

ORIGINAL RESEARCH

Common challenges and suggestions for risk of bias tool development: a systematic review of methodological studies

Eve Tomlinson^{a,*}, Chris Cooper^a, Clare Davenport^{b,c}, Anne W.S. Rutjes^d, Mariska Leeflang^e, Sue Mallett^f, Penny Whiting^a

^aPopulation Health Sciences, Bristol Medical School, University of Bristol, Bristol, UK

^bTest and Prediction Group, Institute of Applied Health Research, University of Birmingham, Birmingham, B15 2TT, UK

^cNIHR Birmingham Biomedical Research Centre, University Hospitals Birmingham NHS Foundation Trust and University of Birmingham, Birmingham B15 2TT, UK

^dDepartment of Medical and Surgical Sciences for Children and Adults (SMECHIMAI), University of Modena and Reggio Emilia, Modena, Italy

^eAmsterdam University Medical Centers, University of Amsterdam, Amsterdam, The Netherlands

^fCentre for Medical Imaging, University College London, London, UK

Accepted 17 April 2024; Published online 24 April 2024

Abstract

Objectives: To review the findings of studies that have evaluated the design and/or usability of key risk of bias (RoB) tools for the assessment of RoB in primary studies, as categorized by the Library of Assessment Tools and Instruments Used to assess Data validity in Evidence Synthesis Network (a searchable library of RoB tools for evidence synthesis): Prediction model Risk Of Bias ASessment Tool (PROBAST), Risk of Bias-2 (RoB2), Risk Of Bias In Non-randomised Studies of Interventions (ROBINS-I), Quality Assessment of Diagnostic Accuracy Studies-2 (QUADAS-2), Quality Assessment of Diagnostic Accuracy Studies-Comparative (QUADAS-C), Quality Assessment of Prognostic Accuracy Studies (QUAPAS), Risk Of Bias in Non-randomised Studies of Exposures (ROBINS-E), and the CONsensus-based Standards for the selection of health Measurement INstruments (COSMIN) RoB checklist.

Study Design and Setting: Systematic review of methodological studies. We conducted a forward citation search from the primary report of each tool, to identify primary studies that aimed to evaluate the design and/or usability of the tool. Two reviewers assessed studies for inclusion. We extracted tool features into Microsoft Word and used NVivo for document analysis, comprising a mix of deductive and inductive approaches. We summarized findings within each tool and explored common findings across tools.

Results: We identified 13 tool evaluations meeting our inclusion criteria: PROBAST (3), RoB2 (3), ROBINS-I (4), and QUADAS-2 (3). We identified no evaluations for the other tools. Evaluations varied in clinical topic area, methodology, approach to bias assessment, and tool user background. Some had limitations affecting generalizability. We identified common findings across tools for 6/14 themes: (1) challenging items (eg, RoB2/ROBINS-I “deviations from intended interventions” domain), (2) overall RoB judgment (concerns with overall risk calculation in PROBAST/ROBINS-I), (3) tool usability (concerns about complexity), (4) time to complete tool (varying demands on time, eg, depending on number of outcomes assessed), (5) user agreement (varied across tools), and (6) recommendations for future use (eg, piloting) and development (add intermediate domain answer to QUADAS-2/PROBAST; provide clearer guidance for all tools). Of the other eight themes, seven only had findings for the QUADAS-2 tool, limiting comparison across tools, and one (“reorganization of questions”) had no findings.

Conclusion: Evaluations of key RoB tools have posited common challenges and recommendations for tool use and development. These findings may be helpful to people who use or develop RoB tools. Guidance is necessary to support the design and implementation of future RoB tool evaluations. © 2024 The Author(s). Published by Elsevier Inc. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

Keywords: Risk of bias; Systematic reviews; Evaluation; Research methods; RoB; Quality assessment

1. Introduction

Systematic reviews seek to answer a specific research question by summarizing results from primary studies that meet prespecified eligibility criteria. Results from primary studies may be biased (ie, deviate from the truth) due to inappropriate design, conduct, analysis, and/or reporting.

Funding: This research did not receive any specific grant from funding agencies in the public, commercial, or not-for-profit sectors.

* Corresponding author. Population Health Sciences, Bristol Medical School, University of Bristol, Canynge Hall, 39 Whatley Road, Bristol, BS8 2PS, UK.

E-mail address: eve.tomlinson@bristol.ac.uk (E. Tomlinson).

<https://doi.org/10.1016/j.jclinepi.2024.111370>

0895-4356/© 2024 The Author(s). Published by Elsevier Inc. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

What is new?

Key findings

- Thirteen studies have evaluated the design and usability of risk of bias (RoB) tools: PROBAST, RoB2, ROBINS-I, and QUADAS-2. The studies highlight common challenges with tool implementation and provide recommendations for tool use and development. We identified no primary studies evaluating QUADAS-C, QUAPAS, ROBINS-E, or the COSMIN RoB checklist.
- Some of the existing tool evaluation studies have methodological limitations restricting the generalizability of their findings.

What this adds to what was known?

- This is the first review to summarize findings from studies evaluating the design and usability of RoB tools.

What is the implication and what should change now?

- Our findings will be helpful to tool users and developers who may benefit from other experiences of applying RoB tools and recommendations for improving the RoB assessment process. Anyone planning to evaluate a RoB tool may also benefit from our review of the current literature on tool evaluation, evaluation methodology used, and limitations of existing research.
- Guidance is necessary to ensure appropriate design and implementation of RoB tool evaluation studies in future.

A review of biased results is likely to provide a misleading conclusion. This is an issue, as conclusions from systematic reviews are widely used to inform clinical, policy, and patient decision-making. Guidance for systematic review conduct therefore emphasizes the importance of completing risk of bias (RoB) assessment in systematic reviews [1], and reporting guidelines (such as the Preferred Reporting Items for Systematic reviews and Meta-Analyses (PRISMA) 2020 [2]) recommend that reviewers transparently report RoB assessments. This information contributes to certainty in the overall evidence [3].

