

AMSTAR-PF: a critical appraisal tool for systematic reviews of prognostic factor studies

Michael L Henry,¹ Neil E O'Connell,² Richard D Riley,^{3,4} Karel G M Moons,^{5,6} Beverley J Shea,^{7,8} Lotty Hooft,^{5,6} Johanna A A Damen,^{5,6} Nicole Skoetz,⁹ Sarah B Wallwork,¹ G Lorimer Moseley¹

For numbered affiliations see end of the article

Correspondence to: N E O'Connell
neil.oconnell@brunel.ac.uk
(or @NeilOConnell on X and @neiloconnell.bsky.social on Bluesky;
ORCID 0000-0003-1989-4537)

Additional material is published online only. To view please visit the journal online.

Cite this as: *BMJ* 2025;391:e085718
<http://dx.doi.org/10.1136/bmj-2025-085718>

Accepted: 13 October 2025

The ability to predict the onset or natural history of an illness, or how people may respond to a treatment, guides clinical decision making. These predictions are commonly based on prognostic factors: clinical, patient, or societal variables that are identified as being predictive of a certain future outcome. Prognostic factor research has increased across fields, with a subsequent increase in the number of systematic reviews of prognostic factors studies. Understanding the quality of such prognostic factor reviews is essential for confidence in their findings, but there is no quality appraisal instrument to specifically assess systematic reviews of prognostic factor studies. A Measurement Tool to Assess Systematic Reviews of Prognostic Factor studies, AMSTAR-PF, has been developed to fill this gap.

Prognosis research encompasses a range of methods and goals,^{1 2} and has been split broadly into four branches, detailed in the PROGRESS (prognosis research strategy) series of papers.³⁻⁶ Fundamental, or overall, prognosis research details the average course of health related conditions in a particular population and setting.³ Prognostic factor research aims to identify factors that are associated with the risk of certain outcomes occurring.⁴ Prognostic model research

concerns the development, validation, or impact of a prognostic model to estimate an individual's risk of a future outcome occurring.^{1 5} Stratified healthcare research uses prognostic information to tailor treatment to individuals.⁶

Knowledge of prognostic factors can form the basis for providing an individual with their likely prognosis, as well as being a building block for prognostic models, and a potential means of stratifying individuals into optimal treatment regimens.⁷ Prognostic factors can be drawn from a range of biopsychosocial domains, including clinical features (eg, tumour grade, blood, imaging or electrophysiology test results, pain levels, omics results), personal characteristics (eg, age, marital status, gender), psychological factors (eg, comorbid mental health problems, injury beliefs, motivation), and societal factors (eg, geographical location, healthcare availability, local environment). However, before implementation into healthcare, the supporting evidence must be reliable.

The number of prognostic factor studies published each year is increasing,^{8 9} many with uncertain or conflicting findings. Consequently, the number of systematic reviews seeking to synthesise, pool, and consolidate prognostic factor knowledge is also increasing.¹⁰ Hoffman and colleagues¹⁰ estimated a 20-fold increase in the total number of systematic reviews between 2000 and 2019, with 80 studies being published daily in 2019 across all research types. Between 2000 and 2004, 3.2% of all systematic reviews were classified as looking at causal links and risk; between 2015 and 2019, this proportion had increased to 22.8%.¹⁰ However, substantial doubts remain about the quality of published systematic reviews, not only of prognostic factors, but also across the full range of research types.¹¹⁻¹³ This uncertainty is important because these reviews often form the basis for policy and practice change, which makes the veracity of their findings immediately relevant to population level health outcomes.

