# Bridging the Gap: Adapting Evidence to Decision Frameworks to support the link between Software Engineering academia and industry

Patrícia G. F. Matsubara
Federal University of Mato Grosso do Sul (UFMS)
Campo Grande, Brazil
patricia.gomes@ufms.br

Tayana Uchôa Conte
Federal University of Amazonas (UFAM)
Manaus, Brazil
tayana@icomp.ufam.edu.br

## Abstract

Over twenty years ago, the Software Engineering (SE) research community have been involved with Evidence-Based Software Engineering (EBSE). EBSE aims to inform industrial practice with the best evidence from rigorous research, preferably from systematic literature reviews (SLRs). Since then, SE researchers have conducted many SLRs, perfected their SLR procedures, proposed alternative ways of presenting their results (such as Evidence Briefings), and profusely discussed how to conduct research that impacts practice. Nevertheless, there is still a feeling that SLRs' results are not reaching practitioners. Something is missing. In this vision paper, we introduce Evidence to Decision (EtD) frameworks from the health sciences, which propose gathering experts in panels to assess the existing best evidence about the impact of an intervention in all relevant outcomes and make structured recommendations based on them. The insight we can leverage from EtD frameworks is not their structure per se but all the relevant criteria for making recommendations to practitioners from SLRs. Furthermore, we provide a worked example based on an SE SLR. We also discuss the challenges the SE research and practice community may face when adopting EtD frameworks, highlighting the need for more comprehensive criteria in our recommendations to industry practitioners.

## CCS Concepts

• **Software and its engineering**; • **Social and professional topics** → **Management of computing and information systems**;

## Keywords

Evidence-Based Software Engineering, Evidence to Decision

## 1 Introduction

Evidence-Based Software Engineering (EBSE) aims to provide the means by which the best evidence from rigorous research could be used to make decisions in industrial practice [15]. An implication from this is the requirement of research in Software Engineering (SE) to be relevant for industry practice. Another implication is the need to search for the currently available best evidence. SE researchers achieve this through systematic literature reviews (SLRs), whose goal is "to search for and identify all relevant material related to a given topic" [14, Section 1.3]. Researchers follow a set of well-defined procedures, which are meant to be "as objective, analytical, and repeatable as possible" [14, Section 1.3].

Nevertheless, finding all the available evidence regarding a topic employing a reliable process does not necessarily mean such evidence is "credible", i.e., evidence we can trust [19] to make decisions and take action in the real world. There might be insufficient evidence to support a strong conclusion, and more research might still be required to raise confidence in one intervention or another. Considering this, Dyba et al. [8] introduced the SE research community to the GRADE (Grades of Recommendation, Assessment, Development, and Evaluation) methodology for assessing how "strong" is the existing evidence supporting research findings from SLRs.

Still, even if SE researchers find strong and reliable evidence supporting an intervention, it must reach practitioners' ears and bring with it all relevant information for decision-making. In part, this requires presenting actionable results from research and describing the implications of findings for practice. In other words, this requires the making of clear recommendations for practitioners. Curiously, GRADE is not only about the strength assessment of a body of evidence but is also about making and assessing recommendations for practitioners—something that we, as SE researchers, seem to have missed from the seminal paper of Dyba et al. [8]. In this paper, we address this gap by briefly summarizing an approach for making such recommendations, thus addressing the connection between academia and industry: Evidence to Decision (EtD) frameworks. We also present a worked example of EtD frameworks and discuss the challenges to its application in the SE domain. In the next section, we present related work.