Several tools exist to facilitate the assessment of RoB for varying study designs [4,5]. The Library of Assessment Tools and Instruments Used to assess Data validity in Evidence Synthesis (LATITUDES) Network provides a searchable library of such tools [6]. LATITUDES identifies RoB tools for inclusion in the library through liaison with experts interested in evidence synthesis or RoB

assessment, screening existing depositories of tools, screening reviews that have focused on tools, and a review of guidance from evidence synthesis organizations. LATITUDES listed eight key tools designed to assess the RoB in primary studies at launch in September 2023: Quality Assessment of Diagnostic Accuracy Studies-2 (QUADAS-2) for diagnostic accuracy studies [7], Risk of Bias-2 (RoB2) for randomized controlled trials [8], Risk Of Bias In Non-randomised Studies of Interventions (ROBINS-I) for nonrandomized studies of interventions [9], Prediction model Risk Of Bias Assessment Tool (PROBAST) for diagnostic and prognostic prediction model studies [10], Quality Assessment of Diagnostic Accuracy Studies-Comparative (QUADAS-C) for comparative diagnostic accuracy studies [11], Quality Assessment of Prognostic Accuracy Studies (QUAPAS) for prognostic accuracy studies [12], Risk Of Bias in Non-randomised Studies of Exposures (ROBINS-E) for exposure studies [13], and the Consensus-based Standards for the selection of health Measurement Instruments (COSMIN) RoB checklist for studies on measurement properties of patient-reported outcome measures [14]. These tools fulfill specific criteria posited by the LATITUDES Network [6]; they (1) focus on RoB, or distinction between items that assess RoB and other aspects of study quality; (2) offer a method to reach either a domain-specific or overall assessment of RoB; (3) have been used in at least one review that none of the tool authors were co-authors on or are an update to a previously recommended LATITUDES key tool; (4) have been developed involving collaborators from different disciplines (eg, methodologists, statisticians, clinicians); and (5) avoid use of summary numerical quality scores.

Reporting guidelines also exist for many of the study designs previously outlined, for example, the Transparent Reporting of a multivariable prediction model for Individual Prognosis Or Diagnosis (TRIPOD) statement [15], the CONSolidated Standards Of Reporting Trials (CONSORT) statement [16], the Strengthening The Reporting of OBServational studies in Epidemiology (STROBE) statement [17], and the STAndards for Reporting of Diagnostic accuracy (STARD) statement [18]. These resources differ to RoB tools—reporting guidelines are checklists of information used to guide authors of primary studies in reporting a specific type of research, whereas RoB tools are used by systematic reviewers to assess the RoB of included primary studies. However, the goal of using reporting guidelines is to improve the quality of the reporting of studies [19], which in turn supports systematic reviewers in conducting RoB assessments.

Studies that have evaluated RoB tools have not previously been summarized in a review. Therefore, this systematic review of methodological studies aims to appraise the findings of studies that have evaluated the design and/or usability of the following RoB tools: PROBAST, RoB2, ROBINS-I, QUADAS-2, QUADAS-C, QUAPAS, ROBINS-E, and the COSMIN RoB checklist. A summary

of the structure (eg, number of domains, number of signaling questions) and features (eg, tool purpose, level of assessment, process for generating an overall RoB rating) of these RoB tools is provided in [supplementary material section 2](#). The tools have many similarities but also some differences, for example, in level of assessment, domain and signaling question answer options, and software to support use. In this review, we highlight challenges faced by tool users, report suggestions for tool development, and discuss tool evaluation methodology.

2. Methods

This review was conducted in line with published guidance on systematic review conduct [1], and is reported according to PRISMA 2020 statement ([supplementary material section 6](#)), with adaptations made for the methodological nature of this review [20]. The protocol was registered on Open Science Framework [21].

2.1. Eligibility criteria

Studies were eligible for inclusion if they stated in the title or abstract, that they aimed to evaluate the design (eg, the number of domains, signaling questions or answer options) and/or usability of the most recent version of any of the following tools: QUADAS-2 [7], RoB2 [8], ROBINS-I [9], PROBAST [10], QUADAS-C [11], QUAPAS [12], ROBINS-E [13], or the COSMIN RoB checklist [14]. We included studies with quantitative or qualitative findings and did not restrict based on publication language or clinical topic area.

2.2. Search strategy

To identify evaluation studies, we undertook a forward citation search in Science Citations Index Expanded (SCI-Expanded, Clarivate, 1990-current) using the primary report of each tool as the source [22]. Tools were searched in two phases as the LATITUDES Network website was developed. In July 2023, we searched for QUADAS-2, PROBAST, ROBINS-I, and RoB2 [7–10]. This included a Google Scholar search to identify eligible evaluations not yet included in SCI-Expanded. In October 2023, we searched the remaining tools added to the LATITUDES Network website, namely QUADAS-C, QUAPAS, and the COSMIN RoB checklist [11,12,14]. ROBINS-E [13] had not been formally published but we checked the recommended citation for this tool. We also checked the bibliographies of all included evaluation studies.

2.3. Study selection

Two reviewers (E.T. and C.C.) undertook study selection. The reviewers independently screened titles and

abstracts of reports identified by the searches, using the web-hosted screening tool Rayyan. Full reports of those considered potentially relevant were obtained and assessed for inclusion. Any potentially relevant reports identified by checking the reference lists of included studies were also assessed for inclusion. Any disagreements were resolved by consensus. No automation tools were used.

2.4. Data extraction

We extracted the following study characteristics from each evaluation study using NVivo [23]: first author, year of publication, country of first author, journal, title, aim, model/study sampling method, details of models/studies assessed (including clinical topic area, study design, number of studies/models), number of individuals using tool, user background, user experience with tool, and guidance used. One reviewer conducted data extraction (E.T.), and where there was uncertainty, this was checked by a second reviewer (C.D.). The team met at regular intervals to discuss findings from data extraction.

2.5. Data synthesis

We tabulated a summary of the included evaluation studies' characteristics and methodology.

We conducted document analysis of the included evaluation studies. One reviewer (E.T.) used NVivo to extract data: relevant information was highlighted from study reports and assigned to themes. We began with a deductive approach and categorized information into the coding structure prespecified in our protocol ([Table 1](#)). This was developed by the review team, based on what we sought to find in the literature. When we found additional relevant information that did not fit into these themes, we moved to an inductive approach and added new themes as necessary, following discussion among the review team.

Following document analysis in NVivo, the information was exported into Microsoft word, providing a list of themes and associated quotes from each study. This information constituted our findings from each included evaluation study. We used this to produce (1) a narrative summary of findings for each study within each included tool and (2) a narrative summary of findings across tools, highlighting common findings. We tabulated any quantitative data.