Therefore, being able to determine the quality of the methods and content of published prognostic factor systematic reviews is critical. Prognostic factor research has a range of areas where quality can be compromised; some areas—such as the presence of a pre-registered protocol; clear research question; comprehensive search; non-selective reporting; and transparent and methodologically sound article selection, data extraction, and risk of bias assessment—are common to a range of research types, albeit sometimes with a different focus or methodology. Other important threats to validity that are not captured by other quality appraisal tools are specific to prognostic factor research, such as

SUMMARY POINTS

Research into prognostic factors is vital for many areas of healthcare
Confidence in the findings of systematic reviews of prognostic factor research can be compromised in a variety of ways
AMSTAR-PF (A Measurement Tool to Assess systematic Reviews of Prognostic Factor studies), based on AMSTAR 2, has been developed, refined, and tested
AMSTAR-PF uses signalling points and 19 questions over 14 domains to assist in coming to an overall judgment of confidence in the results of the review
Providing a standardised and reliable tool to appraise review quality will assist users to ascertain confidence in the review's findings

classification of prognostic factors and outcomes, adjustment factors, comparator and other prognostic factors, and appropriate calculation of prognostic effect sizes. This research directly informs policy and practice, so it is imperative that knowledge and methodology specific to prognostic factors are used to systematically appraise research, in both general areas of potential limitations of quality, as well as the domains specific to prognostic factors.¹⁴ Tools to appraise the quality of primary studies of prognostic factors exist, for example, the QUIPS (quality in prognosis studies) tool,¹⁵ but to our knowledge there is no dedicated tool to critically appraise the quality of systematic reviews of prognostic factor studies. We developed AMSTAR-PF (A Measurement Tool to Assess systematic Reviews of Prognostic Factor studies) to deal with this specific need.

AMSTAR-PF is based on AMSTAR 2 (A Measurement Tool to Assess systematic Reviews, version 2),¹⁶ which was developed to appraise systematic reviews of interventions. AMSTAR 2 evolved from AMSTAR¹⁷ by incorporating important feedback about the original AMSTAR, and adapting the tool to incorporate non-randomised, as well as randomised, interventional studies. AMSTAR-PF resembles AMSTAR 2 in several aspects of its content and guidance notes, but it has been comprehensively remodelled to deal with broad quality issues specific to prognostic factors, and with additions to allow systematic appraisal of issues that are unique to such research. Furthermore, alterations to the answering options allow more nuanced

recording of limitations and align it with a range of other appraisal tools.

Development of AMSTAR-PF

Initial stages and internal development

An overview of the entire development process is outlined in figure 1. In order to optimise efficiency, utility, and confidence in the new tool, we gathered a core research group that included experts who had developed tools and guidance for the following:

- AMSTAR¹⁷ and AMSTAR 2¹⁶ (www.amstar.ca) (BJS)
- ROBINS-I (risk of bias in non-randomised studies of interventions)¹⁸ and ROBINS-E (risk of bias in non-randomised studies of exposures; both www.riskofbias.info),¹⁹ (BJS)
- ROBIS (risk of bias in systematic reviews; www.bristol.ac.uk/population-health-sciences/projects/robis/robis-tool)²⁰ (BJS, KGMM)
- PROGRESS partnership³⁻⁶ (RDR, KGMM)
- PROBAST (prediction model risk of bias assessment tool) and PROBAST+AI (www.probast.org)²¹⁻²³ (KGMM, RDR, LH, JAAD)
- TRIPOD (transparent reporting of a multivariable prediction model for individual prognosis or diagnosis; www.tripod-statement.org) papers²⁴⁻²⁶ (KGMM, RDR, LH)
- CHARMS (checklist for critical appraisal and data extraction for systematic reviews of prediction modelling studies)²⁷ (KGMM)
- QUADAS-2 (updated quality assessment of diagnostic accuracy studies)²⁸ (KGMM)
- Guidance for using GRADE (grading of recommendations, assessment, development, and evaluations; www.gradeworkinggroup.org) for prognosis²⁹⁻³⁰ (RDR, BJS).

