Gruschke, 2009), functionality (Conoscenti et al., 2019), or complexity (Furulund and Molkken-stvold, 2007; Jørgensen and Gruschke, 2009).

At the **execution phase**, distributed development issues also play a role when they are present. Two studies report *cultural differences* as an effort predictor (Altaieb et al., 2020a; Britto et al., 2015). Thus, estimators should consider it a value adjustment characteristic if there are multiple development sites with differing cultures.

The *use of flexible/agile development models* is an accuracy factor regarding project and task characteristics. Molokken-Ostfold and Jorgensen (2005) report that flexible models are associated with lower effort overruns than sequential models. Koch and Turk (2011) also report that the use of agile methods is related to less effort deviation from estimate than rigid models. However, Brown et al. (2013) inform that software developers give lower estimates when the development method is agile than when the development method is a waterfall, suggesting their estimates were optimistic.

Resources dependencies also stood out as one factor affecting estimates. Depending on external resources may lead to delays and/or higher effort that should be considered when estimating (Jørgensen, 2011). Also, dependencies (such as for code reviews) on specific human resources (e.g., product architects) introduce delays (Usman et al., 2018a), and developer resource constraints and external commitments are a reason for estimation inaccuracy (Layman et al., 2008).

Project flexibility is another relevant accuracy factor: a high degree of flexibility in implementing the requirement specification is a reason for accurate estimates (Jorgensen and Molokken-Ostfold, 2004). Another study reports that project flexibility to reduce the scope (functionality/quality) in order to meet plan and budget is a factor more frequent in projects with lower overrun (less than 20% overrun) compared to projects with higher overrun (more than 20% overrun) (Grimstad et al., 2005).

The *occurrence of unforeseen problems* is a factor that impacts estimates negatively. The occurrence of risks, unexpected events, or technical problems leads to a higher need for effort than estimated and estimation errors (Basten and Mellis, 2011; Conoscenti et al., 2019; Jørgensen and Gruschke, 2009).

Two of the factors affect **all phases**. The *business area* has an impact on the effort (He et al., 2010) and productivity (Huang et al., 2015). The other factor is *tool support and availability*. Software development tools have an empirical influence over management and testing efforts (Subramanian et al., 2017). Additionally, insufficient tool support for project management is a reason for estimation error (Jorgensen and Molokken-Ostfold, 2004), and the low availability of required tools leads to higher effort (Lee et al., 2011).

6. Discussion

Our primary research question for this SLM was **RQ 1 - How have researchers investigated the factors that affect expert judgment software estimation?** In this section, we summarize our current answer to this question and discuss our findings.

6.1. The seas of factors that researchers explored the most

The top-five factors in the SEXTAMT regarding the number of articles reporting them are *changes to requirements or scope* (12 articles), *clear requirement specifications* (10 articles), *product size, complexity, and use of historical data* (9 articles each). Most factors (40, representing around 58% of the total) were reported in three or more articles. The remaining 29 factors (around 42%) were investigated in two research articles only, indicating that they could benefit from further investigation.

In addition, many of the top factors were probably investigated extensively because they have a relevant impact on the estimates. Nevertheless, others may have been investigated because of a controversial result. Controversies possibly exist either because of differences in research design or because such factors are more sensitive to the context. Future research efforts should aim to clarify which is the case. For instance, regarding the *combination strategy of individual estimates* most of the results shows that group estimation led to less optimistic estimates compared with averaging. However, one study found the opposite when participants were students (Mahnič and Hovelja, 2012). It is unclear whether this controversial result is due to the difference in choice of participants (software professionals or students) or whether experience interacts with the combination strategy to define which one will bring superior results (more on this in Section 6.5).

Also, if a factor is shown to influence estimates through the employment of varied research strategies, we can more confidently believe that such an effect exists. Each research strategy has its inherent limitations and strengths (Stol and Fitzgerald, 2018). Also, each one has the potential to maximize one research quality criteria at the expense of others. For instance, respondent strategies have the potential to maximize generalizability; field strategies can maximize realism; laboratory strategies can maximize control; and data strategies can maximize precision (Storey et al., 2020). Therefore, we evaluated the existing evidence for the factors in the SEXTAMT by considering the research strategies that researchers employed to investigate them.

Fig. 14 represents only the factors investigated in five or more articles – 21 factors in total, represented by the light gray edges surrounding the top of the circle. We also mapped the factors to the research strategies that researchers employed to investigate them, represented at the bottom of the circle: respondent (R, in dark red), field (F, in blue), data (D, in dark gray), or laboratory (L, in orange). Next, we discuss the type of evidence derived from such studies, considering all these research strategies.

First, six factors have been investigated employing at least three different research strategies: *product size* (1 R, 4 F, 4 D), *complexity* (4 R, 4 F, 1 D), *use of historical data* (2 R, 5 F, 1 D, 1 L), *overall experience* (4 R, 2 F, 1 D), *team size* (1 R, 2 F, 4 D) and *turnover* (2 R, 2 F, 1 D). Most of them were investigated through a combination of research, field, and data strategies – suggesting the generalizability, realism of context, and precision of data regarding the supporting findings. Some of these factors are classic cost drivers, such as *product size* and *complexity*, and software companies may not have much control over them. Other factors are more controllable but may not be so easy to implement. Still, software practitioners and organizations can organize themselves to *use historical data* when estimating, increase their *overall experience*, regulate *team sizes* to keep them small, and reduce *turnover*.

All the remaining factors in Fig. 14 were investigated using two different research strategies. In summary, these factors indicate that improving the estimation process is necessary, but not enough to get better results. Practitioners need to enhance the requirements engineering and management process. For example, we need to work on getting *clear requirements specifications* (6

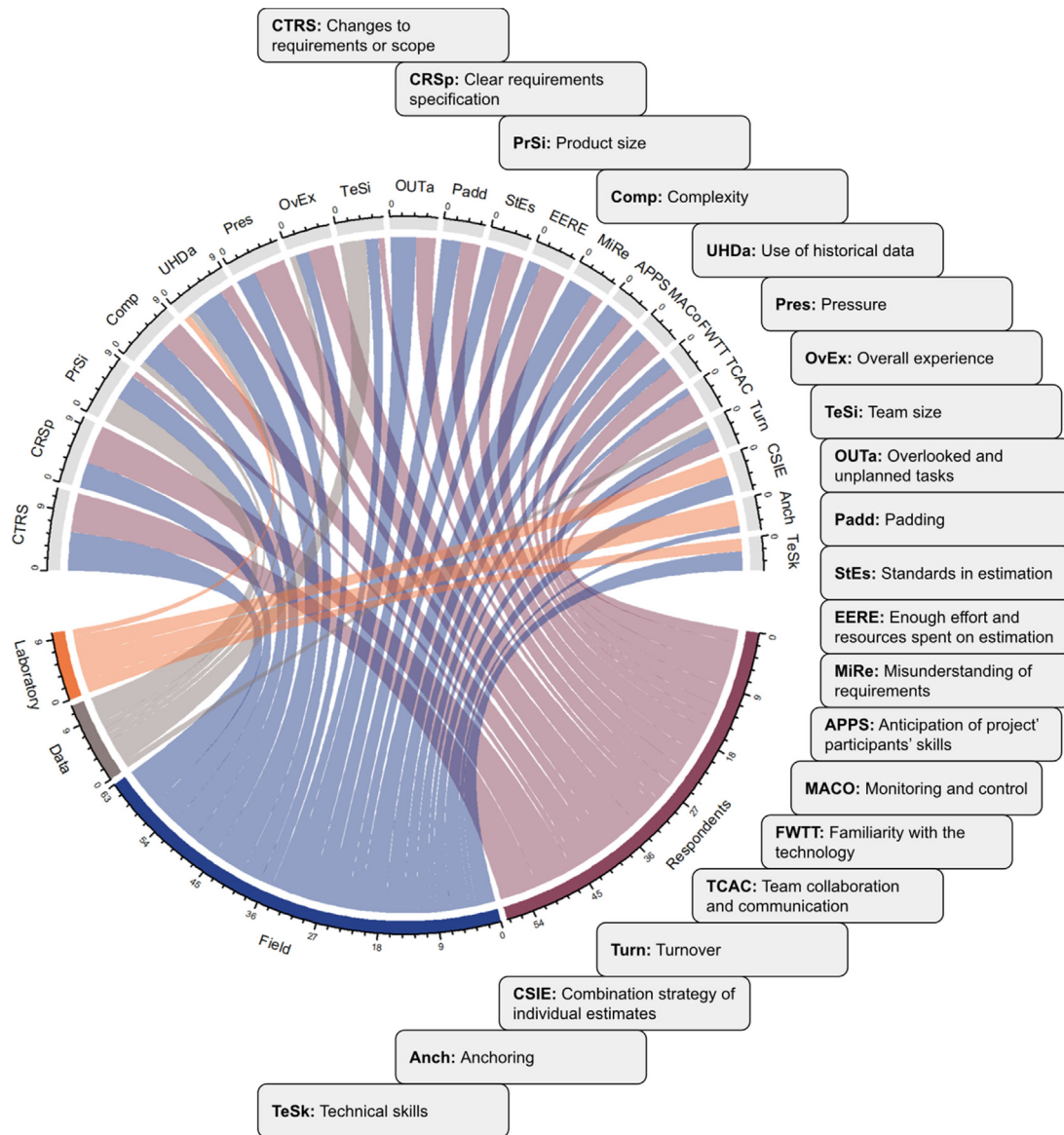


Fig. 14. Top factors and the studies' research strategies.. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

R, 4 F). Moreover, human and social aspects play a crucial role, as highlighted by factors pertaining to the categories of political issues (such as *pressure* – 5 R, 3 F), experience (*familiarity with the technology* – 3 R, 2 F), skill issues (*technical skill* – 3 F, 2 L), biases (*anchoring* – 1 F, 4 L), and team issues (*team communication and collaboration* – 4 R, 1 F). Unexpected events also have their role: *overlooked and unplanned tasks* (3 R, 4 F). Reducing such events is necessary – possibly with the *use of checklists*, another factor from SEXTAMT.

The SEXTAMT factors excluded from Fig. 14 were reported in four or fewer articles and investigated through no more than two research strategies. They can further enrich our understanding of the impact of the requirements and the estimation process, for instance. Nevertheless, they expand our perspectives to other directions as well, such as the impact of product characteristics, client and user issues, environment, attitudes, and maturity, and testing and rework.