2.6. Differences between protocol and review

For completeness, we included all key tools designed to assess the RoB in primary studies listed on the LATITUDES Network website at launch in September 2023 rather than restrict to those considered by LATITUDES before March 2023, as was noted in our protocol.

Table 1. Themes used to organize data

General findings	Overall usability
	Time to complete the tool
	Experience of generating overall risk of bias judgments
	User agreement
	Floor and ceiling effects ^a
Domain and signaling question specific findings	Smallest detectable change ^a
	Standard error of measurement ^a
	Internal consistency ^a
	Items to remove
	Items that are described as being challenging
Future use/development	Items to reword
	Items to add (not study-specific)
	Suggestions for reorganizing signaling questions, for example, into different domains
Future use/development	Recommendations relating to the use or future development of the tool

^a These themes were added using an inductive approach.

3. Results

3.1. Search results

Searches identified 19,512 unique reports. Of these, 13 studies (17 reports) were included in the review (Fig 1).

This included three studies for PROBAST, three for RoB2, three for QUADAS-2, and four for ROBINS-I. We identified no evaluation studies for QUADAS-C, QUAPAS, ROBINS-E, or the COSMIN RoB checklist. References for included studies and studies excluded at full-text assessment are presented in [supplementary material section 1](#), stratified by tool.

3.2. Characteristics of included evaluation studies

Characteristics of the 13 included evaluation studies are presented in [supplementary material section 3](#). Studies were conducted by first authors in Australia [1,24], Brazil [1,25], Canada [1,26], China [1,27], Germany [1,28], Italy (three; led by the same first author) [29–31], Spain [1,32], the Netherlands [2,33,34], and the United Kingdom [2] [35,36].

Most studies focused on one clinical topic area, including cardiology, physical therapy, Alzheimer's disease, uterine cervix examination, cannabis/cannabinoids for multiple sclerosis, melanoma, independence in older people, and the association between depression and risk of stroke. One study sampled studies focusing on three topics: vaccines, opiate abuse, and rehabilitation [31]. Three studies did not restrict to specific topics—one evaluated RoB2 [29], one PROBAST [34], and one ROBINS-I [27].

Sampling methodology varied across evaluation studies: six sourced studies from one or more systematic reviews known to, or in production by, the author team conducting the evaluation study [24,30–32,36,37]; three studies searched databases to identify published systematic reviews from which to sample primary studies [25,27,34]; three

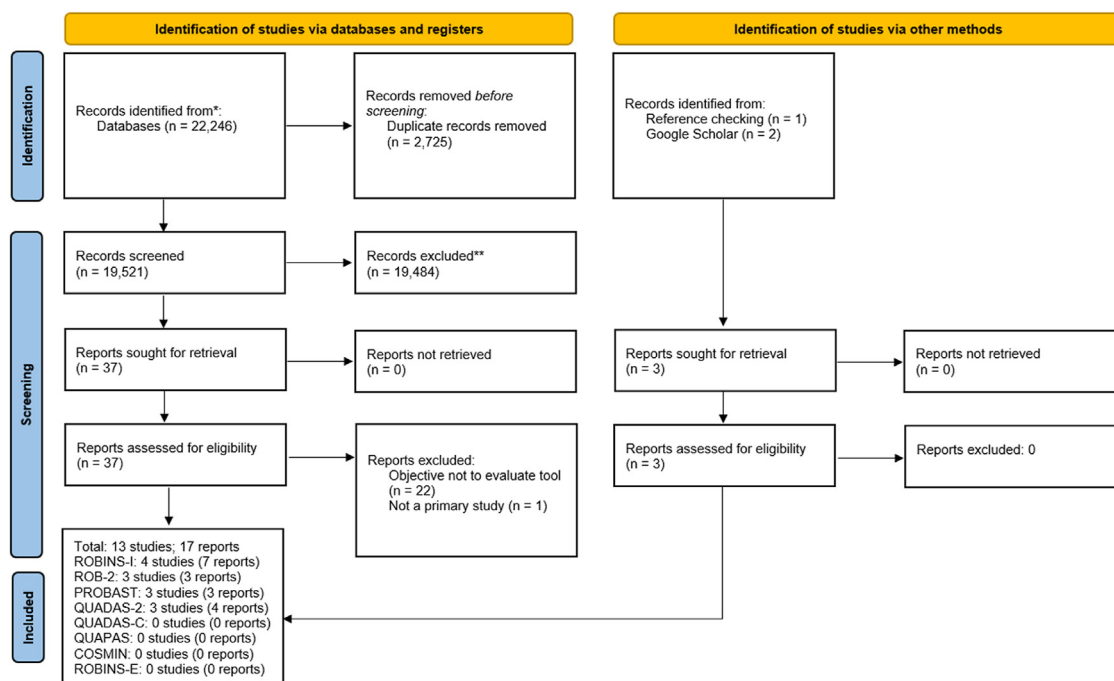


Figure 1. PRISMA diagram showing study selection process for all included risk of bias tools. (For interpretation of the references to color in this figure legend, the reader is referred to the Web version of this article.)

searched databases to identify published studies/models [26,29,33]; and one sampled all melanoma risk prediction studies published until 2021 [28].

In most of the included evaluation studies, two or more individuals independently applied the RoB tool. In one evaluation study, one person completed the RoB assessment and it was checked by another person [37]. In two studies, the authors appraised RoB assessments conducted in other studies: one sampled systematic reviews that had used PROBAST [34] and one sampled meta-analyses that had used QUADAS-2 [25].

Where reported, all study authors used published guidance associated with the tool. One study was conducted as part of the Cochrane RoB2 pilot project; therefore, the authors were in contact with the tool developers [30]. Three studies involved the creation and use of extra project-specific guidance about RoB assessment [26,28,30].

The amount of information reported about the assessors' background and experience varied. Most had research experience but had not used the RoB tool under evaluation. Four studies reported having piloted the tool before beginning the study [24,27,32,37]. Two studies reported that assessors had clinical knowledge relevant to the topic area [24,31]. In other studies, it was either not reported or some assessors had clinical backgrounds but it was not clear if it was relevant to the assessment topic. Cross-checking the author lists of the included studies with the developers of the included RoB tools revealed that none of the evaluation studies were conducted by tool developers.