The core research group had an initial online meeting in July 2020 to discuss the need and optimal development approach for a quality appraisal tool specifically for systematic reviews of prognostic factors. We considered a range of tools on which to base the new tool. Deepest consideration focused on AMSTAR 2¹⁶ and ROBIS.²⁰ AMSTAR 2 is widely used to appraise systematic reviews of interventional studies; ROBIS was designed to appraise risk of bias in systematic reviews of a range of study types including prognosis, but was not considered specific or detailed enough to optimally manage prognostic factor reviews. We first created prognostic factor versions of both AMSTAR 2 and ROBIS. For AMSTAR 2, we added or substituted prognostic factor specific terms and items for the interventional items. For ROBIS, we focused on signalling only relevant to prognostic factors. We reviewed these two versions, considering their applicability and efficiency and comparing commonalities and differences within the domains. The overlap between these two modified tools was considerable, but they primarily differed in their focus—AMSTAR 2 being a quality appraisal tool and ROBIS being a risk of bias tool. Quality appraisal

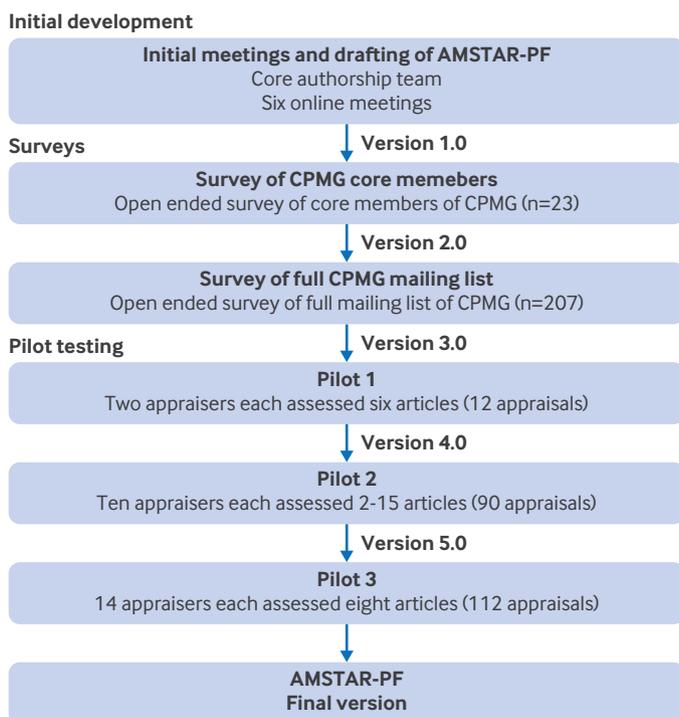


Fig 1 | Diagram of steps to developing AMSTAR-PF (a measurement tool to assess systematic reviews of prognostic factor studies). Feedback received was used at each step of the process to progressively redraft the AMSTAR-PF tool and guidance notes. CPMG=Cochrane Prognostic Methods Group

considers risk of bias, but also the external validity and the reporting quality of the reviewed article.³¹ Through further discussion, we concluded that AMSTAR 2 had a more user friendly format and provided a more complete skeleton for appraising overall quality, not just risk of bias.

The core research group also considered research that compared AMSTAR 2 and ROBIS in interventional studies. One study identified superior agreement among users and faster completion when using AMSTAR or AMSTAR 2 than when using ROBIS for interventional systematic reviews.³² Another study showed a preference among healthcare students for AMSTAR 2, reporting that AMSTAR 2 was clearer and easier to follow than ROBIS.³³ Similar preferences have been reported by other researchers.^{34 35} ROBIS was designed for a range of different review types, and its wider generalisability may come at the cost of utility, especially for content-naïve researchers. AMSTAR 2 has been widely used since its publication in 2017, and we predicted that many researchers would be familiar with it, which may aid ease of use and uptake of AMSTAR-PF. We reached 100% consensus within the core research group to use AMSTAR 2 as the basis for our new tool.