In any case, software organizations and practitioners aiming to diagnose the factors more relevant to their context to improve their estimation results can use the SEXTAMT factors to guide what to include in internal surveys, for instance. Practitioners can

also use the SEXTAMT factors (especially those classified as value adjusting characteristics) to build internal checklists. For instance, Usman et al. (2018b) proposed a process to build checklists to support expert judgment estimation, and the first step is to understand the estimation context. This step has the objective to elicit the factors that should be included in the checklists by, for instance, surveying the literature on the search for effort or cost drivers. The SEXTAMT already provide a map of such factors, and practitioners can save time by using it instead of surveying the literature themselves – a process that involves high costs.

In addition, some of the SEXTAMT factors can be helpful in the debiasing strategy that Kahneman et al. (2021c) proposed to help improve judgments in general: decision observers, i.e., people in charge of observing others making judgments in real-time to identify and alert on the occurrence of biases. Decision observers use checklists to accomplish their tasks, which should be adapted to their specific domain. The SEXTAMT factors can guide such adaptation to the software estimation domain. Particularly, the factors from the bias and the estimation process seas at SEXTAMT can provide valuable items.

Also, practitioners can use the SEXTAMT factors as input for risk analysis for their projects, improving their project planning, monitoring, and control. For instance, projects planned to deliver more extensive or more complex products, with less experienced software teams, or where estimators cannot anticipate the participants' skills when estimating run a larger risk of estimating error and, therefore, of failing to meet their commitments. Thus, project managers of such projects need to be especially caring for monitoring these factors.

Takeaway message 1: There is solid evidence for the factors in the SEXTAMT, with 40 of them reported in three or more articles. A few of those — six in total — were investigated by applying at least three different research strategies. The remaining 29 factors were reported in two studies each, suggesting they can benefit from further investigation.

Takeaway message 2: Practitioners can use the SEXTAMT factors (i) to help diagnose the factors more relevant to change in their contexts, in software process improvement initiatives; (ii) to build supporting checklists for their estimation activities when using expert judgment; (iii) to improve their estimation results in real-time as part of debiasing interventions; or (iv) as input to risk analysis of software projects.

Takeaway message 3: Practitioners interested in improving their estimation can rely on the existing evidence that points to the need for improving the requirements engineering and the estimation process, but also indicates the necessity of considering factors associated with political issues, the management process, experience, team issues, biases, and technical skills.

6.2. Looking through the lenses of the temporal and stakeholder dimensions

When it comes to the process phases in which factors cluster, the planning and executing phases are the ones that stand out. It is natural to have factors at the **planning phase** because estimating occurs primarily during such stage. At the **executing phase**, factors emerge because the dynamics of projects impact estimating error and accuracy. For instance, our software projects have a moving target (Magazinius et al., 2012), and we found in our SLM that *changes to requirements or scope* are an error factor, especially if the original estimates are not modified to reflect the changes. *Overlooked and unplanned tasks* may also be revealed by project execution dynamics, leading to a higher need for effort, costs, and duration than expected.

It is noticeable that only one factor emerged at the **initiating phase** and none at the **closing phase**. However, when looking for the factors reported in one article only, we can find more about such phases. For instance, *bidding situations* are relevant at the **initiating phase**, with one field experiment reporting that companies selected on the criteria of the low bid have higher cost overruns, a phenomenon known as the “winner’s curse” (Jørgensen and Grimstad, 2005). Therefore, estimators might need to pay special attention to the initiating phase in bidding contexts.

Additionally, more investigation on learning and feedback has the potential to shed some light on what is relevant at the **closing phase**. For instance, at least four studies (Jørgensen and Molokken-Ostfold, 2004; Matos et al., 2013; Jørgensen et al., 2007; Jørgensen and Gruschke, 2009) suggest that estimation error, feedback, and learning from past projects and tasks might be beneficial to reducing overconfidence and improving estimates.

Regarding stakeholders, many of the factors are related to **estimators**, which is expected once they are the primarily responsible people for estimates. Our results also indicate the power of other roles that might not be directly involved with the estimating process, such as the client and managers.

Takeaway message 4: Most factors cluster at the planning phase, because estimating occurs primarily at this stage. Many factors also pertain to the execution phase because project dynamics can alter the assumptions on which estimates were generated.

Takeaway message 5: The initiating and closing phases are less explored, and we can benefit from investigating more factors regarding such phases.

Takeaway message 6: Many factors are related to estimators, and many others indicate the power that people playing other roles also have over the estimates, showing that improvement initiatives in the industry must account for them too.

6.3. The strategies researchers employed to explore the seas

As for the project variables, most studies focused on **effort**, which is understandable — as (McConnell, 2006a) suggested by his flow of well-estimated projects that the effort is an intermediary estimate in software projects, ideally used as input to cost and duration estimates. Therefore, factors that impact effort estimates indirectly impact both cost and duration, and because of that, researchers may consider it more beneficial to focus on them.

The mechanism for measuring the impact of the factors that researchers applied the most is rather indirect: the **participants' perceptions** of reasons for errors and accuracy. Such an approach may provide rich insights into the phenomena that cause errors when estimating or promote accuracy in field settings. Considering that many participants in respondents and field studies in our SLM are experts in software development and maintenance tasks, we cannot overlook their opinions about the factors affecting estimates. Nevertheless, the approach has drawbacks also. For instance, people may attribute different meanings to the term “estimate”, even when they work at the same company (Jørgensen, 2014a), making it difficult to interpret the results of surveys (Jørgensen, 2007a).

Another widely employed mechanism for measuring the impact of factors over the estimates was the **difference of estimates** between groups. The difference of estimates does not provide direct evidence about accuracy, but it can evidence when a factor causes an estimate to increase or decrease for reasons beyond the estimation process. This allows us to identify factors that can induce optimism in estimators, leading them to provide low estimates instead of realistic ones. Considering that extensive projects tend to be underestimated with a median time overrun of 20% (Halkjelsvik and Jørgensen, 2018a), identifying such factors can be very useful.

Additionally, researchers have used objective error measures, such as **MRE**, **MREBias**, **BRE**, and **BREBias**. Nevertheless, since the 90's at least, MRE has been criticized because it has the disadvantage of weighing differently under and overestimations. Underestimations are not weighted sufficiently, leading to higher penalization of overestimations (Jørgensen, 2007a). MREBias suffer from this same problem. BRE and BREBias are balanced metrics in this sense (Molokken-Ostfold and Jørgensen, 2005). In Fig. 15, we grouped MRE and MREBias under the label “Unbalanced” and BRE and BREBias under “Balanced”. It shows that, gradually, researchers are moving to the use of more balanced metrics over the years.

Also, researchers prefer accuracy metrics over bias: with 19 occurrences for MRE and BRE together versus 15 occurrences of MREBias and BREBias. Accuracy is the average unsigned error, irrespective of whether the estimate is too high or too low; bias is the average tendency to generate too high or too low estimates (Halkjelsvik and Jørgensen, 2018b).

In any case, using MRE or BRE and similar metrics can be misleading because they depend on actual values, and work can be adjusted to fit an initial estimate (Jørgensen and Sjøberg, 2001), leading to a “moving target problem” (Jørgensen, 2007a) and to a distorted perception of accuracy. For instance, this makes

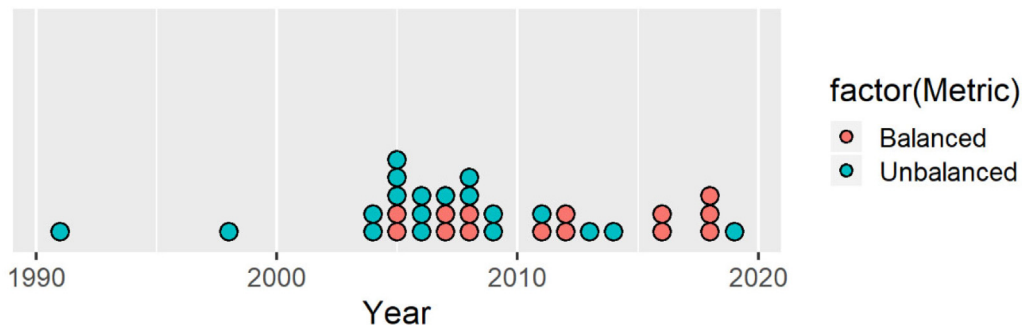


Fig. 15. Balanced (BRE & BREBias) x unbalanced (MRE & MREBias) over the Years.

it harder to understand exactly whether a factor contributed effectively to improving estimation accuracy, or whether a software team just took advantage of a higher project flexibility to create an illusion of accuracy. A possible solution comes from the literature about judgment in general: the measurement of noise instead of bias or accuracy. Noise is the random scatter of judgments that should ideally be identical – or in other words, unwanted variability, a significant component contributing to judgment error, along with bias (Kahneman et al., 2021d). The advantage of measuring noise over bias or accuracy is that we do not need to know actual values. One issue that emerges from this discussion is how to measure noise. A common measure from statistics is the standard deviation (Kahneman et al., 2021f).

Nevertheless, we found very few studies discussing the variability of estimates in the software domain. Only one study explores explicitly the issue, showing a high level of inconsistency when software practitioners estimate the same task, based on the same information and under the same conditions, but at different times (Grimstad and Jørgensen, 2007). In addition, very few studies in our SLM report the standard deviation of estimates, when using the difference of estimates as a measurement strategy (see Shepperd et al., 2018; Passing and Shepperd, 2003). This reveals a low awareness of researchers in our community regarding noise, its relationship with error in expert judgment estimation, and the benefits of measuring and reducing it. Regarding software estimation practice, it is unclear whether practitioners share the perspective of researchers about this concept. In any case, software organizations can benefit from investigating how much disagreement there is among their professionals estimating the same tasks independently.

Regarding research strategies, researchers employed the **laboratory research strategy** widely, and the **respondents' strategy** was quite popular too. Laboratory research strategies favor the investigation of only a few factors at once. In contrast, the articles employing respondents strategies tended to reveal much more factors in each study, contributing significantly to the wide variety of factors we found. The factors with more articles using a laboratory experiment strategy were also the ones that researchers refined the most by investigating relevant variations. For instance, researchers investigated different nuances of the *anchoring effect*, assessing the impact of both numerical and textual anchors (Jørgensen and Grimstad, 2012), as well as of single and interval anchors (Løhre and Jørgensen, 2016). Another refinement was the investigation of the moderating effect of the expertise of the source and of the receiver of the anchor value (Løhre and Jørgensen, 2016) and the impact of one intervention to reduce its effects (Shepperd et al., 2018). Another example is the *sequence*

effect, whose impact over the estimates varies with the size of the tasks estimates in the sequence (Jørgensen and Halkjelsvik, 2020). Researchers perceived an assimilation effect (the estimate become more similar to the one of a previously estimated task) for tasks of different sizes, and a contrast effect (the estimate become more different than the previous one) for tasks of similar sizes.