3.3. Findings from included evaluation studies

Table 2 presents common findings across the RoB tools for which we identified evaluation studies. Study findings should be interpreted in light of the study characteristics previously summarized and tabulated in [supplementary material section 3](#), as the studies varied in clinical topic area, number of individuals using the tool, methodology, and user background. Further detail, including user agreement at the domain and signaling level and all coded study findings by theme for each tool, is provided in [supplementary material section 4](#) and [section 5](#). Additionally, as previously noted, an overview of the included RoB tools, and the similarities and differences between them, can be found in [supplementary material section 2](#).

We identified common findings across tools for six of 14 themes (Table 2): challenging items, overall RoB judgment across domains, tool usability, time to complete tool, user agreement, and future use and development of tools. Of the other eight themes, six only had findings for the QUADAS-2 tool, limiting comparison across tools, and we identified no findings for the theme "reorganization of questions".

3.3.1. Specific tool items

The domain "deviations from intended interventions" was found to be difficult to complete in two studies of

RoB2 [27,31] and two studies of ROBINS-I [29,36] (Table 2).

3.3.2. Design of tool

Tool users raised concerns about tool ceiling effects and the way that overall RoB is calculated within PROBAST [33,34] and ROBINS-I [27], noting that the tool structure does not allow users to distinguish between studies of varying poor quality, as studies will be rated at high RoB overall, whether they are high risk in one or all domains.

3.3.3. Tool usability

Many studies reported that the RoB tools were challenging and time-consuming to use, due to their length and complexity. The time it took to implement the tools varied and depended on the number of outcomes being assessed and on the tasks included in measurement (eg, including consensus discussions or not). It reduced with the use of project-specific guidance and training [26,28,30].

User agreement for overall RoB varied across tools. QUADAS-2 [24] and RoB2 [29,30] had low inter-rater agreement. It varied across studies for PROBAST [28,33] and ROBINS-I [26,27,31,32]. Inter-rater agreement at the domain and signaling question level for each tool is outlined in [supplementary material section 5](#), showing variation between studies.

3.3.4. Training and guidance

Five studies of PROBAST, ROBINS-I, and RoB2 noted that the guidance provided by tool developers lacks clarity [26,29,31–33,36]. Studies called for clearer guidance [24,27,31,36,37], and simplified versions of the tools [27,31,33]. Three studies recommended that QUADAS-2 and PROBAST should include an intermediate domain-level answer option between low and high risk [33,34,37].

Studies (all tools) commonly posited two pieces of advice for tool users: (1) form a team with methodological, statistical, and clinical expertise [28–31,37] and (2) conduct preparatory work before beginning RoB assessments, such as training, piloting the tool, and creating project-specific RoB guidance [26–28,31,32,34].

4. Discussion

4.1. Summary of main results

Thirteen studies have evaluated the design and usability of the PROBAST, RoB2, ROBINS-I, and QUADAS-2 RoB tools. These constitute four of eight key RoB tools specified by the LATITUDES Network. No study evaluations were identified for the other tools: QUADAS-C, QUAPAS, ROBINS-E, and the COSMIN RoB checklist. With the exception of the COSMIN RoB checklist, this may be because they were created more recently.

Table 2. Findings common to more than one RoB tool

Theme	Key finding	Tool	Details from studies
Challenging items	The domain “deviations from intended interventions” is challenging to complete	ROBINS-I	<ul style="list-style-type: none"> • Challenging to interpret the signaling questions in this domain [31]. • Most time-consuming domain; struggled to decipher if an intervention was a cointervention, whether the cointervention was important and whether it resulted in substantial imbalances between groups [27].
		RoB2	<ul style="list-style-type: none"> • Barriers to assessing domain included difficulty in deciding when deviations from the intended intervention did not reflect usual practice and whether they were likely to have affected the outcome, new terminology and approach used in the tool, poor clarity of the guidance, and lack of subject matter and statistics knowledge [29]. • Uncertainty regarding what evidence should be considered to indicate “probably no” to SQ2.3 which asks whether there are deviations from the intended intervention that arose because of the trial context [36].
		PROBAST	<ul style="list-style-type: none"> • Guidance leaves questions open to interpretation which complicates assessment [33].
Guidance lacks clarity		ROBINS-I	<ul style="list-style-type: none"> • Lengthy detailed guidance may contribute to limited clarity in overall bias assessment [26]. • Insufficient instructions in guidance on decision trees of conditional signaling questions when the answer to the previous question is “no information” may contribute to poor domain-level agreement between assessors [31]. • Clarity of instructions and items in the tool rated as “poor” [32].
		RoB2	<ul style="list-style-type: none"> • Lack of clarity in the guidance concerning terminology and approach introduced since RoB1 contributed to difficulties in assessment. Guidance is an improvement to the original version of the tool, but it is complex and its application is demanding [29]. • Extensive guidance but in some places it lacks specificity sufficient to operationalize it, and some parts are unnecessarily discursive and theoretical with insufficient practical advice for interpreting the signaling questions [36].
		PROBAST	<ul style="list-style-type: none"> • Add intermediate option to facilitate “less stringent assessment” [33]. • Add intermediate option to help distinguish between RoB in studies [34].
Future use and development of tool – suggestions for tool developers	Add intermediate domain-level answer option between “low” and “high” risk (QUADAS-2; PROBAST)	QUADAS-2	<ul style="list-style-type: none"> • Add intermediate option to promote consistency [37].
		ROBINS-I	<ul style="list-style-type: none"> • Suggestions to make guidance clearer included integrate practical examples of biases, and add more examples from different fields of medicine to clarify the questions and highlight situations where actual bias may arise [27,31].
		RoB2	<ul style="list-style-type: none"> • Refine the guidance with a focus on operationalizing the tool, for example, (1) add specific examples directly related to

(Continued)