Having drafted an initial AMSTAR-PF, replacing AMSTAR 2's interventional focus with prognostic factor specific adaptations, we appraised this draft through extensive group discussion and considered further changes. We considered the relevance of each AMSTAR 2 domain in turn; additional domains or signalling points that needed to be created; and the structure and format of response options for both the individual domain questions and overall judgment.

Through a comprehensive iterative process that included email collaboration and six online meetings of the core research group between July 2020 and July 2021, drafts of AMSTAR-PF were created, reviewed, and revised. In July 2021, we had a finalised version 1.0 of the AMSTAR-PF tool, and then developed guidance notes. This version 1.0 of AMSTAR-PF tool and guidance notes was then sent for external review and feedback.

External review and feedback

The external feedback process and pilot testing were approved by the human research ethics committee of the University of South Australia (ID 203954, 206951), and participants gave informed consent to participate. We emailed the draft tool and guidance notes to the core members of the Cochrane Prognostic Methods Group (n=23), because this group was considered to have substantive expertise in prognostic research. The tool and guidance notes were accompanied by an anonymous survey (supplementary file 1A) that covered areas of redundancy, missing domains or questions, and the relevance and utility of each domain and of the tool as a whole. Responses were free text boxes. No questions were compulsory. Survey data were collected between 25 October and 27 November 2021 and managed using REDCap (Research Electronic

Data Capture) electronic data capture tools hosted at the University of South Australia.^{36 37} We received 17 responses; nine responses provided written feedback on the tool or guidance notes or both, and the remaining eight responses contained no information. The core research group then met online to discuss this feedback and made further revisions to the tool and guidance notes, thereby generating version 2.0 of the AMSTAR-PF tool and guidance notes.

AMSTAR-PF version 2.0 was then distributed, with the same anonymous REDCap survey described above, to the full mailing list of the Cochrane Prognostic Methods Group (n=207), who were permitted to share the email with colleagues they believed would have relevant expertise. Feedback was collected between 29 April and 27 May 2022. Twenty one responses were received; nine provided written feedback on AMSTAR-PF, the guidance notes or both. The authorship team reviewed the feedback, and as a result we developed a Microsoft Excel spreadsheet with an auto-populating summary page of question responses, as part of AMSTAR-PF version 3.0 and guidance notes.

Pilot testing

Three rounds of pilot testing occurred using 26 different appraisers. Twenty five of the appraisers were independent of the tool development; they volunteered following an email inviting participation that had been sent to a research group email list, or were asked because they were in the process of completing an umbrella review of prognostic factor studies (supplementary file 1B includes background information of the testers and supplementary file 1C includes details of the systematic reviews appraised).

Initial testing (pilot 1) of version 3.0 of AMSTAR-PF was performed by two PhD students who were independent of the tool development process. They independently appraised a convenience sample of six systematic reviews of prognostic factors for cardiovascular disease, and then provided appraisal results and qualitative written feedback on their experience of using the tool. The students had previously performed quality appraisals as part of systematic and umbrella reviews, including of prognostic factor research. This feedback led to very minor changes and generation of AMSTAR-PF version 4.0 and guidance notes.

A testing protocol for pilot 2 was pre-registered on the Open Science Framework (osf.io/3cgz9). Pilot 2 involved 10 appraisers (PhD students, postdoctoral researchers, and clinicians). One appraiser had been involved in developing AMSTAR-PF (MLH); nine appraisers had not. The appraisers used AMSTAR-PF version 4.0 to appraise 23 systematic reviews, which investigated prognostic factors for chronic pain (n=14), brain injury (n=6), post-traumatic stress disorder (n=2), and lymphoma (n=1), and included both narrative syntheses and meta-analyses. All the review articles were appraised independently by two or more appraisers, with 90 completed appraisals in total. All appraisers were invited to provide qualitative

Table 1 | Gwet's agreement coefficient for interrater, interpair, and intrapair agreement in pilot 3³⁸ using AMSTAR-PF version 5.0