When considering the taxonomy of Stol and Fitzgerald (2018) for research strategies, it is interesting to notice that the studies employing the field strategy, there are very few field experiments – a total of 10. In other words, when it comes to factors affecting estimates, researchers are more likely to enter natural settings to collect data without manipulating variables. Probably such manipulations are hard to be approved by administrative staff or to be adequately carried out. Thus, they restrict the manipulations of variables to the lab, reinforcing the need for triangulation of strategies (Stol and Fitzgerald, 2018) to evaluate further the impact of factors investigated.

Additionally, considering that the potential for generalizability from respondent studies and the potential for realism from field studies can be taken as proxies of the relevance of research results for practice, from all the 69 factors from the SEXTAMT, most (62) have this type of evidence. From the seven factors with no evidence from respondent or field studies, three are related to biases on estimation and were investigated through lab studies only: *sequence effects*, *time frame size*, and *unit effects*. The *client's expectation* was a factor investigated only through lab studies. The *programming language*, *business area*, and *longer projects* emerged from data studies only. Nevertheless, the lack of evidence from respondents and field studies for these factors does not mean they are irrelevant. For instance, practitioners are not aware of the biases affecting them in many cases, which makes it hard for them to identify this kind of factor in respondent studies. Therefore, combining research strategies reveals complementary findings in research topics so complex as this one. This has been highlighted before in the study of reasons for software effort estimation error in one single company: combining information sources, data collection methods, and data analysis methods leads to complementary insights (Jørgensen and Molokken-Ostfold, 2004).

Takeaway message 7: The participants' perceptions can provide a rich picture of factors affecting estimates in practice, even though it provides a subjective perspective. For more objective measurements of impact, the difference of estimates between a control and an experimental group has been largely adopted.

Takeaway message 8: Despite the criticism over metrics such as MRE, researchers are still gradually moving to use more balanced metrics such as BRE to assess the accuracy of estimates.

Takeaway message 9: Researchers are not fully aware of the concept of noise and its contribution to estimation error, even though it can reveal estimation problems with the benefit that we do not need to know actual values to measure it. It is not clear whether practitioners are unaware of it as well. In any case, software organizations can benefit from noise audits as starting points to improvement initiatives and noise measurements to assess the effectiveness of interventions to their estimation processes.

Takeaway message 10: Respondents strategies allowed for discovering many factors relevant in practice, while laboratory strategies allowed for the refinement of factors.

Takeaway message 11: The combination of different research strategies provides complementary factors, allowing for a richer map of the factors affecting expert judgment estimates.

6.4. Into the wild – part 1: underexplored seas

We excluded from the SEXTAMT a total of 166 factors reported in one research article only each.⁸ Therefore, we consider they are in a gray area, and there is a need to execute more research to strengthen the evidence about their impact. Some of them have the potential to enlarge the territory of existing seas in the SEXTAMT. In contrast, others have the potential to reveal new seas of their own.

Such a myriad of factors does not mean that all factors reported by unique articles are worthy of further investigation. We need some filtering on them to decide which ones are good candidates for more studies. For instance, *luck* is a factor reported in a respondents study. However, what does it mean? Also, the presence of other factors we identified in our SLM might explain luck to some extent: we can consider that a software project was luckier because requirements did not change, for example.

In addition, we classified some of these factors as a satellite to others, meaning they are somewhat related, even though not enough to be united to create a final factor to include at the SEXTAMT. One example is “team process experience” and “expertise of new team members”. Although they are related, the first can relate to all team members, not only to new ones, while the latter is very specific in including only new people. Therefore, we cannot unify them to form a single factor investigated in three articles, allowing its inclusion among the SEXTAMT factors. Therefore, we kept them as part of the unique factors, marking them as satellites of each other.

Another example is the case of the factor *forcing to stay within the estimate*. It is a satellite of one SEXTAMT factor: *project flexibility*. For instance, software practitioners need the flexibility to deliver less polished features when they are forced to stay within a deadline, no matter what. We can also argue that *forcing to stay within the estimate* is a repercussion of other factors from unique articles, such as *estimates interpreted as commitments* or the *use of uncertain estimates as baselines*. Nevertheless, researchers have not validated such relationships. In any case, we indicate the satellite factors as part of our supplementary material.

Factors investigated through laboratory research strategies are good candidates for field experiments to assess whether their impact is kept in real-life contexts. Take the *format* factor, for example, which is about using the traditional request format – “How much effort is required to complete X?” – versus using an alternative format – “How much can be completed in Y work-hours?”. In the laboratory, the alternative format has led to more

optimistic estimates (Jørgensen and Halkjelsvik, 2010). However, it is precisely the format we expect when using agile methodologies. Does it impact estimates negatively in the trenches, making them more optimistic? Another factor whose effect is relevant in the same context is the *use of Fibonacci scales* that, compared to linear scales, led to lower estimates when using Tamrakar and Jørgensen (2012).

Takeaway message 12: Researchers have investigated a large and varied set of factors affecting estimates when using expert judgment. Most of such factors were reported in one article only, needing more research to strengthen the evidence about their impact.

6.5. Into the wild – part 2: validated relationships among the factors

The discussion of satellite factors leads us to another underexplored issue: the relationship among different factors. Therefore, after answering the main research questions that we presented in Section 3.1, we decided to extract and analyze data for an additional question: “SQ 1.6 – What are the validated relationships among the factors affecting expert judgment expert estimates?”

Only nine articles had results regarding such relationships. We illustrate the relationships we found in Fig. 16, where each light blue rounded rectangle represents one SEXTAMT factor. Each gray rounded rectangle represents one factor we did not include in the SEXTAMT because it was investigated in only one article.

Fig. 16 shows that *overall experience* moderates the impact of the *combination strategy of individual estimates*: more experience is connected with less optimistic estimates when using Planning Poker compared to when using a statistical combination of estimates. In the context where estimators have less experience, the result is inverted: the statistical combination leads to less optimism (Mahnič and Hovelja, 2012). This result suggests that without experience discussions lack the benefit of meaningful divergent perspectives about the task complexity, or the wisdom to recognize forgotten tasks, or other flaws in judgment. It also seems that higher experience is needed to overcome the effect of the social influence bias (Lorenz et al., 2011) in group discussions of the estimates. Nevertheless, we should consider this result carefully because researchers contrasted a sample of students (representing less experience) with a sample of software professionals (representing more experience).

Overall experience also reduces the impact of the *anchoring effect* over the estimates (Løhre and Jørgensen, 2016), as well as *debiasing workshops* (Shepperd et al., 2018) and the use of subsequent anchors aimed at neutralizing first impressions caused by the first anchor (Jørgensen and Løhre, 2012). However, in none of these studies, the anchoring effect was completely removed: only reduced. In another study, the researchers showed that mixed-handers were more strongly influenced by anchors compared with strong-handers (Jørgensen, 2007c), revealing how *handedness* can influence estimates.

Handedness also impacts the *effect of more and/or irrelevant information*: mixed-handers also are more impacted by irrelevant information (Jørgensen, 2007c). In addition, people who score high in *interdependence* are also more strongly influenced by more and/or irrelevant information than people who score low (Jørgensen and Grimstad, 2012). Higher interdependence refers to higher connectedness to others and higher importance to social context and relationships. In addition, another study showed that higher technical skill reduces the impact of more and/or irrelevant information (Jørgensen and Grimstad, 2008).

Also, the *time frame size* moderates the impact caused by using an alternative *format* for requesting estimates. Smaller time frames increase the impact by leading to more optimistic estimates (Halkjelsvik and Jørgensen, 2011; Jørgensen and Halkjelsvik, 2010).

⁸ A complete list of such factors, together with their codes and categories can be found in our supplementary material (Matsubara et al., 2021).

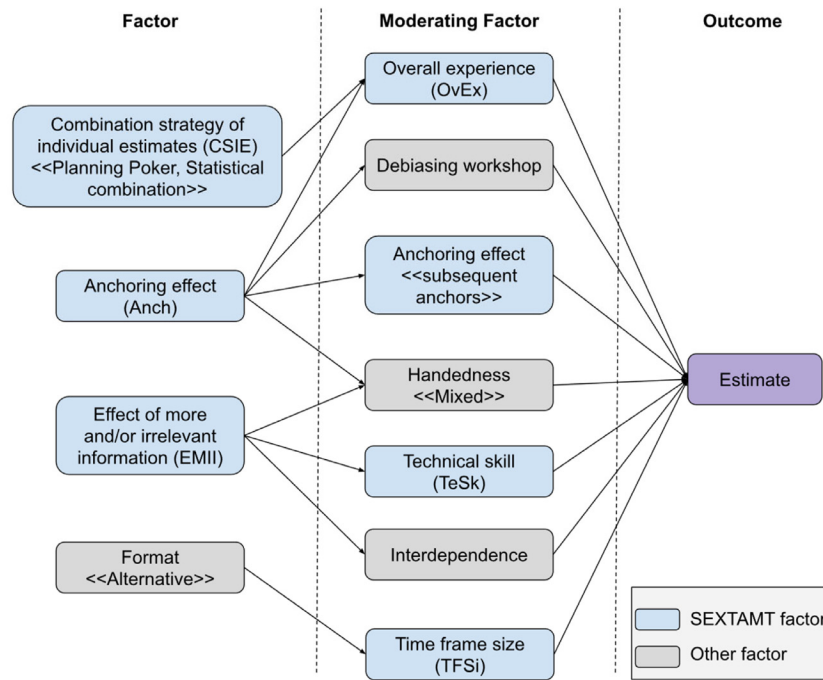


Fig. 16. Relationships among factors.. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

Interestingly, many of the factors on the leftmost side of Fig. 16 are related to a psychological or social bias. For instance, combination strategies of individual estimates are subject to social biases. The anchoring effect is a psychological bias. The presence of more and/or irrelevant information can also bias judgment, leading people to think the task is larger than it truly is, for example. Additionally, the moderating factors give us hints about interventions to deal with such biases. For instance, if we know estimators have low experience, it might be wiser to use the statistical combination of estimates instead of Planning Poker. Another example is composing estimating teams to include people with higher experience and technical skills whenever possible because this helps reduce the effects of psychological biases.

These results show the relevance of studying the relationships among factors. Such relationships can reveal the paths of interaction among factors and the ones that can trigger chains of negative or positive effects over the estimates. Therefore, in software process improvement initiatives regarding the estimation process, focusing in factors mediating or moderating others can be a cost-effective strategy to improve accuracy.

Takeaway message 13: Although the set of factors affecting estimates when using expert judgment is large and varied, we could benefit from more studies exploring the relationships among such factors. This investigation can help narrow down the factors to focus to keep a good balance of costs and benefits when dealing with estimation problems.