Table 2. Continued

Theme	Key finding	Tool	Details from studies
			signaling questions and more examples of “judgment calls” rather than extremes; (2) provide greater emphasis on the application of, and location of dividing lines for, each signaling question; (3) move some of the theoretical background and empirical evidence to an appendix; and (4) provide more guidance regarding a suitable cut-off point for judging whether nearly all data are available in a study [36].
		QUADAS-2	<ul style="list-style-type: none"> • Give specific guidance on what constitutes low/high domain-level bias [37]. • Give clearer instructions on how to assess studies including multiple index tests, reference standards, or pathologies [24].
	Produce simplified versions of tool	PROBAST	<ul style="list-style-type: none"> • Shorter version of the tool would be useful, particularly when conducting high-volume assessments where the classification into high/low risk is the main objective. Note: the evaluation study authors created a short version in this study [33].
		ROBINS-I	<ul style="list-style-type: none"> • Tool is overcomplicated and a simpler version is needed [27]. • Number of signaling questions should be reduced [31].
Future use and development of tool – suggestions for tool users	Form a team with relevant experience and knowledge	PROBAST	<ul style="list-style-type: none"> • Include experienced epidemiologists specialized to the area of research and involve them to create valid decision rules for RoB ratings [28].
		ROBINS-I	<ul style="list-style-type: none"> • Form a team containing methodology, statistics, and clinical expertise [31].
		RoB2	<ul style="list-style-type: none"> • Ensure team has clinical, methodological, and statistical knowledge [29,30].
		QUADAS-2	<ul style="list-style-type: none"> • Ensure team includes clinicians to assist with applicability assessment [37].
	Conduct adequate preparatory work before starting RoB assessments including training on how to use the tool and creating project-specific guidance	PROBAST	<ul style="list-style-type: none"> • Create customized guidance for the use of the tool specific to the topic of study (and present this in the systematic review report) and undertake disease-specific and study-type –specific training [28]. • Complete specialized training in the tool before use [34].
		ROBINS-I	<ul style="list-style-type: none"> • Undertake training on the use of the tool, have a preliminary discussion on how to apply the tool using a common approach, and pilot assessment [31]. • Undertake training in the use of the tool and develop an implementation document containing clear decision rules customized to project, to agree on how the assessments will be conducted before beginning [26]. • Create a list of confounders, including time-varying confounding and cointervention, to save assessment time [27,32].
		RoB2	<ul style="list-style-type: none"> • Undertake formal training on how to use tool and have a preliminary discussion on how the tool will be applied, and pilot tool [29]. • Agree on the validity and appropriateness of outcome measures at protocol stage, create

(Continued)

Table 2. Continued

Theme	Key finding	Tool	Details from studies
			<p>an implementation document tailored to review with instructions on how to answer each signaling question and a list of possible co-interventions that could lead to bias (to help with the “deviation from intended interventions” domain), and pilot tool [30].</p> <ul style="list-style-type: none"> • Develop specific guidance tailored to the review to help to overcome challenges and ensure consistency across assessments [36].
Floor and ceiling effects	Presence of floor and/or ceiling effects in the tools	PROBAST	<ul style="list-style-type: none"> • Ceiling effect identified [34] – Users unable to distinguish between studies of varying poor quality, for example, a study scoring poorly in one domain will have the same overall score as a study scoring poorly in multiple domains [33,34].
		ROBINS-I	<ul style="list-style-type: none"> • If all domains are judged as “moderate”, the overall judgment will be “moderate”, which may lead to inaccurate results/cover up a serious study issue [27].
		QUADAS-2	<ul style="list-style-type: none"> • Floor and ceiling effects not identified in the tool [24].
Overall RoB judgment	Crude categorization of overall RoB	PROBAST	<ul style="list-style-type: none"> • Concerns about overall RoB due to stringent assessment and lack of “intermediate” answer option [33,34].
		ROBINS-I	<ul style="list-style-type: none"> • Discriminative ability of the tool is limited by the “crude” categorization of overall RoB as “low,” “moderate,” “serious,” or “critical” [27].
	Tailor how signaling questions contribute to overall judgment	QUADAS-2	<ul style="list-style-type: none"> • Useful to tailor how signaling questions within a domain contribute to the overall domain-level RoB judgment. This provides the opportunity to weight signaling questions to topic specific sources of bias [37].
Overall usability	Tools are complex to use	ROBINS-I	<ul style="list-style-type: none"> • Tool is too comprehensive to provide a concise critical appraisal [27].
		RoB2	<ul style="list-style-type: none"> • Tool and guidance are comprehensive and although it is an improvement on the last version of the tool, it is complex and demanding to apply [29]. • Tool is complex, particularly the signaling questions [30]. • Conducting assessments with RoB2 is a substantial and challenging undertaking [36].
Time to complete tool	Large time investment to conduct RoB assessments but this varies depending on factors such as the number of outcomes and reports assessed and the individuals conducting assessments. Time to apply tool was mostly found to reduce with use of project-specific guidance and training	PROBAST	<ul style="list-style-type: none"> • Estimated that applying tool took less than one hour [33]. • Before training sessions (but after having read the tool and guidance document), the first 20 studies had 8 h of consensus discussions (2 × 4 h meetings). Following training/customized guidance, the next 22 studies had nine hours of consensus discussions (6 × 1.5 h meetings) [28].
		ROBINS-I	<ul style="list-style-type: none"> • Mean time to apply tool (read article and assess relevant outcomes) was 27.8 min (standard deviation [SD] 12.6) [31].

(Continued)