	Interrater agreement	Interpair agreement	Intrapair
Domains (all answering options)	0.59 (0.21-0.90)*	0.61 (0.24-0.91)*	0.75 (0.45-0.95)†
Domains (collapsed)	0.72 (0.37-0.96)*	0.73 (0.36-0.96)*	0.83 (0.63-0.97)†
Overall appraisal	0.46 (0.30 to 0.62)‡	0.46 (0.17 to 0.74)‡	0.68 (0.22-1.00)*

Gwet's agreement coefficient was calculated unweighted for nominal data (questions 2b, 7c, 9a, 9b, 10, and 12, which include a "not applicable" option), and with linear weighting for ordinal data (remaining 13 questions, and the overall appraisal).

AMSTAR-PF=a measurement tool to assess systematic reviews of prognostic factor studies.

*Data are mean (range).

†Data are mean (range of means).

‡Data are the coefficient and 95% confidence interval.

feedback on the usability and acceptability of the tool, with eight providing feedback on their experiences. There were deviations from protocol (supplementary file 1D). Three reviewers were unable to complete their full quota of appraisals. Feedback received during pilot 2 suggested that further revisions would be beneficial, so a planned agreement analysis was not undertaken on pilot 2 in favour of making updates to the tool and guidance notes, then completing another pilot test and calculating agreement on the updated version. Reaching a final conclusion on the quality of each review was reported as being difficult, so it became a focus of revision. We had consciously avoided being too prescriptive in our instructions for deciding the overall quality of the review, but feedback from many appraisers suggested that a certain amount of guidance would be helpful as a starting point.

Additionally, two expert researchers, with substantial experience in prognostic research, reviewed AMSTAR-PF version 4.0 and guidance notes, and provided feedback. Minor revisions were then made, thereby generating AMSTAR-PF version 5.0.

Finally, we undertook a third pilot test (pilot 3) on AMSTAR-PF version 5.0. The protocol for pilot 3 was pre-registered on the Open Science Framework (osf.io/acrwf) and is described in detail elsewhere.³⁸ We deviated from the original protocol by our use of Gwet's agreement coefficient rather than kappa scores (supplementary file 1D). Pilot 3 involved eight systematic reviews, two each on prognostic factors for cancer, brain injury, lower back pain, and covid-19. The articles were appraised by 14 appraisers (11 female and three male, ranging in experience from undergraduate student to postdoctoral researcher) using AMSTAR-PF version 5.0 and guidance notes. No appraisers had been involved in the development of AMSTAR-PF, nor in any earlier feedback or testing rounds. On completion of pilot 3, these appraisers were also given the opportunity to provide feedback on any aspect of the tool and guidance notes. That feedback, along with analysis of each domain's agreement scores,³⁸ was used to make final minor changes to the guidance notes, signalling points and wording of the questions, and thus develop the final AMSTAR-PF tool and guidance notes.

Agreement and usage time

Interrater, interpair, and intrapair agreement data from pilot 3 are summarised in table 1 and provided

in detail elsewhere, alongside individual domain data and information on agreement over time.³⁸ Agreement was calculated twice for the domains, using Gwet's agreement coefficient³⁹: once with full answering options (yes; probably yes; probably no; no; and, for some questions, not applicable), and once with "yes" and "probably yes" collapsed into a single positive category and "no" and "probably no" collapsed into a single negative category. Agreement was moderate, substantial, or almost perfect for most domains, according to Gwet's application⁴⁰ of Landis and Koch's benchmarks.⁴¹ Variability was noted between raters and between pairs.³⁸ Agreement was higher with collapsed answering options, and intrapair agreement was routinely higher than interrater or interpair agreement. Mean time to complete quality appraisal across all articles was 39 (standard deviation 15) minutes. Mean time to reach consensus was 11 (7) minutes.³⁸

AMSTAR-PF tool

AMSTAR-PF contains 19 questions grouped into 14 domains. These domains are listed below, with justification and explanations about their content. The domains are arranged similarly to AMSTAR 2, and broadly follow the regular layout of a systematic review, from methods to results and discussion. The complete tool and guidance notes are included in supplementary materials 2A-C, and we encourage readers to refer to the guidance notes for more comprehensive detail on the included questions.