7. Threats to validity

We analyzed the validity threats to this SLM, considering threats to the study selection validity, threats to data validity, and threats to research validity (Ampatzoglou et al., 2019). One of the threats for study selection validity is the adequacy of initial relevant publications identification, addressed with an automatic search in known digital libraries. Another mitigation action to this threat was the use of a known set of papers to evaluate the search strategy (Zhang et al., 2011). The goal of this evaluation was to reach a sensitivity of 70% in automated search (Zhang et al., 2011). A final mitigation action to this threat was snowballing

procedures to enlarge the number of retrieved relevant papers, reaching a sensitivity of 100% afterward. Another threat to study selection validity for this SLM is the study inclusion/exclusion bias, addressed through the definition of study inclusion and exclusion criteria in the research protocol. Additionally, the authors executed the selection process over a sample of the articles, discussing any inclusion or exclusion conflicts. Their agreement level was measured with the kappa statistic, leading to the refinement of the inclusion and exclusion criteria.

A threat to data validity in this SLM is the data extraction bias, addressed through a pilot data extraction. The authors reviewed and discussed a pilot data extraction sample to improve the data extraction form. Another threat is the bias of classification schema. To avoid it, we relied on previous existing classifications when possible, such as the research strategies framework of Storey et al. (2020). We used the process groups from PM-BOK (Project Management Institute, 2017a) for the phases and familiar stakeholders' roles regarding the factors. We aggregated similar findings under labels that reflected the articles' original texts for naming the factors affecting software estimates. The authors held meetings for reviewing the factors and the categories in the SEXTAMT, and the types of effects of each factor.

As for research validity, there is the threat of lack of repeatability. One of the mitigation actions for this threat was involving more than one researcher during the process. Another action is to make all the SLM data publicly available, including decisions about inclusion and exclusion of papers, extracted data from primary studies, among others. Finally, we developed a research protocol to ensure replications or updates to this SLM. The protocol we developed and the discussions among the researchers involved helped mitigate the research method bias, another threat to research validity.

8. Conclusion

In this article, we presented an SLM about factors affecting expert judgment software estimates. We present such factors by three dimensions: the project phase they are likely to happen or to cause an impact over the estimates; the stakeholder that is

responsible for a task or process to which the factor is linked, that directly causes the factor or that is directly impacted by the effects of the factor; and type of effect the factor causes. Some factors can have a negative effect, leading to errors when they are present, while others may have a positive or neutral effect. Such dimensions allow for easier navigation through the myriad of factors we found.

Most of the factors clustered at the planning and executing phases. It is natural to have factors at the planning phase because estimating occurs primarily during such stage. At the executing phase, factors emerge because the dynamics of projects impact estimating error and accuracy. Moreover, most of the studies employed a research strategy of laboratory experiments, investigating one factor in a controlled setting with an experimental and control group. Also, they evaluated the difference of estimates between these groups to assess the impact of the factors.

Top factors – those that emerged in a higher number of studies – revealed the importance of issues beyond the estimation process. It is also necessary to improve the requirements engineering process, to deal with political issues, to consider the product characteristics, among others. Researchers have investigated a wide and varied set of factors. Therefore, we created a map to support readers in navigation through them: the SEXTAMT. If an interested reader desires to identify all factors that affect only one project phase, we provide them a classification through this dimension. If the reader desires to identify all factors given one stakeholder, we also provide this. Finally, if the reader wants to find out a class of factors given a specific effect – for instance, all factors that lead to improved accuracy – our map also has a dimension regarding this.

Our research confirms and aggregates existing results about factors affecting expert judgment estimates, a relevant contribution to move knowledge forward, especially when we organize such knowledge to facilitate understanding and future uses (for both research and practice). Also, the classification of measurement strategies is an additional relevant contribution. This enabled us to spot that our research community is missing the benefits of investigating more of noise as component of error.

The SEXTAMT can have many valuable uses in practice and software practitioners can employ its factors as part of many different initiatives, such as:

- Diagnosing improvement opportunities to their estimation processes through the investigation of the most relevant factors in their contexts, considering their types of effects;
- building checklists to support estimators, considering especially the value adjustment characteristics;
- adapting checklists to aid debiasing interventions, considering especially the factors from the bias and estimation process categories;
- analyzing project risks by identifying the factors leading to larger risks of estimating error in their contexts and, therefore, of leading to failures to meet their commitments.

As for future work, we need to keep the SEXTAMT updated. Special care is due to the factors coming from unique articles: more investigation about them is needed. However, some filtering to identify the best candidates for more assessment is also necessary. Another critical issue is investigating the relationships among the factors to enrich the map with relevant mediation and moderation connections. A more complex framework can be helpful to identify the factors more likely to cause a more considerable impact over the estimates, focusing on them to adopt more cost-effective interventions during software improvement initiatives regarding estimation processes.

We highlight that another research issue comes from the software project dynamics that allows practitioners to adjust

their work to fit an estimate when they need to, creating a “moving target” problem. This makes it harder to measure error and accuracy correctly. It also makes it harder to understand whether a factor contributed effectively to improving estimation accuracy or whether a software team just took advantage of higher project flexibility to create an illusion of accuracy. The solution to this comes from the judgment literature: measuring noise – unwanted variability from judgments that ideally should be identical (Kahneman et al., 2021d). Nevertheless, few studies from our SLM discuss this issue, revealing that our research community can benefit from understanding and using more of this concept.

CRedit authorship contribution statement

Patrícia Gomes Fernandes Matsubara: Conceptualization, Methodology, Investigation, Formal Analysis, Writing – original draft, Visualization. **Bruno Freitas Gadelha:** Conceptualization, Methodology, Formal Analysis, Validation, Writing – review & editing, Supervision. **Igor Steinmacher:** Formal Analysis, Validation, Writing – review & editing. **Tayana Uchôa Conte:** Conceptualization, Methodology, Formal Analysis, Validation, Writing – review & editing, Supervision, Funding Acquisition.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgments

This research, carried out within the scope of the Samsung-UFAM Project for Education and Research (SUPER), according to Article 48 of Decree no 6.008/2006 (SUFRAMA), was funded by Samsung Electronics of Amazonia Ltda., under the terms of Federal Law no 8.387/1991, through agreement 001/2020, signed with Federal University of Amazonas and FAEPI, Brazil and through agreement no 003/2019 (PROPPGI), signed with ICOMP/UFAM. Also supported by Universidade Federal do Amazonas (UFAM) and Coordenação de Aperfeiçoamento de Pessoal de Nível Superior (CAPES) - Financing Code 001, CNPq processes 314174/2020-6 and 313067/2020-1, FAPEAM process 062.00150/2020, and grant #2020/05191-2 São Paulo Research Foundation (FAPESP).

Appendix A. Supplementary data

Supplementary material related to this article can be found online at <https://doi.org/10.1016/j.jss.2021.111148>.

References

- Altaieb, A., Altherwi, M., Gravell, A., 2020a. An industrial investigation into effort estimation predictors for mobile app development in agile processes. In: 2020 9th International Conference on Industrial Technology and Management (ICITM). pp. 291–296. <http://dx.doi.org/10.1109/ICITM48982.2020.9080362>.
- Altaieb, Abdullah, Altherwi, Muna, Gravell, Andy, 2020b. A pair estimation technique of effort estimation in mobile app development for agile process. In: Proceedings of the 2020 the 3rd International Conference on Information Science and System. Retrieved November 26, 2020 from <https://dl.acm.org/doi/abs/10.1145/3388176.3388212>.
- Altaieb, Abdullah, Gravell, Andrew, 2019. An empirical investigation of effort estimation in mobile apps using agile development process. *J. Softw.* 14 (8), 356–369.
- Ampatzoglou, Apostolos, Bibi, Stamatia, Avgeriou, Paris, Verbeek, Marijn, Chatzigeorgiou, Alexander, 2019. Identifying, categorizing and mitigating threats to validity in software engineering secondary studies. *Inf. Softw. Technol.* 106, 201–230. <http://dx.doi.org/10.1016/j.infsof.2018.10.006>.