Table 2. Continued

Theme	Key finding	Tool	Details from studies
			<ul style="list-style-type: none"> • Mean time to assess RoB and reach a consensus was 48.45 min (95% CI 45.61–51.29). It reduced by 12.9 min (95% CI –16.4 to –9.4) after guidance/training [26]. • Time to assess a single study (excluding time to read study) reduced from 7 to 3 h as familiarity increased. Time to reach consensus reduced from 40 min to 14 min [27].
		RoB2	<ul style="list-style-type: none"> • Mean time to apply tool (one outcome per study; not using Excel tool; not including consensus decisions) was 28 min (SD 13.4) [29]. • Mean time to apply tool (multiple outcomes per study; not reported what was included) was 168.5 min (SD: 68.7). This reduced to 41 min (SD: 18.39) when using an implementation document, which itself took 40 h over the course of 3 mo to develop [30]. • Mean time for an individual assessment of a study (peer reviewer, with multiple results and reports) was 127 min (2 h 7 min; SD 67) and 54 min (SD 43) for the consensus meeting. In total (for two complete individual assessments and a consensus meeting), the resource for the overall process was 358 min (SD 183). For 99 studies, it took 35,472 min of worktime (591.2 person-hours or 2.1 mo of 2× FTE work) including each individual assessment and two people in a consensus meeting (5 h 58 min per study, 47 min per result). Number of results and reports per study increased the time to conduct assessments and the former also increased time for consensus meetings. There was substantial variation between individual reviewers and experience reduced time taken to conduct individual assessments and consensus meetings. For two reviewers who had already previously assessed at least 25 studies with RoB2, the overall worktime to assess 32 studies was 5 h 15 min per study or 44 min per result on average. Regression analysis estimated this to be 178 minutes per study (2 h 58 min; 95% CI 139–218) plus 19 min per result (95% CI 15–24) (adjusted R 2.73) [36].
		QUADAS-2	<ul style="list-style-type: none"> • More time-consuming than original version of the tool because more free text was required, but this was helpful for forming overall judgments [37].
User agreement – overall RoB	Low inter-rater agreement for overall RoB in QUADAS-2 and RoB2; varying agreement for PROBAST and ROBINS-I	PROBAST	<ul style="list-style-type: none"> • 90% agreement between two assessors before final consensus meeting (Kappa 0.33) [33]. • Poor agreement before customized training (mean pairwise AC1: 0.098/mean pairwise Cohen's k: 0.132, multirater AC1: 0.071), then moderate after training (mean pairwise AC1: 0.474/mean pairwise Cohen's k: 0.261, multirater AC1: 0.476) [28].
		ROBINS-I	<ul style="list-style-type: none"> • Moderate agreement (PABAK 0.57) [32].

(Continued)

Table 2. Continued

Theme	Key finding	Tool	Details from studies
			<ul style="list-style-type: none"> • Slight agreement (Kappa 0.06 [95% CI 0.001–0.12]) [31]. • Substantial agreement (AC1 statistic 95% CI 0.61 [0.42–0.79]) [27]. • Poor agreement before implementation document (AC1 statistic 95% CI 0.00 [0.00–0.18]), then fair agreement after implementation document used (AC1 statistic 95% CI 0.38 [0.18–0.58]) [26].
		RoB2	<ul style="list-style-type: none"> • Slight agreement (Kappa 0.16, 95% CI 0.08–0.24) [29]. • No agreement (Kappa –0.15) before implementation document used, then moderate agreement (Kappa 0.42) after implementation document used [30].
		QUADAS-2	<ul style="list-style-type: none"> • Poor intertester reliability of summary scores (ICC = 0.36; 95% CI 0.08–0.59) [24].

Note: Study characteristics including how many people used the tool, their experience, and topics of studies assessed are outlined in [supplementary material section 3](#).

CI, confidence interval; ICC, intraclass correlation coefficient; RoB, risk of bias.

We have shown that there are common findings from evaluation studies across RoB tools concerning tool design and usability (eg, concerns with tool length and complexity, overall bias calculations, and varying user agreement), and recommendations for tool users (eg, conduct adequate preparatory work prior to RoB assessment) and tool developers (eg, produce clearer guidance and simplify tools).

4.2. Strengths and limitations

This paper is the first to provide a review of evaluations of RoB tools, limiting our ability to interpret the findings within the context of other available evidence. Strengths of this review include that it followed a preregistered protocol and adhered to prespecified inclusion criteria. The author team includes tool developers and tool users; therefore, we had a robust understanding of how the tools are intended to be used. We used objective tool inclusion criteria and members of the team who have been involved in the development of included tools did not undertake screening, data extraction, or synthesis. Nonetheless, we acknowledge that our prior experience may have influenced how we interpreted results.

A potential limitation of this review is that studies were only included if they had reported an objective to evaluate one of the included RoB tools within the title or abstract. We may therefore have missed studies that evaluated a tool but not as a main objective. Additionally, we may have missed evaluations by focusing only on published work. In future, it would be useful to supplement this review by

interviewing author teams who have used RoB tools and by contacting tool developers for their own tool evaluations. It may also be insightful to explore evaluations of a wider sample of RoB tools. This review focused on eight tools specified as key tools by the LATITUDES Network; however, these are only a handful of those available [4,6].

Our findings should be interpreted within the context of the characteristics of the included evaluation studies, which varied in methodology, approach to bias assessment, clinical topic area, and tool user background. The challenges faced by tool users in the studies may be topic-dependent and therefore the generalizability of our findings may be limited. However, we have included a range of clinical topic areas in this review.

We identified limitations in the evaluation methodology employed by included studies. Not all evaluations were conducted in “real life” settings, for example, one team received support from tool developers as part of a Cochrane RoB2 pilot and some other studies assessed only one outcome per study, which may not often be the case in practice. One study that evaluated ROBINS-I used a sample of cohort studies of exposure rather than nonrandomized studies of interventions, for which the tool is intended to be used [32]. Additionally, experience and topic expertise varied across evaluation teams making generalizability difficult to judge, and the design of the evaluations differed between studies. Some studies focused on one clinical topic area and some spanned multiple areas, with one noting that the generalizability of its findings was limited due to the inclusion of 19 clinical specialities [29]. However, despite this variation, themes emerged across

included studies, adding confidence to the validity and applicability of our findings.

4.3. Implications

The findings of this review will be useful to individuals planning to develop, update, use, or evaluate a RoB tool. Tool developers may benefit from incorporating the findings of existing evaluations into tool updates, as is currently being done by members of the present author team when updating QUADAS-2 into QUADAS-3. Tool developers may also find it helpful to appraise the similarities and differences between these widely used tools, with a view to aligning (where relevant) with other tools, to support reviewers who are likely to use more than one tool. We have shown that there are common themes across tools that could be addressed as part of aligning key features in design and support/guidance available for tool users. Additionally, as this review has shown that the process of using RoB tools can be complicated and RoB tool guidance can lack clarity, tool users may benefit from reading our summary of the recommendations for use of the specific tools given by the authors of the included evaluation studies (eg, to develop review-specific guidance to assist the team in interpreting the generic RoB tool guidance). Finally, anyone seeking to evaluate a RoB tool in the future may be keen to learn from the methodology used in existing evaluation studies and the limitations of previous work.

To support future evaluations of RoB tools, guidance on good practice for tool evaluation is needed. Learning from the limitations outlined in existing evaluation studies, we suggest a few starting points to consider within tool evaluations. First, consider clinical topic area. While it is useful to apply the tool to studies from a range of topics to improve the generalizability of results, it is also important to reduce noise that is not related to tool features and usability. Therefore, it may be useful for authors of evaluation studies to limit to one or a few topics. Second, consider the clinical and RoB assessment experience of the research team and try to mirror the “real life” use of the tool as much as possible. Finally, consider which outcomes and measures are most helpful to the evaluation. A core set of outcomes could be developed to improve uniformity across evaluation studies, thus improving our ability to compare study findings.