## 2 Background and Related Work

In 2008, Dyba et al. [8] raised concerns about how much confidence one can have in SLR results, introducing the SE community to the need of evaluating the quality of primary studies. Since then, a growing number of SLRs have assessed the quality of primary studies in SE[20]. Dyba et al. [8] also presented the GRADE methodology for

assessing the strength of a body of evidence. GRADE suggests that the analysis of the strength of the evidence supporting a research finding is more than just the quality of the primary studies. It should include other relevant dimensions: publication bias, inconsistency of results, imprecision of results, and the indirection coming from differences in population, intervention, and outcome of interest targeted for answering a research question. Issues in these dimensions are reasons to downgrade the confidence of otherwise strong evidence [5]. The approach recommends the summarization of research findings and the judgments for each relevant dimension of the strength of evidence in Evidence Profiles (where judgments for all GRADE dimensions are detailed) or Summary of findings (where only the overall judgment for strength is presented) tables [11]. However, very few SLR in the SE domain have assessed the strength of the evidence they gathered [16].

Cartaxo et al. [7] proposed Evidence Briefings (EBs), brief documents summarizing SLR results, to disseminate SLR' findings to practitioners. Despite being a clear and understandable medium, EBs may not fully answer the practical questions practitioners have [7]. Also, EBs may not include all the information required to support a decision from a practitioner in their daily practice, such as the resources needed for adopting a method or tool [1].

Back to GRADE, the approach also aims to assist the development of guidelines with recommendations for practitioners based on the evidence from SLRs [11]. The working group responsible for developing GRADE proposed the idea of Evidence to Decision (EtD) frameworks, which aggregate information (evidence included) to assist practitioners' decision-making regarding interventions. We explain this framework in detail in the next section.

## 3 The elements of an EtD Framework

Originally, an EtD is meant to help creating and assessing a recommendation or guideline regarding a clinical question about what to do when treating a patient, considering the evidence about existing interventions. Its creation involves two major stages, as Figure 1 illustrates: an evidence synthesis and a recommendation stage.

During the evidence synthesis stage, reviewers gather the best available evidence to answer the question through typical SLR procedures. This means synthesizing evidence and concluding about the targeted intervention's effects on critical outcomes of interest. Reviewers also critically appraise the confidence (or strength) on such evidence, i.e. on the estimates of the effects of the intervention.

During the recommendation stage, a panel of experts knowledgeable in SLR procedures, the GRADE system, and the topic of interest are responsible for making the recommendation. They will use the information from the SLR, along with other criteria defined by EtD frameworks to do so. The panel will follow three steps, which correspond to the three sections of the EtD structure [3], which we explore in the next sections:

(1) Formulating the question (Section 3.1);
(2) Making an assessment (Section 3.2); and
(3) Drawing a conclusion (Section 3.3).

## 3.1 Formulating the question

The first step in the EtD framework is formulating a question to address. Table 1 presents all the items that a panel of experts must consider when defining the question to structure their judgments and to reach a recommendation. It also shows a worked example of a question for the SE domain on the topic of pair programming in comparison with solo programming. The worked example is inspired by the SLR results of Hannay et al. [13], but not necessarily wholly faithful to the most current and best-existing evidence on the topic. When needed, we also inserted fictitious decisions for the panel to increase the force of the example and provide means for reflection on what we can represent with EtD frameworks.

The question uses the PICO (population, intervention, comparison, outcome) format, and a few additional details [3]. When defining the population, the panel must clarify any relevant characteristics of the targeted people for the intervention. For example, only people at a high-risk of a catastrophic outcome (e.g. death) might be targeted. Determining how broad is the population or the intervention is a challenging decision when formulating the question, and a practical guideline is to consider variations of magnitudes of effects of the treatment across subgroups of the population or across differences in treatment options [12]. The panel should also be specific about the intervention and the comparison, i.e., the alternative treatment. In medicine, this can involve decisions about doses, and a combination of drugs. The outcomes of interest also clarify the variables used in the assessment of the intervention and the comparison.