AMSTAR-PF domains

1. Research question

Systematic reviews of prognostic factors should have a clear, predefined research question with appropriate depth of information. PICOTS^{7 27} provides a framework for important information in prognostic factor research, and represents the population, index prognostic factors, comparator or other prognostic factors, outcome, timing, and setting. Having this information clearly detailed allows for a logical step to inclusion and exclusion factors, and for readers to judge the relevance of the review for their own practice or research.

The elements of PICOTS should be described in a way that makes the review question clear and reproducible, whether the authors choose to strictly define the parameters of a certain element or leave it more open and inclusive.

2. Pre-registered protocol

It is recommended that a systematic review is clearly planned, and that this process is documented in a protocol before commencing the review. Documented plans help as a safeguard to ensure that important methodological or analytical decisions are not based on findings made during the conduct of the review, and hence are an important step to minimise bias. The protocol should be registered and made publicly available by the time the systematic review is published, to enable readers to compare the protocol to the review conducted. Any deviations from the protocol should be documented and justified in the systematic review.

The AMSTAR-PF tool divides this domain into two questions, to clearly deal with both the presence of a publicly available pre-registered protocol (question 2a), and how any deviations from the protocol are handled (question 2b).

3. Included study designs and types

Authors should state which study designs are eligible for inclusion in the review. Prospective and retrospective cohort studies, case control studies, registry data, and data from randomised trials may be included, because all can give information on prognostic factors, although with potential different methodological issues to consider. Prospective cohort studies are seen as the ideal design for prognostic factor research, because researchers have maximum control over timing and types of prognostic factor and outcome measurements. Other study designs are more likely to have weaknesses, although are often easier and cheaper to run, and so are likely to be encountered in the reported research.

As well as different study designs, studies with different primary objectives may be included. Most commonly, studies of prognostic factors will be the source of evidence in a systematic review of prognostic factors, although in certain situations, prognostic modelling studies (ie, that aim to develop a model for individual risk prediction) or causal research can provide information about prognostic factors.¹

Decisions made by review authors on the types and designs of studies eligible for inclusion in a review can have repercussions on bias and heterogeneity and may necessitate separate synthesis when calculating prognostic factor effect sizes.

4. Search strategy

A comprehensive search strategy is a vital component of a systematic review of prognostic factors.^{42 43} The search strategy should be wide-ranging enough to ensure all relevant studies are captured and allow for a complete synthesis of available evidence. A systematic review founded on an inadequate search may be missing important data and hence give an unreliable picture of the current state of evidence on a topic.

What constitutes a comprehensive search strategy will be different in different topic areas, as well as for different prognostic factors and outcomes. Appraisers

are encouraged to consider the signalling points listed as well as any other potential sources of data relevant to the systematic review topic and aims.

5. Process for evaluating and including studies

At least two review authors should independently determine the eligibility of each study found in the search for inclusion into the systematic review. This process commonly involves a two stage process, whereby an initial screen is performed on title and abstract, followed by a secondary screening of the full text of all potentially relevant articles. While the process may differ in different reviews, at least two independent reviewers should be involved at each point of the decision making process. A plan for resolution of disagreements should be described.

6. Excluded studies

All potentially relevant studies that were read at the full text stage should be listed with a justification for exclusion of each. This step is important to allow readers to judge the validity of the systematic review and applicability to their own research or practice.

7. Data extraction

At least two assessors should independently extract the necessary data from the included studies, and a plan should be followed to resolve any disagreements. Having two independent extractors helps to minimise the risk of errors in this process.

The data extracted from the included studies is likely to vary depending on the aims of the prognostic factor review.²⁷ However, all included studies should be described in adequate detail to allow readers to assess the appropriateness of each study's inclusion and relevance.