5. Conclusion

Studies have evaluated the design and usability of RoB tools, highlighting common challenges and recommendations for tool use and development. These findings may be useful for people using, developing, or evaluating tools. Guidance is needed to encourage appropriate design and implementation

of RoB tool evaluation studies. Such guidance could help to enhance the quality of future evaluations of RoB tools.

CRedit authorship contribution statement

Eve Tomlinson: Writing – review & editing, Writing – original draft, Visualization, Project administration, Methodology, Investigation, Formal analysis, Conceptualization. **Chris Cooper:** Writing – review & editing, Methodology, Conceptualization. **Clare Davenport:** Writing – review & editing, Supervision, Methodology, Conceptualization. **Anne W.S. Rutjes:** Writing – review & editing, Supervision, Methodology, Conceptualization. **Mariska Leeflang:** Writing – review & editing, Supervision, Methodology, Conceptualization. **Sue Mallett:** Writing – review & editing, Supervision, Methodology, Conceptualization. **Penny Whiting:** Writing – review & editing, Supervision, Methodology, Conceptualization.

Data availability

The full summarized data are provided in supplementary material. A data extraction sheet containing the verbatim quotes extracted from included studies is available upon reasonable request.

Declaration of competing interest

Eve Tomlinson: Since working on this review, Eve has joined the QUADAS-3 steering group.

Chris Cooper: None.

Clare Davenport: Clare Davenport is a member of the Diagnostic test Accuracy Editorial Team and is an author on the Cochrane Handbook for Systematic Reviews of Diagnostic Test Accuracy Studies. Clare was a member of the QUADAS-2 working group and a member of the QUADAS-C steering group. She is currently a member of the QUADAS-3 steering group.

Anne W.S. Rutjes: Anne is regularly hired by the Federal office of Public Health in Bern, Switzerland as a methodological consultant for Health Technology Assessment reports and receives teaching fees from the academic sector on topics related to research designs, risk of bias tools, and reporting tools. Anne was an editor for the Cochrane Dementia and Cognitive Improvement group up to March 2023 and is an author on the Cochrane Handbook for Systematic Reviews on Diagnostic Test Accuracy Studies. She was a panel member of PROBAST and various COSMIN tools, part of the steering group of QUADAS and QUADAS-2, is part of the core group developing a RoB tool for prevalence studies, and one developing a RoB and reporting tool for use in studies evaluating the incidence of adverse events detected with record review methods.

Mariska Leeflang: Mariska Leeflang is an author of the following risk of bias tools: QUADAS-C, QUADAS-2, and PROBAST. Mariska is an editor of the Cochrane Diagnostic Test Accuracy Handbook.

Sue Mallett: Sue Mallett is an author of the following risk of bias tools: QUADAS-C, QUADAS-2, and PROBAST. She is also on the leadership board for the LATITUDES Network. No payments are associated with this.

Penny Whiting: Penny is an author on the Cochrane Handbook for Systematic Reviews on Diagnostic Test Accuracy Studies. She led the development of QUADAS and QUADAS-2, and was a steering group member for the development of PROBAST. She contributed to the development of RoB 2 and ROBINS-I. She is also on the leadership board for the LATITUDES Network.

Acknowledgments

We thank Amanda Owen-Smith at the University of Bristol for providing guidance concerning qualitative analysis.

Supplementary data

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.jclinepi.2024.111370>.