- Aranda, Jorge, Easterbrook, Steve, 2005. Anchoring and adjustment in software estimation. In: Proceedings of the 10th European Software Engineering Conference Held Jointly with 13th ACM SIGSOFT International Symposium on Foundations of Software Engineering (ESEC/FSE-13). ACM, New York, NY, USA, pp. 346–355. <http://dx.doi.org/10.1145/1081706.1081761>.
- Arifin, H.H., Daengdej, J., Khanh, N.T., 2017. An empirical study of effort-size and effort-time in expert-based estimations. In: 2017 8th International Workshop on Empirical Software Engineering in Practice (IWSEEP). pp. 35–40. <http://dx.doi.org/10.1109/IWSEEP.2017.21>.
- Arnuphaptrairong, Tharwon, 2018. The state of practice of software cost estimation: evidence from thai software firms. In: Proceedings of the International MultiConference of Engineers and Computer Scientists. Hong Kong, p. 6.
- Atas, Muesluem, Reiterer, Stefan, Felfernig, Alexander, Trang Tran, Thi Ngoc, Stettinger, Martin, 2018. Polarization effects in group decisions. In: Adjunct Publication of the 26th Conference on User Modeling, Adaptation and Personalization (UMAP '18). Association for Computing Machinery, New York, NY, USA, pp. 305–310. <http://dx.doi.org/10.1145/3213586.3225242>.
- Basten, Dirk, Mellis, Werner, 2011. A current assessment of software development effort estimation. In: 2011 International Symposium on Empirical Software Engineering and Measurement. IEEE, Banff, AB, Canada, pp. 235–244. <http://dx.doi.org/10.1109/ESEM.2011.32>.
- Basten, Dirk, Sunyaev, Ali, 2014. A systematic mapping of factors affecting accuracy of software development effort estimation. Commun. Assoc. Inf. Syst. 34, 4. <http://dx.doi.org/10.17705/1CAIS.03404>.
- Benschop, Nick, Hilhorst, Cokky A.R., Nuijten, Arno L.P., Keil, Mark, 2020. Detection of early warning signals for overruns in IS projects: linguistic analysis of business case language. Eur. J. Inf. Syst. 29 (2), 190–202. <http://dx.doi.org/10.1080/0960085X.2020.1742587>.
- Bergeron, François, St-Arnaud, Jean-Yves, 1992. Estimation of information systems development efforts: A pilot study. Inf. Manage. 22 (4), 239–254. [http://dx.doi.org/10.1016/0378-7206\(92\)90026-C](http://dx.doi.org/10.1016/0378-7206(92)90026-C).
- Bhatt, P., Shroff, G., K, W., Misra, A.K., 2006. An empirical study of factors and their relationships in outsourced software maintenance. In: 2006 13th Asia Pacific Software Engineering Conference (APSEC'06). pp. 301–308. <http://dx.doi.org/10.1109/APSEC.2006.21>.
- Boehm, Barry W., 1984. Software engineering economics. IEEE Trans. Softw. Eng. SE-10 (1), 4–21. <http://dx.doi.org/10.1109/TSE.1984.5010193>.
- Boetticher, Gary, Lokhandwala, Nazim, 2007. Assessing the Reliability of a Human Estimator. USA. Retrieved January 7, 2020 from <https://dl.acm.org/doi/pdf/10.1109/PROMISE.2007.2>.
- Branco, Daniel Tadeu Martínez C., Oliveira, Edson Cesar Cunhade, Galvão, Leandro, Prikladnicki, Rafael, Conte, Tayana, 2015. An empirical study about the influence of project manager personality in software project effort. In: Proceedings of the 17th International Conference on Enterprise Information Systems - Volume 2 (ICEIS 2015). SCITEPRESS - Science and Technology Publications, Lda, Portugal, pp. 102–113. <http://dx.doi.org/10.5220/0005373001020113>.
- Bratthall, Lars, Arisholm, Erik, Jørgensen, Magne, 2001. Program understanding behavior during estimation of enhancement effort on small Java programs | SpringerLink. Retrieved February 11, 2021 from https://link.springer.com/chapter/10.1007/3-540-44813-6_30.
- Britto, Ricardo, Mendes, Emilia, Börstler, Jürgen, 2015. An empirical investigation on effort estimation in agile global software development. In: 2015 IEEE 10th International Conference on Global Software Engineering. IEEE, pp. 38–45. <http://dx.doi.org/10.1109/ICGSE.2015.10>.
- Brown, M., Dirska, Henry, Pelosi, M., Assadullah, M., 2013. Agile method software development estimation biases. Int. J. Adv. Res. Comput. Sci. Softw. Eng. 3 (2013), 10.
- Bukhari, S., Malik, A.A., 2012. Determining the factors affecting the accuracy of effort estimates for different application and task types. In: 2012 10th International Conference on Frontiers of Information Technology. pp. 41–45. <http://dx.doi.org/10.1109/FIT.2012.16>.
- Cao, Lan, 2008. Estimating Agile Software Project Effort: An Empirical Study. 11.
- Conoscenti, Marco, Besner, Veronika, Vetrò, Antonio, Fernández, Daniel Méndez, 2019. Combining data analytics and developers feedback for identifying reasons of inaccurate estimations in agile software development. J. Syst. Softw. 156 (2019), 126–135. <http://dx.doi.org/10.1016/j.jss.2019.06.075>.
- Davis, J.S., 1989. Investigation of predictors of failures and debugging effort for large MIS. Inf. Softw. Technol. 31 (4), 170–174. [http://dx.doi.org/10.1016/0950-5849\(89\)90034-7](http://dx.doi.org/10.1016/0950-5849(89)90034-7).
- Dyba, T., Dingsoyr, T., Hanssen, G.K., 2007. Applying systematic reviews to diverse study types: An experience report. In: First International Symposium on Empirical Software Engineering and Measurement (ESEM 2007). pp. 225–234. <http://dx.doi.org/10.1109/ESEM.2007.59>.
- Furulund, Kristian Marius, Molkken-stvold, Kjetil, 2007. Seventh International Conference on Quality Software (QSIC 2007). In: Increasing Software Effort Estimation Accuracy Using Experience Data, Estimation Models and Checklists, pp. 342–347. <http://dx.doi.org/10.1109/QSIC.2007.4385518>.
- Gandomani, T.J., Faraji, H., Radnejad, M., 2019. Planning Poker in cost estimation in Agile methods: Averaging vs. Consensus. In: 2019 5th Conference on Knowledge Based Engineering and Innovation (KBEI). pp. 066–071. <http://dx.doi.org/10.1109/KBEI.2019.8734960>.
- Glass, R.L., Rost, J., Matook, M.S., 2008. Lying on software projects. IEEE Softw. 25 (6), 90–95. <http://dx.doi.org/10.1109/MS.2008.150>.
- Grapenthin, S., Book, M., Richter, T., Gruhn, V., 2016. Supporting feature estimation with risk and effort annotations. In: 2016 42th Euromicro Conference on Software Engineering and Advanced Applications (SEAA). pp. 17–24. <http://dx.doi.org/10.1109/SEAA.2016.24>.
- Gray, A.R., MacDonell, S.G., Shepperd, M.J., 1999. Factors systematically associated with errors in subjective estimates of software development effort: the stability of expert judgment. In: Proceedings Sixth International Software Metrics Symposium (Cat. No. PR00403). pp. 216–227. <http://dx.doi.org/10.1109/METRIC.1999.809743>.
- Gren, L., Svensson, R.B., Unterkalmsteiner, M., 2017. Is it possible to disregard obsolete requirements? An initial experiment on a potentially new bias in software effort estimation. In: 2017 IEEE/ACM 10th International Workshop on Cooperative and Human Aspects of Software Engineering (CHASE). pp. 56–61. <http://dx.doi.org/10.1109/CHASE.2017.10>.
- Grimstad, S., Jørgensen, M., 2007. The impact of irrelevant information on estimates of software development effort. In: 2007 Australian Software Engineering Conference (ASWEC'07). pp. 359–368. <http://dx.doi.org/10.1109/ASWEC.2007.48>.
- Grimstad, Stein, Jørgensen, Magne, 2007. Inconsistency of expert judgment-based estimates of software development effort. J. Syst. Softw. 80, 11. <http://dx.doi.org/10.1016/j.jss.2007.03.001>.
- Grimstad, S., Jørgensen, M., 2009. Preliminary study of sequence effects in judgment-based software development work-effort estimation. IET Softw. 3 (5), 435–441. <http://dx.doi.org/10.1049/iet-sen.2008.0110>.
- Grimstad, S., Jørgensen, M., Molokken-Ostvold, K., 2005. The clients' impact on effort estimation accuracy in software development projects. In: 11th IEEE International Software Metrics Symposium (METRICS'05). pp. 10 pp.–10. <http://dx.doi.org/10.1109/METRICS.2005.30>.
- Gruschke, Tanja M., Jørgensen, Magne, 2008. The role of outcome feedback in improving the uncertainty assessment of software development effort estimates. ACM Trans. Softw. Eng. Methodol. 17 (4), 20:1–20:35. <http://dx.doi.org/10.1145/13487689.13487693>.
- Halkjelsvik, Torleif, Jørgensen, Magne, 2011. To read two pages, I need 5 minutes, but give me 5 minutes and I will read four: how to change productivity estimates by inverting the question. Appl. Cogn. Psychol. 25 (2), 314–323.
- Halkjelsvik, Torleif, Jørgensen, Magne, 2012. From origami to software development: A review of studies on judgment-based predictions of performance time. Psychol. Bull. 138 (2), 238–271. <http://dx.doi.org/10.1037/a0025996>.
- Halkjelsvik, Torleif, Jørgensen, Magne, 2018a. Time Predictions: Understanding and Avoiding Unrealism in Project Planning and Everyday Life. Springer International Publishing, <http://dx.doi.org/10.1007/978-3-319-74953-2>.
- Halkjelsvik, Torleif, Jørgensen, Magne, 2018b. How we predict time usage. In: Halkjelsvik, Torleif, Jørgensen, Magne (Eds.), Time Predictions: Understanding and Avoiding Unrealism in Project Planning and Everyday Life. Springer International Publishing, Cham, pp. 5–11. http://dx.doi.org/10.1007/978-3-319-74953-2_2.
- Halstead, Susanne, Ortiz, Rosario, Córdova, Mario, Seguí, Miguel, 2012. The impact of lack in domain or technology experience on the accuracy of expert effort estimates in software projects. In: Proceedings of the 13th International Conference on Product-Focused Software Process Improvement (PROFES'12). Springer-Verlag, Berlin, Heidelberg, pp. 248–259. http://dx.doi.org/10.1007/978-3-642-31063-8_19.
- Haugen, N.C., 2006. An empirical study of using planning poker for user story estimation. In: AGILE 2006 (AGILE'06). pp. 9 pp.–34. <http://dx.doi.org/10.1109/AGILE.2006.16>.
- He, Mei, Zhang, He, Yang, Ye, Wang, Qing, Li, Mingshu, 2010. Understanding the influential factors to development effort in Chinese software industry. In: Product-Focused Software Process Improvement. In: Lecture Notes in Computer Science, Springer, Berlin, Heidelberg, pp. 306–320. http://dx.doi.org/10.1007/978-3-642-13792-1_24.
- Henry, Raymond M., McCray, Gordon E., Purvis, Russell L., Roberts, Tom L., 2007. Exploiting organizational knowledge in developing IS project cost and schedule estimates: An empirical study. Inf. Manage. 44 (6), 598–612. <http://dx.doi.org/10.1016/j.im.2007.06.002>.
- Hill, J., Thomas, L.C., Allen, D.E., 2000. Experts' estimates of task durations in software development projects. Int. J. Proj. Manage. 18 (1), 13–21. [http://dx.doi.org/10.1016/S0263-7863\(98\)00062-3](http://dx.doi.org/10.1016/S0263-7863(98)00062-3).
- Host, M., Wohlin, C., 1998. An experimental study of individual subjective effort estimations and combinations of the estimates. In: Proceedings of the 20th International Conference on Software Engineering. pp. 332–339. <http://dx.doi.org/10.1109/ICSE.1998.671386>.
- Huang, J., Sun, H., Li, Y., 2015. An empirical study of the impact of project factors on software economics. In: 2015 IEEE International Conference on Industrial Engineering and Engineering Management (IEEM). pp. 43–47. <http://dx.doi.org/10.1109/IEEM.2015.7385605>.