Depending on the question, the definition of the setting is crucial. It helps identify constraints that might limit the recommendations, e.g., the wealth of a country can impose limitations on feasible treatments to provide in health care systems. The perspective helps to create a "practitioner" focus, clarifying whether the panel has to answer with an individual or group point of view in mind. It will influence the outcomes of interest to analyze and even economic evaluations [6]. For instance, an individual might have low concerns regarding the costs of a treatment if a health system pays for it. However, the cost information will be most important when assuming a health system perspective. The subgroup aids in specifying sub-populations for which recommendations might differ from the overall population, given their specific characteristics and the existing evidence.

Also, the panel needs to inform whether any of its members have conflicts of interest—for instance, on promoting the intervention or the comparison for any reason. Valid reasons can include participation in organizations that provide tools or training associated with the intervention, for example. Identifying conflicts will aid the panel in dealing with them, either by just declaring the conflict or excluding members from discussions about specific questions [3].

In the worked example, the panel's recommendation is made without specific constraints, applicable to organizations or teams developing software. It notes that the recommendation might differ for inexperienced developers, who may achieve different results than experienced ones. One of the panel members declared a conflict of interest, as he is a member of an organization that provides training in the topic of interest. Therefore, he might have a more favorable view of the effects of the intervention. In summary, the question provides the details to guide the assessment criteria for making a recommendation. The next section presents these criteria and continues with our worked example.
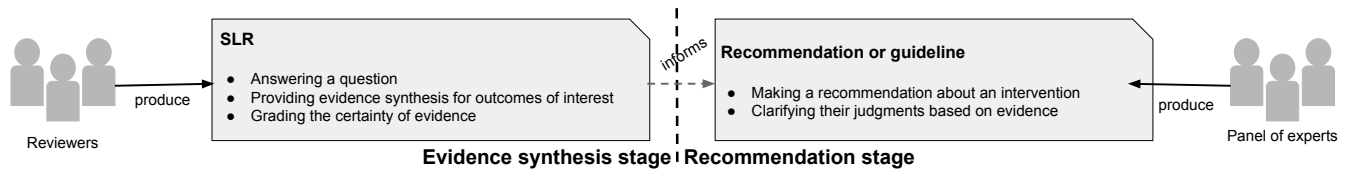
**Figure 1: The relationship between SLRs and recommendations/guidelines.**

| **Question:** The question using the PICO (population, intervention, comparison, outcome) format. **Example:** Should pair programming be adopted in software development projects? | | |
|---|---|---|
| **Question details:** | | |
| **(P) Population:** | The characteristics of the targeted population for the intervention/comparison. | **Ex.:** Software developers. |
| **(I) Intervention:** | The intervention that the recommendation focuses on. | **Ex.:** Pair programming. |
| **(C) Comparison:** | An alternative treatment to the intervention. | **Ex.:** Solo programming. |
| **(O) Main outcomes:** | The critical outcomes of interest to analyze when recommending an intervention or its comparison. | **Ex.:** Duration, effort, quality. |
| **Setting:** | The characteristics of the environment where the recommendation is to be implemented. | **Ex.:** A software development organization in general. |
| **Perspective:** | The perspective taken for the recommendation. It can be of an individual practitioner deciding to adopt the intervention, or of an organization that can also adopt it. | **Ex.:** The organization/team. |
| **Subgroups:** | Description of subgroups for which the recommendation is likely to vary due to special characteristics they possess, if there is any. | **Ex.:** Inexperienced software developers. |
| **Conflicts of interest:** | Any conflict of interest that any panel member might have that can lead to the favouring of the intervention or of the comparison. | **Ex.:** Mr. James Smith is a member of an organization that provides training for pair programming dynamics. The remaining panel members declared no conflict. |

**Table 1: Items (from [3]) for formulating the question, and their descriptions (by the authors).**

## 3.2 Making an assessment

In the second step, the panel judges a set of criteria relevant to practitioners' decision-making. These criteria vary with the type of decision. In the health domain, types of decisions include (i) clinical recommendations from an individual perspective, (ii) clinical recommendations from a population perspective, (iii) coverage decisions[1], (iv) health system and public health system recommendations/decisions, and (v) diagnostic, screening, or other tests [3].