References

- Higgins JP, Thomas J, Chandler J, Cumpston M, Li T, Page MJ, et al. *Cochrane handbook for systematic reviews of interventions*. Chichester, UK: John Wiley & Sons; 2019.
- Page MJ, Moher D, Bossuyt PM, Boutron I, Hoffmann TC, Mulrow CD, et al. PRISMA 2020 explanation and elaboration: updated guidance and exemplars for reporting systematic reviews. *BMJ* 2021;372:n160.
- Guyatt G, Oxman AD, Akl EA, Kunz R, Vist G, Brozek J, et al. GRADE guidelines: 1. Introduction—GRADE evidence profiles and summary of findings tables. *J Clin Epidemiol* 2011;64:383–94.
- Ma L-L, Wang Y-Y, Yang Z-H, Huang D, Weng H, Zeng X-T. Methodological quality (risk of bias) assessment tools for primary and secondary medical studies: what are they and which is better? *Mil Med Res* 2020;7(1):1–11.
- Otzen T, Manterola C, Mora M, Quiroz G, Salazar P, García N. Statements, recommendations, proposals, guidelines, checklists and scales available for reporting results in biomedical research and quality of conduct. A systematic review. *Int J Morphol* 2020;38(3):774–86.
- Network L. LATITUDES Network. 2023. Available at: <https://www.latitudes-network.org/>. Accessed October 10, 2023.
- Whiting PF, Rutjes AW, Westwood ME, Mallett S, Deeks JJ, Reitsma JB, et al. QUADAS-2: a revised tool for the quality assessment of diagnostic accuracy studies. *Ann Intern Med* 2011;155:529–36.
- Sterne JA, Savović J, Page MJ, Elbers RG, Blencowe NS, Boutron I, et al. RoB 2: a revised tool for assessing risk of bias in randomised trials. *BMJ* 2019;366:14898.
- Sterne JA, Hernán MA, Reeves BC, Savović J, Berkman ND, Viswanathan M, et al. ROBINS-I: a tool for assessing risk of bias in non-randomised studies of interventions. *BMJ* 2016;355:i4919.
- Wolff RF, Moons KG, Riley RD, Whiting PF, Westwood M, Collins GS, et al. PROBAST: a tool to assess the risk of bias and applicability of prediction model studies. *Ann Intern Med* 2019;170:51–8.
- Yang B, Mallett S, Takwoingi Y, Davenport CF, Hyde CJ, Whiting PF, et al. QUADAS-C: a tool for assessing risk of bias in comparative diagnostic accuracy studies. *Ann Intern Med* 2021;174:1592–9.
- Lee J, Mulder F, Leeflang M, Wolff R, Whiting P, Bossuyt PM. QUA-PAS: an adaptation of the QUADAS-2 tool to assess prognostic accuracy studies. *Ann Intern Med* 2022;175:1010–8.
- Higgins JMR, Rooney A, Taylor K, Thayer K, Silva R, et al, ROBINS-E Development Group. Risk of bias in non-randomized studies - of exposure (ROBINS-E). 2023. Available at: <https://www.riskofbias.info/welcome/robins-e-tool>. Accessed October 10, 2023.
- Mokkink LB, De Vet HC, Prinsen CA, Patrick DL, Alonso J, Bouter LM, et al. COSMIN risk of bias checklist for systematic reviews of patient-reported outcome measures. *Qual Life Res* 2018;27:1171–9.
- Collins GS, Reitsma JB, Altman DG, Moons KG. Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (TRIPOD) the TRIPOD statement. *Circulation* 2015;131:211–9.
- Schulz KF, Altman DG, Moher D. CONSORT 2010 statement: updated guidelines for reporting parallel group randomised trials. *J Pharmacol Pharmacother* 2010;1(2):100–7.
- Von Elm E, Altman DG, Egger M, Pocock SJ, Gøtzsche PC, Vandenbroucke JP. The Strengthening of Reporting of Observational Studies in Epidemiology (STROBE) statement: guidelines for reporting observational studies. *Lancet* 2007;370(9596):1453–7.
- Bossuyt PM, Reitsma JB, Bruns DE, Gatsonis CA, Glasziou PP, Irwig L, et al. Stard 2015: an updated list of essential items for reporting diagnostic accuracy studies. *Radiology* 2015;277:826–32.
- Moher D, Jones A, Lepage L, Group C, Group C. Use of the CONSORT statement and quality of reports of randomized trials: a comparative before-and-after evaluation. *JAMA* 2001;285:1992–5.
- Page MJ, McKenzie JE, Bossuyt PM, Boutron I, Hoffmann TC, Mulrow CD, et al. The PRISMA 2020 statement: an updated guideline for reporting systematic reviews. *Syst Rev* 2021;10(1):1–11.
- Tomlinson E. Risk of bias tools: a systematic review of usability. 2023. Available at: <https://osf.io/jf3xp/>. Accessed May 7, 2023.
- Cooper C, Booth A, Britten N, Garside R. A comparison of results of empirical studies of supplementary search techniques and recommendations in review methodology handbooks: a methodological review. *Syst Rev* 2017;6(1):1–16.
- Lumivero. NVivo (Version 14). 2023. Available at: https://techcenter.qsrinternational.com/Content/welcome/toc_welcome.htm. Accessed January 4, 2023.
- Kaizik MA, Garcia AN, Hancock MJ, Herbert RD. Measurement properties of quality assessment tools for studies of diagnostic accuracy. *Braz J Phys Ther* 2020;24(2):177–84.
- Venazzi A, Swardfager W, Lam B, Siqueira JO, Herrmann N, Cogomoraireira H. Validity of the QUADAS-2 in assessing risk of bias in Alzheimer's disease diagnostic accuracy studies. *Front Psychiatr* 2018;9:221.
- Jeyaraman MM, Robson RC, Copstein L, Al-Yousif N, Pollock M, Xia J, et al. Customized guidance/training improved the psychometric properties of methodologically rigorous risk of bias instruments for non-randomized studies. *J Clin Epidemiol* 2021;136:157–67.
- Zhang Y, Huang L, Wang D, Ren P, Hong Q, Kang D. The ROBINS-I and the NOS had similar reliability but differed in applicability: a random sampling observational studies of systematic reviews/meta-analysis. *J Evid Based Med* 2021;14(2):112–22.
- Kaiser I, Pfahlberg AB, Mathes S, Uter W, Diehl K, Steeb T, et al. Inter-rater agreement in assessing risk of bias in melanoma prediction studies using the prediction model risk of bias assessment tool (PROBAST): results from a controlled experiment on the effect of specific rater training. *J Clin Med* 2023;12(5):1976.
- Minozzi S, Cinquini M, Gianola S, Gonzalez-Lorenzo M, Banzi R. The revised Cochrane risk of bias tool for randomized trials (RoB 2) showed low interrater reliability and challenges in its application. *J Clin Epidemiol* 2020;126:37–44.
- Minozzi S, Dwan K, Borrelli F, Filippini G. Reliability of the revised Cochrane risk-of-bias tool for randomised trials (RoB2) improved with the use of implementation instruction. *J Clin Epidemiol* 2022;141:99–105.

- [31] Minozzi S, Cinquini M, Gianola S, Castellini G, Gerardi C, Banzi R. Risk of bias in nonrandomized studies of interventions showed low inter-rater reliability and challenges in its application. *J Clin Epidemiol* 2019;112:28–35.
- [32] Losilla JM, Oliveras I, Marin-Garcia JA, Vives J. Three risk of bias tools lead to opposite conclusions in observational research synthesis. *J Clin Epidemiol* 2018;101:61–72.
- [33] Venema E, Wessler BS, Paulus JK, Salah R, Raman G, Leung LY, et al. Large-scale validation of the prediction model risk of bias assessment Tool (PROBAST) using a short form: high risk of bias models show poorer discrimination. *J Clin Epidemiol* 2021;138:32–9.
- [34] Langenhuijsen LF, Janse RJ, Venema E, Kent DM, van Diepen M, Dekker FW, et al. Systematic meta-review of prediction studies demonstrates stable trends in bias and low PROBAST inter-rater agreement. *J Clin Epidemiol* 2023;159:159–73.
- [35] Wade R, Spackman E, Corbett M, Walker S, Light K, Naik R, et al. Adjunctive colposcopy technologies for examination of the uterine cervix - DySIS, LuViva Advanced Cervical Scan and Niris Imaging System: a systematic review and economic evaluation. *Health Technol Assess* 2013;17:1–240.
- [36] Crocker TF, Lam N, Jordão M, Brundle C, Prescott M, Forster A, et al. Risk-of-bias assessment using RoB2 was useful but challenging and resource-intensive: observations from a systematic review. *J Clin Epidemiol* 2023;161:39–45.
- [37] Wade R, Corbett M, Eastwood A. Quality assessment of comparative diagnostic accuracy studies: our experience using a modified version of the QUADAS-2 tool. *Res Synth Methods* 2013;4(3):280–6.