- Hugh, Chrisholm, 1911. *The Encyclopaedia Britannica - Dictionary of Arts, Sciences, Literature and General Information*, eleventh ed.
- IEEE, 2017b. ISO/IEC/IEEE International Standard - Systems and Software Engineering-Vocabulary. ISO/IEC/IEEE 24765:2017(E), pp. 1–541. <http://dx.doi.org/10.1109/IEEESTD.2017.8016712>.
- Javed, Ali, Ullah, M.A., Aziz-Ur, Rehman, 2013. Factors affecting software cost estimation in developing countries. *Int. J. Inf. Technol. Comput. Sci.* 5, 54–59.
- Jørgensen, M., 2004. Realism in assessment of effort estimation uncertainty: it matters how you ask. *IEEE Trans. Softw. Eng.* 30 (4), 209–217. <http://dx.doi.org/10.1109/TSE.2004.1274041>.
- Jørgensen, Magne, 2007a. A critique of how we measure and interpret the accuracy of software development effort estimation. In: *First International Workshop on Software Productivity Analysis and Cost Estimation*. Japan.
- Jørgensen, Magne, 2007b. Forecasting of software development work effort: Evidence on expert judgement and formal models. *Int. J. Forecast.* 23 (3), 449–462. <http://dx.doi.org/10.1016/j.ijforecast.2007.05.008>.
- Jørgensen, Magne, 2007c. Individual differences in how much people are affected by irrelevant and misleading information. In: *Hellenic Cognitive Science Society*. pp. 347–352.
- Jørgensen, Magne, 2010a. Identification of more risks can lead to increased over-optimism of and over-confidence in software development effort estimates. *Inf. Softw. Technol.* 52 (5), 506–516. <http://dx.doi.org/10.1016/j.infsof.2009.12.002>.
- Jørgensen, Magne, 2010b. Selection of strategies in judgment-based effort estimation - *ScienceDirect*. *J. Syst. Softw.* 83 (6), 1039–1050.
- Jørgensen, Magne, 2011. Contrasting ideal and realistic conditions as a means to improve judgment-based software development effort estimation. *Inf. Softw. Technol.* 53 (12), 1382–1390. <http://dx.doi.org/10.1016/j.infsof.2011.07.001>.
- Jørgensen, Magne, 2013a. The influence of selection bias on effort overruns in software development projects. *Inf. Softw. Technol.* 55 (9), 1640–1650. <http://dx.doi.org/10.1016/j.infsof.2013.03.001>.
- Jørgensen, M., 2013b. Relative estimation of software development effort: It matters with what and how you compare. *IEEE Softw.* 30 (2), 74–79. <http://dx.doi.org/10.1109/MS.2012.70>.
- Jørgensen, Magne, 2014. What we do and don't know about software development effort estimation. *IEEE Softw.* 31 (2), 37–40. <http://dx.doi.org/10.1109/MS.2014.49>.
- Jørgensen, Magne, 2014a. Communication of software cost estimates. In: *Proceedings of the 18th International Conference on Evaluation and Assessment in Software Engineering (EASE '14)*. ACM, New York, NY, USA, pp. 28:1–28:5. <http://dx.doi.org/10.1145/2601248.2601262>.
- Jørgensen, M., 2014b. Fallacies and biases when adding effort estimates. In: *2014 40th EUROMICRO Conference on Software Engineering and Advanced Applications*. pp. 277–284. <http://dx.doi.org/10.1109/SEAA.2014.16>.
- Jørgensen, Magne, 2014c. The ignorance of confidence levels in minimum-maximum software development effort intervals. *LNSE 2* (4), 327–330. <http://dx.doi.org/10.7763/LNSE.2014.V2.144>.
- Jørgensen, M., 2015. The effect of the time unit on software development effort estimates. In: *2015 9th International Conference on Software, Knowledge, Information Management and Applications (SKIMA)*. pp. 1–5. <http://dx.doi.org/10.1109/SKIMA.2015.7399992>.
- Jørgensen, Magne, 2016a. Unit effects in software project effort estimation: Work-hours gives lower effort estimates than workdays. *J. Syst. Softw.* 117, 274–281. <http://dx.doi.org/10.1016/j.jss.2016.03.048>.
- Jørgensen, Magne, 2016b. The use of precision of software development effort estimates to communicate uncertainty. In: *Software Quality, the Future of Systems and Software Development*. In: *Lecture Notes in Business Information Processing*, Springer International Publishing, Cham, pp. 156–168. http://dx.doi.org/10.1007/978-3-319-27033-3_11.
- Jørgensen, M., 2018. Looking back on previous estimation error as a method to improve the uncertainty assessment of benefits and costs of software development projects. In: *2018 9th International Workshop on Empirical Software Engineering in Practice (IWESep)*. pp. 19–24. <http://dx.doi.org/10.1109/IWESep.2018.00012>.
- Jørgensen, M., Bergersen, G.R., Liestol, K., 2020. Relations between effort estimates, skill indicators, and measured programming skill. *IEEE Trans. Softw. Eng.* 1. <http://dx.doi.org/10.1109/TSE.2020.2973638>.
- Jørgensen, Magne, Carelius, Gunnar J., 2004. An empirical study of software project bidding. *IEEE Trans. Softw. Eng.* 30 (12), 953–969. <http://dx.doi.org/10.1109/TSE.2004.92>.
- Jørgensen, Magne, Faugli, Bjørn, Gruschke, Tanja, 2007. Characteristics of software engineers with optimistic predictions. *J. Syst. Softw.* 80 (9), 1472–1482. <http://dx.doi.org/10.1016/j.jss.2006.09.047>.
- Jørgensen, M., Grimstad, S., 2005. Over-optimism in software development projects: the winner's curse. In: *15th International Conference on Electronics, Communications and Computers (CONIELECOMP'05)*. pp. 280–285. <http://dx.doi.org/10.1109/CONIEL.2005.58>.
- Jørgensen, Magne, Grimstad, Stein, 2008. Avoiding irrelevant and misleading information when estimating development effort. *IEEE Softw.* 25 (3), 78–83. <http://dx.doi.org/10.1109/MS.2008.57>.
- Jørgensen, M., Grimstad, S., 2011. The impact of irrelevant and misleading information on software development effort estimates: A randomized controlled field experiment. *IEEE Trans. Softw. Eng.* 37 (5), 695–707. <http://dx.doi.org/10.1109/TSE.2010.78>.
- Jørgensen, M., Grimstad, S., 2012. Software development estimation biases: The role of interdependence. *IEEE Trans. Softw. Eng.* 38 (3), 677–693. <http://dx.doi.org/10.1109/TSE.2011.40>.
- Jørgensen, Magne, Gruschke, Tanja M., 2009. The impact of lessons-learned sessions on effort estimation and uncertainty assessments. *IEEE Trans. Softw. Eng.* 35 (3), 368–383. <http://dx.doi.org/10.1109/TSE.2009.2>.
- Jørgensen, Magne, Halkjelsvik, Torleif, 2010. The effects of request formats on judgment-based effort estimation. *J. Syst. Softw.* 83 (1), 29–36. <http://dx.doi.org/10.1016/j.jss.2009.03.076>.
- Jørgensen, Magne, Halkjelsvik, Torleif, 2020. Sequence effects in the estimation of software development effort. *J. Syst. Softw.* 159, 110448. <http://dx.doi.org/10.1016/j.jss.2019.110448>.
- Jørgensen, Magne, Løhre, Erik, 2012. *First Impressions in Software Development Effort Estimation: Easy to Create and Difficult To Neutralize*. IET, pp. 216–222.
- Jørgensen, Magne, Moløkken, Kjetil, 2002. Combination of software development effort prediction intervals: why, when and how? In: *Proceedings of the 14th International Conference on Software Engineering and Knowledge Engineering (SEKE '02)*. ACM, Ischia, Italy, pp. 425–428. <http://dx.doi.org/10.1145/568760.568833>.
- Jørgensen, Magne, Moløkken, Kjetil, 2004. Eliminating over-confidence in software development effort estimates. In: *Product Focused Software Process Improvement*. In: *Lecture Notes in Computer Science*, Springer, Berlin, Heidelberg, pp. 174–184. http://dx.doi.org/10.1007/978-3-540-24659-6_13.
- Jørgensen, M., Moløkken-Ostfold, K., 2004. Reasons for software effort estimation error: impact of respondent role, information collection approach, and data analysis method. *IEEE Trans. Softw. Eng.* 30 (12), 993–1007. <http://dx.doi.org/10.1109/TSE.2004.103>.
- Jørgensen, Magne, Shepperd, Martin, 2007. A systematic review of software development cost estimation studies. *IEEE Trans. Softw. Eng.* 33 (1), 33–53. <http://dx.doi.org/10.1109/TSE.2007.256943>.
- Jørgensen, Magne, Sjøberg, Dag I.K., 2001. Impact of effort estimates on software project work. *Inf. Softw. Technol.* 43 (2001), 10.
- Jørgensen, Magne, Sjøberg, Dag I.K., 2004. The impact of customer expectation on software development effort estimates. *Int. J. Proj. Manage.* 22 (4), 317–325. [http://dx.doi.org/10.1016/S0263-7863\(03\)00085-1](http://dx.doi.org/10.1016/S0263-7863(03)00085-1).
- Jørgensen, Magne, Teigen, Karl-Halvor, 2002. Uncertainty Intervals Versus Interval Uncertainty: an Alternative Method for Eliciting Effort Prediction Intervals in Software Development Projects. Singapore, pp. 343–352.
- Jørgensen, Magne, Teigen, Karl Halvor, Moløkken, Kjetil, 2004. Better sure than safe? Over-confidence in judgement based software development effort prediction intervals. *J. Syst. Softw.* 70 (1), 79–93. [http://dx.doi.org/10.1016/S0164-1212\(02\)00160-7](http://dx.doi.org/10.1016/S0164-1212(02)00160-7).
- Kahneman, Daniel, Sibony, Olivier, Sunstein, Cass R., 2021a. *Noise: A Flaw in Human Judgment*. Little, Brown Spark, New York.
- Kahneman, Daniel, Sibony, Olivier, Sunstein, Cass R., 2021b. *Crime and noisy punishment*. In: *Noise: A Flaw in Human Judgment*. Little, Brown Spark, New York, pp. 18–25.
- Kahneman, Daniel, Sibony, Olivier, Sunstein, Cass R., 2021c. *Debiasing and decision hygiene*. In: *Noise: A Flaw in Human Judgment*. Little, Brown Spark, New York, pp. 203–209.
- Kahneman, Daniel, Sibony, Olivier, Sunstein, Cass R., 2021d. *Introduction: two kinds of error*. In: *Noise: A Flaw in Human Judgment*. Little, Brown Spark, New York, pp. 9–15.
- Kahneman, Daniel, Sibony, Olivier, Sunstein, Cass R., 2021e. *Sequencing information in forensic science*. In: *Noise: A Flaw in Human Judgment*. Little, Brown Spark, New York, pp. 210–220.
- Kahneman, Daniel, Sibony, Olivier, Sunstein, Cass R., 2021f. *Your mind is a measuring instrument*. In: *Noise: A Flaw in Human Judgment*. Little, Brown Spark, New York, pp. 39–41.
- Karna, H., Gotovac, S., 2014. Estimators characteristics and effort estimation of software projects. In: *2014 9th International Conference on Software Engineering and Applications (ICSOFT-EA)*. pp. 26–35.
- Keaveney, Siobhan, Conboy, Kieran, Cost estimation in agile development projects. In: *ECIS 2006 Proceedings*. p. 16.
- Kitchenham, Barbara Ann, Budgen, David, Brereton, Pearl, 2015. *Evidence-Based Software Engineering and Systematic Reviews*, first ed. CRC Press.
- Koch, Stefan, Turk, Gerhard, 2011. Human resource related problems in agile and traditional software project process models. *Int. J. Inf. Technol. Proj. Manag.* 2 (2), 1–13. <http://dx.doi.org/10.4018/jitpm.2011040101>.
- Lagerström, Robert, Würtemberg, Liv Marcks von, Holm, Hannes, Luczak, Oscar, 2012. Identifying factors affecting software development cost and productivity. *Softw. Qual. J.* 20 (2), 395–417. <http://dx.doi.org/10.1007/s11219-011-9137-8>.