Table 2 shows a few of the criteria for two decision types in Alonso-Coello et al. [2]: population-based and individual recommendations. Some criteria are type-specific, while others apply to both types. The supplementary material [17] shows the judgments of a fictitious panel of experts for some of these relevant criteria from the perspective of organizations or teams pondering whether to adopt pair programming, aligned with the question in Section 3.1.

The judgments are guided by the question Section 3.1 presents. For instance, for the criterion "How substantial are the desirable anticipated effects?" (from Table 2) the panel used the Summary of Findings table with the evidence for all the outcomes of interest

from the question: quality, duration, and effort. The judgment for the criterion "Is the intervention acceptable to key stakeholders?" also consider the perspective informed in the question: organizations and teams, instead of the perspective of individual developers. The panel could even explore the relevant subgroup included in the question during their judgments by creating a separate Summary of Findings with evidence for this specific subgroup or additional considerations if there is enough evidence that the results for the outcomes of interest were different for these groups.

## 3.3 Drawing a conclusion

In the third step, the panel draws a conclusion from their judgments for all criteria and deliver an actionable, specific, and clear recommendation for practitioners. [3]. The panel informs the recommendation's direction: in favor or against the intervention. The panel must also declare the recommendation's strength, which can be either strong—when it is highly likely that a patient will accept the intervention/treatment in any circumstances—or discretionary— when a patient's acceptance of the intervention/treatment will vary depending on the balance between desirable and undesirable effects, and their values and preferences, among other things [4].

---

[1]For health insurance providers.

| Population perspective | Individual perspective |
|---|---|
| Is the problem a priority (from a population perspective)? | Is the problem a priority (from an individual perspective)? |
| How substantial are the desirable anticipated effects? ||
| What is the overall certainty of the evidence of effects? ||
| Does the balance between desirable and undesirable effects favour the intervention or the comparison? ||
| What is the certainty of the evidence of resource requirements (costs)? | Does the cost effectiveness of the intervention (the out-of-pocket cost relative to the net desirable effect) favour the intervention or the comparison? |
| Is the intervention acceptable to key stakeholders?* | Is the intervention acceptable to patients, their caregivers, and healthcare providers? |

**Table 2: Examples of criteria for making assessments for population and individual perspectives (from [2]).**

The supplementary material also shows the panel's recommendations for our worked example [17]. Such recommendation is discretionary, as different organizations and teams will likely vary in their decision to adopt pair programming for the many reasons explained in the recommendation. For instance, for some organizations, the undesirable effect of higher effort might be unacceptable. The recommendation also includes details that can help practitioners implement the intervention and monitor and evaluate its efficacy. In the worked example, the panel considers it relevant for organizations to have a productivity measurement system in place to help assess the impact of pair programming on productivity. The recommendation also presents further research directions to increase confidence in the evidence or explore gaps likely necessary for better comprehension of the intervention's effects on relevant outcomes of interest. In the worked example, the panel highlights the importance of considering the individual perspective when making decisions by identifying software developers' values and preferences towards working solo or in pairs. The panel also considers the need to perform primary research on the effects of pair programming on other critical and unexplored outcomes of interest, such as learning, especially when onboarding new team members. Moreover, they consider that measuring developers' satisfaction and well-being helps to understand productivity better [9], so it is also a critical outcome to further research.

The recommendation created with an EtD framework will stem from a clear question defined with practitioners in mind. Then, the ETD framework guides the panel in judging the existing evidence, accounting for all relevant issues practitioners will likely consider in making a decision about an intervention. Such evidence is unlikely to come from a single SLR [21], so EtD frameworks provide the means for gathering it in a single place. The framework also allows the panel to highlight the strength of their recommendations based on the objective criteria assessed, which is unusual in SLRs and can be a relevant factor for decision-making. Panels are also expected to update EtD frameworks and their recommendations occasionally or when new relevant evidence emerges. This can also benefit practitioners, enabling them to update their knowledge in their fields of expertise more easily.

## 4 Envisioned Challenges for SE

While the adoption of EtD frameworks in SE may present challenges, it also holds the potential to significantly impact the industry. This section delves into the potential hurdles and issues the SE community must resolve to leverage their benefits.

**Who are the entities entitled to make a recommendation for industry practitioners, if any?** Clinicians can turn to national specialty associations and even worldwide organizations such as WHO (World Health Organization) for guidelines. SE chapters of professional associations such as ACM SIGSOFT (ACM Special Software Engineering Group) and IEEE TCSE (IEEE Technical Community on Software Engineering) could assume a similar role in Software Engineering by organizing expert panels to develop and issue guidelines. Additionally, invited panels of technical experts could contribute to guidelines published through Special Issues in selected venues.

**We need platforms where practitioners can look for answers to their questions.** Cartaxo et al. [7] found that most survey participants preferred platforms like StackExchange over static EBs. This suggests we should design and promote engaging platforms where the information practitioners need is readily available in the form of answers that EtD frameworks provide.

**What are industry practitioners' different types of decisions and their relevant criteria?** In Section 3.2, we present the different types of decisions discussed in EtD frameworks for the health domain. SE researchers must also identify the types of decisions for which SE research can provide evidence, as this connects to the criteria for making a recommendation. SE clearly can make recommendations from an individual perspective—to software developers, for instance. Kitchenham et al. [15] highlighted the perspectives of project managers interested in results for specific projects and senior managers interested in results for a department or the whole organization. Such perspectives suggest there is likely a need for recommendations for group perspectives, such as teams and organizations, in SE too.

**We must aggregate evidence around outcomes of interest.** Identifying different perspectives is paramount to defining relevant criteria for recommendations and clarifying the different outcomes of interests to consider when conducting primary studies. For instance, while individual developers and software teams might focus on maintainability when considering quality as an outcome of interest, stakeholders with business or managerial points of view might focus on user experience or customer satisfaction [10]. Additionally, acceptability and feasibility may vary between group and individual perspectives, depending on the intervention.

**We also need to define the equivalent of "clinical questions".** Perhaps we can call such questions "practice questions" to adapt to the SE context and differentiate them from research questions. This connects with the finding of Cartaxo et al. [7] that,

though participants might find SLR findings important, not necessarily such findings answer *their* questions. This is critical and current for primary research also. Winters [18] have recently remarked how research in a highly relevant SE venue still addresses questions of no interest to the software industry or provides solutions with no practical context.

## 5 Conclusion

The thing with EtD frameworks is not their structure per se but all the information they show practitioners need for making meaningful decisions against or in favor of an intervention. Evidence about interventions' effectiveness is insufficient, as the intervention might be too costly to be adopted [1], or unacceptable to a relevant set of stakeholders. In this paper, we describe the EtD frameworks and present a worked example of their application in the SE domain. We show that EtD frameworks include information lacking in individual SLRs in SE. EtD frameworks can also aid the SE research community in directing research efforts, by clarifying missing evidence to derive guidelines for practitioners or subgroups of the population who can benefit from more specific recommendations. Therefore, EtD frameworks are a promising method to help us close the academia-industry gap.

## Acknowledgments

## References

[1] Nauman bin Ali. 2016. Is effectiveness sufficient to choose an intervention? Considering resource use in empirical software engineering. In *Proceedings of the 10th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM '16)*. Association for Computing Machinery, New York, NY, USA, 1–6. doi:10.1145/2961111.2962631

[2] Pablo Alonso-Coello, Andrew D. Oxman, Jenny Moberg, Romina Brignardello-Petersen, Elie A. Akl, Marina Davoli, Shaun Treweek, Reem A. Mustafa, Per O. Vandvik, Joerg Meerpohl, Gordon H. Guyatt, Holger J. Schünemann, and the GRADE Working Group. 2016. GRADE Evidence to Decision (EtD) frameworks: a systematic and transparent approach to making well informed healthcare choices. 2: Clinical practice guidelines. *BMJ* 353 (June 2016), i2089. doi:10.1136/bmj.i2089 Publisher: British Medical Journal Publishing Group Section: Research Methods &amp; Reporting.