- Layman, Lucas, Nagappan, Nachiappan, Guckenheimer, Sam, Beehler, Jeff, Begel, Andrew, 2008. Mining software effort data: preliminary analysis of visual studio team system data. In: Proceedings of the 2008 International Working Conference on Mining Software Repositories (MSR '08). Association for Computing Machinery, New York, NY, USA, pp. 43–46. <http://dx.doi.org/10.1145/1370750.1370762>.
- Lederer, Albert, Mirani, Rajesh, 1990. Information system cost estimating: A management perspective. MIS Q. 14 (2), Retrieved from <https://aisel.aisnet.org/misq/vol14/iss2/3>.
- Lederer, Albert L., Prasad, Jayesh, 1991. The validation of a political model of information systems development cost estimating. ACM SIGCPR Comput. Pers. 13 (2), 47–57. <http://dx.doi.org/10.1145/122393.122398>.
- Lederer, Albert L., Prasad, Jayesh, 1995a. Causes of inaccurate software development cost estimates. J. Syst. Softw. 31 (2), 125–134. [http://dx.doi.org/10.1016/0164-1212\(94\)00092-2](http://dx.doi.org/10.1016/0164-1212(94)00092-2).
- Lederer, Albert L., Prasad, Jayesh, 1995b. Perceptual congruence and information systems cost estimating. In: Proceedings of the 1995 ACM SIGCPR Conference on Supporting Teams, Groups, and Learning Inside and Outside the IS Function Reinventing IS (SIGCPR '95). Association for Computing Machinery, Nashville, Tennessee, USA, pp. 50–59. <http://dx.doi.org/10.1145/212490.212504>.
- Lederer, A.L., Prasad, J., 1998. A causal model for software cost estimating error. IEEE Trans. Softw. Eng. 24 (2), 137–148. <http://dx.doi.org/10.1109/32.666827>.
- Lee, Michael, Rothenberger, Marcus, Peffers, Ken, 2011. Identifying effort estimation factors for corrective maintenance in object-oriented systems. In: AMCIS 2011 Proceedings. Retrieved from https://aisel.aisnet.org/amcis2011_submissions/186.
- Lenarduzzi, Valentina, 2015. Could social factors influence the effort software estimation? In: Proceedings of the 7th International Workshop on Social Software Engineering (SSE 2015). Association for Computing Machinery, New York, NY, USA, pp. 21–24. <http://dx.doi.org/10.1145/2804381.2804385>.
- Lind, M.R., Sulek, J.M., 1998. Undersizing software systems: third versus fourth generation software development. Eur. J. Inf. Syst. 7 (4), 261–268. <http://dx.doi.org/10.1038/sj.ejis.3000308>.
- Little, T., 2006. Schedule estimation and uncertainty surrounding the cone of uncertainty. IEEE Softw. 23 (3), 48–54. <http://dx.doi.org/10.1109/MS.2006.82>.
- Löhre, Erik, Jørgensen, Magne, 2016. Numerical anchors and their strong effects on software development effort estimates. J. Syst. Softw. 116, 49–56. <http://dx.doi.org/10.1016/j.jss.2015.03.015>.
- Lorenz, Jan, Rauhut, Heiko, Schweitzer, Frank, Helbing, Dirk, 2011. How social influence can undermine the wisdom of crowd effect. Proc. Natl. Acad. Sci. USA 108 (22), 9020–9025. <http://dx.doi.org/10.1073/pnas.1008636108>.
- Magazinius, Ana, Börjesson, Sofia, Feldt, Robert, 2012. Investigating intentional distortions in software cost estimation – An exploratory study. J. Syst. Softw. 85 (8), 1770–1781. <http://dx.doi.org/10.1016/j.jss.2012.03.026>.
- Magazinius, A., Feldt, R., 2011. Confirming distortional behaviors in software cost estimation practice. In: 2011 37th EUROMICRO Conference on Software Engineering and Advanced Applications. pp. 411–418. <http://dx.doi.org/10.1109/SEAA.2011.61>.
- Magazinius, A., Svensson, R.B., 2014. Effects of feature complexity on software effort estimates – an exploratory study. In: 2014 40th EUROMICRO Conference on Software Engineering and Advanced Applications. pp. 301–304. <http://dx.doi.org/10.1109/SEAA.2014.69>.
- Magazinovic, Ana, Pernstål, Joakim, 2008. Any other cost estimation inhibitors? In: Proceedings of the Second ACM-IEEE International Symposium on Empirical Software Engineering and Measurement - ESEM '08. ACM Press, Kaiserslautern, Germany, p. 233. <http://dx.doi.org/10.1145/1414004.1414042>.
- Mahnič, Viljan, Hovelja, Tomaž, 2012. On using planning poker for estimating user stories. J. Syst. Softw. 85 (9), 2086–2095. <http://dx.doi.org/10.1016/j.jss.2012.04.005>.
- Matos, Olavo, Fortaleza, Luiz, Conte, Tayana, Mendes, Emilia, 2013. Realising web effort estimation. In: Proceedings of the 17th International Conference on Evaluation and Assessment in Software Engineering. Association for Computing Machinery, Porto de Galinhas, pp. 12–23. <http://dx.doi.org/10.1145/2460999.2461002>.
- Matsubara, Patricia, Gadelha, Bruno, Steinmacher, Igor, Conte, Tayana, 2021. Supplementary Material for the SEXTAMT. figshare, Retrieved April 29, 2021 from <https://doi.org/10.6084/m9.figshare.14502405.v2>.
- McConnell, Steve, 2006a. Flow of software estimates on a well-estimated project. In: Software Estimation: Demystifying the Black Art. Microsoft Press, Redmond, pp. 171–180.
- McConnell, Steve, 2006b. Introduction to estimation techniques. In: Software Estimation: Demystifying the Black Art. Microsoft Press, Redmond, pp. 171–180.
- McConnell, Steve, 2006c. What is an estimate? In: Software Estimation: Demystifying the Black Art. Microsoft Press, Redmond, pp. 3–14.
- McDonald, James, 2005. The impact of project planning team experience on software project cost estimates. Empir. Softw. Eng. 10 (2), 219–234. <http://dx.doi.org/10.1007/s10664-004-6192-9>.
- McGarry, F., Burke, S., Decker, B., 1998. Measuring the impacts individual process maturity attributes have on software products. In: Proceedings Fifth International Software Metrics Symposium. Metrics (Cat. No. 98TB100262), pp. 52–60. <http://dx.doi.org/10.1109/METRIC.1998.731226>.
- Mendes, Emilia, Mosley, Nile, Counsell, Steve, 2005. Investigating Web size metrics for early Web cost estimation. J. Syst. Softw. 77 (2), 157–172. <http://dx.doi.org/10.1016/j.jss.2004.08.034>.
- Molokken, K., Jørgensen, M., 2003. A review of software surveys on software effort estimation. In: 2003 International Symposium on Empirical Software Engineering, 2003. ISESE 2003. Proceedings. pp. 223–230. <http://dx.doi.org/10.1109/ISESE.2003.1237981>.
- Moløkken, Kjetil, Jørgensen, Magne, 2005. Expert estimation of web-development projects: Are software professionals in technical roles more optimistic than those in non-technical roles? Empir. Softw. Eng. 10 (1), 7–30. <http://dx.doi.org/10.1023/B:EMSE.0000048321.46871.2e>.
- Molokken-Ostvold, K., Furulund, K.M., Customer Collaboration, 2007. The relationship between and Software Project Overruns. In: Agile 2007 (AGILE 2007). pp. 72–83. <http://dx.doi.org/10.1109/AGILE.2007.57>.
- Moløkken-Ostvold, Kjetil, Haugen, Nils Christian, Benestad, Hans Christian, 2008. Using planning poker for combining expert estimates in software projects. J. Syst. Softw. 81 (12), 2106–2117. <http://dx.doi.org/10.1016/j.jss.2008.03.058>.
- Moløkken-Ostvold, Kjetil, Jørgensen, Magne, 2004. Group processes in software effort estimation. Empir. Softw. Eng. 9 (4), 315–334. <http://dx.doi.org/10.1023/B:EMSE.0000039882.39206.5a>.
- Molokken-Ostvold, K., Jørgensen, M., 2005. A comparison of software project overruns - flexible versus sequential development models. IEEE Trans. Softw. Eng. 31 (9), 754–766. <http://dx.doi.org/10.1109/TSE.2005.96>.
- Morgenshtern, Ofer, Raz, Tzvi, Dvir, Dov, 2007. Factors affecting duration and effort estimation errors in software development projects. Inf. Softw. Technol. 49 (8), 827–837. <http://dx.doi.org/10.1016/j.infsof.2006.09.006>.
- Nugroho, Ariadi, Lange, Christian F.J., 2008. On the relation between class-count and modeling effort. In: Models in Software Engineering (Lecture Notes in Computer Science). Springer, Berlin, Heidelberg, pp. 93–104. http://dx.doi.org/10.1007/978-3-540-69073-3_11.
- Ohlsson, M.C., Wohlin, C., Regnell, B., 1998. A project effort estimation study. Inf. Softw. Technol. 40 (14), 831–839. [http://dx.doi.org/10.1016/S0950-5849\(98\)00097-4](http://dx.doi.org/10.1016/S0950-5849(98)00097-4).
- Passing, Ursula, Shepperd, Martin, 2003. An experiment on software project size and effort estimation. In: Proceedings of the 2003 International Symposium on Empirical Software Engineering (ISESE '03). IEEE Computer Society, Washington, DC, USA, 120–. Retrieved December 19, 2019 from <http://dl.acm.org/citation.cfm?id=942801.943632>.
- Petersen, Kai, Vakkalanka, Sairam, Kuzniarz, Ludwik, 2015. Guidelines for conducting systematic mapping studies in software engineering: An update. Inf. Softw. Technol. 64, 1–18.
- Project Management Institute, 2017a. A Guide to the Project Management Body of Knowledge (PMBOK Guide), sixth ed. Project Management Institute, Newtown Square, PA.
- Rahikkala, Jurka, Hyrynsalmi, Sami, Leppänen, Ville, 2015a. Accounting Testing in Software Cost Estimation: A Case Study of the Current Practice and Impacts. Tampere, Finland, p. 15.
- Rahikkala, Jurka, Hyrynsalmi, Sami, Leppänen, Ville, Porres, Ivan, 2018. The role of organisational phenomena in software cost estimation: A case study of supporting and hindering factors. e-Inf. Softw. Eng. J. 12 (1), 167–198. <http://dx.doi.org/10.5277/e-inf180107>.
- Rahikkala, Jurka, Leppänen, Ville, Ruohonen, Jukka, Holvitie, Johannes, 2015b. Top management support in software cost estimation: A study of attitudes and practice in Finland. Int. J. Manag. Proj. Bus. 8 (3), 513–532. <http://dx.doi.org/10.1108/IJMPB-11-2014-0076>.
- Ramesur, Melvina Autar, Nagawah, Soulakshme Devi, 2020. Factors affecting sprint effort estimation. In: Advanced Computing and Intelligent Engineering (Advances in Intelligent Systems and Computing). Springer, Singapore, pp. 507–518. http://dx.doi.org/10.1007/978-981-15-1483-8_43.
- Rozalina, Rianti, Mansor, Zulkefli, 2018. Validated software cost estimation factors for government projects using rasch measurement model. Int. J. Adv. Sci. Eng. Inf. Technol. 8 (5), 1890–1896–1896. <http://dx.doi.org/10.18517/ijaseit.8.5.6386>.
- Sehra, Sumeet Kaur, Brar, Yadwinder Singh, Kaur, Navdeep, Sehra, Sukhjit Singh, 2017. Research patterns and trends in software effort estimation. Inf. Softw. Technol. 91, 1–21. <http://dx.doi.org/10.1016/j.infsof.2017.06.002>.
- Shepperd, Martin, Mair, Carolyn, Jørgensen, Magne, 2018. An experimental evaluation of a de-biasing intervention for professional software developers. In: Proceedings of the 33rd Annual ACM Symposium on Applied Computing (SAC '18). ACM, New York, NY, USA, pp. 1510–1517. <http://dx.doi.org/10.1145/3167132.3167293>.
- Shmueli, Ofira, Pliskin, Nava, Fink, Lior, 2016. Can the outside-view approach improve planning decisions in software development projects? Inf. Syst. J. 26 (4), 395–418. <http://dx.doi.org/10.1111/isj.12091>.