[3] Pablo Alonso-Coello, Holger J. Schünemann, Jenny Moberg, Romina Brignardello-Petersen, Elie A. Akl, Marina Davoli, Shaun Treweek, Reem A. Mustafa, Gabriel Rada, Sarah Rosenbaum, Angela Morelli, Gordon H. Guyatt, Andrew D. Oxman, and the GRADE Working Group. 2016. GRADE Evidence to Decision (EtD) frameworks: a systematic and transparent approach to making well informed healthcare choices. 1: Introduction. *BMJ* 353 (June 2016), i2016. doi:10.1136/bmj.i2016 Publisher: British Medical Journal Publishing Group Section: Research Methods &amp; Reporting.

[4] Jeffrey C. Andrews, Holger J. Schünemann, Andrew D. Oxman, Kevin Pottie, Joerg J. Meerpohl, Pablo Alonso Coello, David Rind, Victor M. Montori, Juan Pablo Brito, Susan Norris, Mahmoud Elbarbary, Piet Post, Mona Nasser, Vijay Shukla, Roman Jaeschke, Jan Brozek, Ben Djulbegovic, and Gordon Guyatt. 2013. GRADE guidelines: 15. Going from evidence to recommendation-determinants of a recommendation's direction and strength. *Journal of Clinical Epidemiology* 66, 7 (July 2013), 726–735. doi:10.1016/j.jclinepi.2013.02.003

[5] Howard Balshem, Mark Helfand, Holger J. Schünemann, Andrew D. Oxman, Regina Kunz, Jan Brozek, Gunn E. Vist, Yngve Falck-Ytter, Joerg Meerpohl, Susan Norris, and Gordon H. Guyatt. 2011. GRADE guidelines: 3. Rating the quality of evidence. *Journal of Clinical Epidemiology* 64, 4 (April 2011), 401–406. doi:10.1016/j.jclinepi.2010.07.015 Publisher: Elsevier.

[6] Massimo Brunetti, Ian Shemilt, Silvia Pregno, Luke Vale, Andrew D. Oxman, Joanne Lord, Jane Sisk, Francis Ruiz, Suzanne Hill, Gordon H. Guyatt, Roman Jaeschke, Mark Helfand, Robin Harbour, Marina Davoli, Laura Amato, Alessandro Liberati, and Holger J. Schünemann. 2013. GRADE guidelines: 10. Considering resource use and rating the quality of economic evidence. *Journal of Clinical Epidemiology* 66, 2 (Feb. 2013), 140–150. doi:10.1016/j.jclinepi.2012.04.012 Publisher: Elsevier.

[7] Bruno Cartaxo, Gustavo Pinto, Elton Vieira, and Sergio Soares. 2016. Evidence Briefings | Proceedings of the 10th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement. In *ESEM '16: Proceedings of the 10th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement*. Association for Computing Machinery, Ciudad Real, Spain, Article 57, 1–10. doi:10.1145/2961111.2962603

[8] Tore Dybå and Torgeir Dingsøyr. 2008. Strength of evidence in systematic reviews in software engineering. In *Proceedings of the Second ACM-IEEE international symposium on Empirical software engineering and measurement (ESEM '08)*. Association for Computing Machinery, New York, NY, USA, 178–187. doi:10.1145/1414004.1414034

[9] Nicole Forsgren, Margaret-Anne Storey, Chandra Maddila, Thomas Zimmermann, Brian Houck, and Jenna Butler. 2021. The SPACE of Developer Productivity: There's more to it than you think. *Queue* 19, 1 (March 2021), Pages 10:20–Pages 10:48. doi:10.1145/3454122.3454124

[10] Collin Green, Ciera Jaspan, Maggie Hodges, and Jessica Lin. 2024. Developer Productivity for Humans, Part 7: Software Quality. *IEEE Software* 41, 1 (Jan. 2024), 25–30. doi:10.1109/MS.2023.3324830 Conference Name: IEEE Software.

[11] Gordon Guyatt, Andrew D. Oxman, Elie A. Akl, Regina Kunz, Gunn Vist, Jan Brozek, Susan Norris, Yngve Falck-Ytter, Paul Glasziou, Hans deBeer, Roman Jaeschke, David Rind, Joerg Meerpohl, Philipp Dahm, and Holger J. Schünemann. 2011. GRADE guidelines: 1. Introduction—GRADE evidence profiles and summary of findings tables. *Journal of Clinical Epidemiology* 64, 4 (April 2011), 383–394. doi:10.1016/j.jclinepi.2010.04.026 Publisher: Elsevier.

[12] Gordon H. Guyatt, Andrew D. Oxman, Regina Kunz, David Atkins, Jan Brozek, Gunn Vist, Philip Alderson, Paul Glasziou, Yngve Falck-Ytter, and Holger J. Schünemann. 2011. GRADE guidelines: 2. Framing the question and deciding on important outcomes. *Journal of Clinical Epidemiology* 64, 4 (April 2011), 395–400. doi:10.1016/j.jclinepi.2010.09.012

[13] Jo E. Hannay, Tore Dybå, Erik Arisholm, and Dag I. K. Sjøberg. 2009. The effectiveness of pair programming: A meta-analysis. *Information and Software Technology* 51, 7 (July 2009), 1110–1122. doi:10.1016/j.infsof.2009.02.001

[14] Barbara Ann Kitchenham, David Budgen, and Pearl Brereton. 2015. *Evidence-Based Software Engineering and Systematic Reviews* (1 ed.). CRC Press.

[15] Barbara A. Kitchenham, Tore Dyba, and Magne Jorgensen. 2004. Evidence-Based Software Engineering. In *Proceedings of the 26th International Conference on Software Engineering (ICSE '04)*. IEEE Computer Society, Washington, DC, USA, 273–281. http://dl.acm.org/citation.cfm?id=998675.999432

[16] Patrícia G. F. Matsubara and Tayana Conte. 2025. The unfinished business of evidence strength in software engineering: Current practices and future directions. *Empirical Software Engineering* 31, 1 (Oct. 2025), 4. doi:10.1007/s10664-025-10728-9

[17] Patricia Matsubara and Tayana Conte. 2026. Supplementary Material for Bridging the Gap: Adapting Evidence to Decision Frameworks to support the link between Software Engineering academia and industry. https://doi.org/10.6084/m9.figshare.31039837

[18] Titus Winters. 2024. Thoughts on applicability. *Journal of Systems and Software* 215 (Sept. 2024), 112086. doi:10.1016/j.jss.2024.112086

[19] Claes Wohlin and Austen Rainer. 2021. Challenges and recommendations to publishing and using credible evidence in software engineering. *Information and Software Technology* 134 (June 2021), 106555. doi:10.1016/j.infsof.2021.106555

[20] Lanxin Yang, He Zhang, Haifeng Shen, Xin Huang, Xin Zhou, Guoping Rong, and Dong Shao. 2021. Quality Assessment in Systematic Literature Reviews: A Software Engineering Perspective. *Information and Software Technology* 130 (Feb. 2021), 106397. doi:10.1016/j.infsof.2020.106397

[21] Yuan Zhang, Elie A. Akl, and Holger J. Schünemann. 2019. Using systematic reviews in guideline development: The GRADE approach. *Research Synthesis Methods* 10, 3 (2019), 312–329. doi:10.1002/jrsm.1313 _eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1002/jrsm.1313.