- Silva-de-Souza, Thiago, Travassos, Guilherme Horta, 2017. Observing effort factors in the test design & implementation process of web services projects. In: Proceedings of the 2nd Brazilian Symposium on Systematic and Automated Software Testing (SAST). Association for Computing Machinery, New York, NY, USA, pp. 1–10. <http://dx.doi.org/10.1145/3128473.3128480>.
- Stol, Klaas-Jan, Fitzgerald, Brian, 2018. The ABC of software engineering research. *ACM Trans. Softw. Eng. Methodol.* 27 (3), 11:1–11:51. <http://dx.doi.org/10.1145/3241743>.
- Storey, Margaret-Anne, Ernst, Neil A., Williams, Courtney, Kalliamvakou, Eirini, 2020. The who, what, how of software engineering research: a socio-technical framework. *Empir. Softw. Eng.* 25 (5), 4097–4129. <http://dx.doi.org/10.1007/s10664-020-09858-z>.
- Subramanian, Girish, Pendharkar, Parag C., Pai, Dinesh, 2017. An examination of determinants of software testing and project management effort. *J. Comput. Inf. Syst.* 57 (2), 123–129.
- Subramanian, Girish H., Pendharkar, Parag C., Wallace, Mary, 2006. An empirical study of the effect of complexity, platform, and program type on software development effort of business applications. *Empir. Softw. Eng.* 11 (4), 541–553. <http://dx.doi.org/10.1007/s10664-006-9023-3>.
- Suliman, Safa, Kadoda, Gada, 2017. Factors that Influence Software Project Cost and Schedule Estimation. Elnihood, Sudan, Retrieved February 10, 2021 from <https://ieeexplore.ieee.org/document/8293053>.
- Taff, L.M., Borcherding, J.W., Hudgins, W.R., 1991. Estimeetings: development estimates and a front-end process for a large project. *IEEE Trans. Softw. Eng.* 17 (8), 839–849. <http://dx.doi.org/10.1109/32.83918>.
- Tamrakar, R., Jørgensen, M., 2012. Does the use of Fibonacci numbers in planning poker affect effort estimates? In: 16th International Conference on Evaluation Assessment in Software Engineering (EASE 2012), pp. 228–232. <http://dx.doi.org/10.1049/ic.2012.0030>.
- Tanveer, Binish, Guzmán, Liliana, Engel, Ulf Martin, 2017. Effort estimation in agile software development: Case study and improvement framework. *J. Softw.: Evol. Process* 29 (11), e1862. <http://dx.doi.org/10.1002/smr.1862>.
- Trendowicz, Adam, Münch, Jürgen, Jeffery, Ross, 2011. State of the practice in software effort estimation: A survey and literature review. In: *Software Engineering Techniques (Lecture Notes in Computer Science)*. Springer Berlin Heidelberg, pp. 232–245.
- Tripathi, N., Seppänen, P., Oivo, M., Similä, J., Liukkunen, K., 2017. The effect of competitor interaction on startup's product development. In: 2017 43rd Euromicro Conference on Software Engineering and Advanced Applications (SEAA), pp. 125–132. <http://dx.doi.org/10.1109/SEAA.2017.34>.
- Usman, Muhammad, Börstler, Jürgen, Petersen, Kai, 2017. An effort estimation taxonomy for agile software development. *Int. J. Softw. Eng. Knowl. Eng.* 27 (04), 641–674. <http://dx.doi.org/10.1142/S0218194017500243>.
- Usman, Muhammad, Britto, Ricardo, Damm, Lars-Ola, Börstler, Jürgen, 2018a. Effort estimation in large-scale software development: An industrial case study. *Inf. Softw. Technol.* 99, 21–40. <http://dx.doi.org/10.1016/j.infsof.2018.02.009>.
- Usman, Muhammad, Mendes, Emilia, Börstler, Jürgen, 2015. Effort estimation in agile software development: a survey on the state of the practice. In: Proceedings of the 19th International Conference on Evaluation and Assessment in Software Engineering (EASE '15). Association for Computing Machinery, Nanjing, China, pp. 1–10. <http://dx.doi.org/10.1145/2745802.2745813>.
- Usman, Muhammad, Petersen, Kai, Börstler, Jürgen, Neto, Pedro Santos, 2018b. Developing and using checklists to improve software effort estimation: A multi-case study. *J. Syst. Softw.* 146, 286–309. <http://dx.doi.org/10.1016/j.jss.2018.09.054>.
- Valerdi, R., 2007. Cognitive limits of software cost estimation. In: First International Symposium on Empirical Software Engineering and Measurement (ESEM 2007), pp. 117–125. <http://dx.doi.org/10.1109/ESEM.2007.85>.
- Vicinanza, Steven S., Mukhopadhyay, Tridas, Prietula, Michael J., 1991. Software-effort estimation: An exploratory study of expert performance. *Inf. Syst. Res.* 2 (4), 243–262. <http://dx.doi.org/10.1287/isre.2.4.243>.
- Vijayakumar, S., 1997. Use of historical data In software cost estimation. *Comput. Control Eng. J.* 8 (3), 113–119. <http://dx.doi.org/10.1049/cce:19970303>.
- Wohlin, Claes, 2014. Guidelines for snowballing in systematic literature studies and a replication in software engineering. In: Proceedings of the 18th International Conference on Evaluation and Assessment in Software Engineering (EASE '14). ACM, New York, NY, USA, pp. 38:1–38:10. <http://dx.doi.org/10.1145/2601248.2601268>.
- Yang, Da, Wang, Qing, Li, Mingshu, Yang, Ye, Ye, Kai, Du, Jing, 2008. A survey on software cost estimation in the chinese software industry. In: Proceedings of the Second ACM-IEEE International Symposium on Empirical Software Engineering and Measurement - ESEM '08. ACM Press, Kaiserslautern, Germany, p. 253. <http://dx.doi.org/10.1145/1414004.1414045>.
- Zapata, A.H., Chaudron, M.R.V., 2013. An empirical study into the accuracy of it estimations and its influencing factors. *Int. J. Softw. Eng. Knowl. Eng.* 23 (04), 409–432. <http://dx.doi.org/10.1142/S0218194013400081>.
- Zarour, Ahmed, Zein, Samer, 2019. Software development estimation techniques in industrial contexts: An exploratory multiple case-study. *Int. J. Technol. Educ. Sci.* 3 (2), 72–84.
- Zhang, He, Babar, Muhammad Ali, Tell, Paolo, 2011. Identifying relevant studies in software engineering. *Inf. Softw. Technol.* 53 (6), 625–637. <http://dx.doi.org/10.1016/j.infsof.2010.12.010>.



Patrícia Matsubara is a Ph.D. student at the Federal University of Amazonas (UFAM – Brazil). She is also a Lecturer at the Federal University of Mato Grosso do Sul (UFMS – Brazil). She received her master degree in Computer Science at the Federal University of Goiás (UFG – Brazil) in 2010. Her research interests include Human Aspects of Software Engineering, Software Effort Estimation, and Empirical Software Engineering.



Bruno Gadelha is Associate Professor at the Federal University of Amazonas (UFAM – Brazil). He holds a Ph.D. in Informatics from the Pontifical Catholic University of Rio de Janeiro (PUC RIO – Brazil). He is co-leader of the Usability and Software Engineering Group (USES) lab at UFAM. His research focuses on Virtual Learning Environments, Software Product Lines, and Computer Supported Cooperative Work.



Igor Steinmacher is an Assistant Professor at the Federal University of Technology, Paraná (UTFPR – Brazil). He received a Ph.D. in Computer Science from the University of São Paulo (USP – Brazil). He researches the intersections of Software Engineering (SE) and Computer Supported Cooperative Work (CSCW). Currently, his research focuses on the behavior of developers in Open Source Communities, including support of newcomers, the impact of Bots in the community, and gender bias in Open Source Software. His interests include Open Source Software, Human Aspects of Software Engineering, Empirical Software Engineering, and Mining Software Repositories techniques.



Tayana Conte holds a Ph.D. in Systems Engineering and Computer from the Federal University of Rio de Janeiro (UFRJ). She is an Associate Professor at the Federal University of Amazonas (UFAM – Brazil), heading the Usability and Software Engineering (USES) lab. Her research interests include the intersection between Software Engineering and Human-Computer Interaction, Software Quality, Human-Centered Computing, and Empirical Software Engineering.