Prognostic factor effect estimates and their precision from primary studies may be derived in different ways and presented in different formats, which can make extraction for the review challenging. In certain situations, authors of systematic reviews may have needed to calculate prognostic factor effect estimates and their precision from the results provided in the primary studies.⁷ Where these calculations have been done, the methods used should be clearly reported and appropriate. A major topic is whether unadjusted or adjusted prognostic factor results were extracted and, if adjusted, the adjustment factors of interest. Adjusted results examine the contribution of a prognostic factor over and above (ie, after adjusting for) other variables, which should typically refer to well known existing prognostic factors. However, extracting or deriving estimates adjusted for the same set of factors in each study is very challenging, and often impossible. How this issue was dealt with should be explained.

For clarity, AMSTAR-PF subdivides the data extraction domain into three questions, dealing with the process of data extraction (question 7a), the detail of report of extracted data (question 7b), and the appropriateness of calculating effect sizes from reported data, where relevant (question 7c).

8. Assessing risk of bias

Two assessors should independently assess the risk of bias of each included article in the systematic review and have a plan in place for resolving any disagreements during the consensus process.

The method chosen for assessing risk of bias should be pre-planned and clear, which will often involve the use of a recognised risk of bias tool; preferably one designed for prognostic factor studies such as QUIPS.¹⁵

This domain is subdivided into two questions in AMSTAR-PF, to cover both the process (question 8a) and technique used (question 8b) to complete assessment for risk of bias.

9. Data synthesis

The review protocol should have stated the principles on which review authors based their decision to synthesise data from included studies, and how it is planned to do this. There are many areas in the included studies where differences may arise (eg, lengths of follow-up; cut-off points used to dichotomise continuous factors; adjustment factors), and PICOTS^{7,27} provides a good framework for exploring many of these factors. Additionally, different ways of calculating and presenting prognostic effect sizes and their precision provides another area in which differences can arise, and where compatibility of studies decrease. Even in non-quantitative syntheses, the principles behind ensuring interpretability of results remain, and separate summaries may need to be presented (eg, by unadjusted and adjusted results; for each cut-off point) when included studies differ in key areas.

If meta-analysis was performed, authors are expected to have used appropriate methods, such as a random effects model to account for unexplained heterogeneity in prognostic factor effects. Authors should have also pre-planned how they will investigate heterogeneity, and quantify it appropriately. As well as statistical heterogeneity, an examination of potential clinical or methodological reasons for heterogeneity may be necessary.

This domain is subdivided into two questions in the AMSTAR-PF tool, focusing on the approach taken to ensure interpretability of results (question 9a), and on the analysis methods (in the case meta-analysis was performed; question 9b).

10. Small study effects

Small study effects is the phenomenon that occurs when studies with smaller sample sizes demonstrate systematically different prognostic factor effects than studies with larger sample sizes. This effect might be due to publication bias, such that smaller prognostic factor studies are likely to be published when they provide significant results, compared to non-significant results. However, heterogeneity may also cause small study effects, for example, if smaller studies have a shorter length of follow-up than larger studies, and a prognostic factor effect genuinely varies from short to longer term. Hence, prognostic factor reviews investigating asymmetry in funnel plots

should refer to the examination of small study effects, rather than publication bias, because it is not possible to be certain that publication bias is the reason for any asymmetry observed, particularly in the absence of formal registries for prognostic factor studies.

Commonly, statistical tests for asymmetry may accompany a funnel plot, although these tests often have low power, so they are not always a solution. If funnel plots or asymmetry tests do suggest that small study effects may be present, it is expected that this observation forms part of the discussion about the review's findings and conclusions.

11. Discussing risk of bias

Consideration of the potential effects of bias on the review's results and conclusions are important. In cases where only studies judged to have a low risk of bias were included in the review, there may be relatively little discussion.

If studies of varying quality were included (which is often the situation with prognostic factor reviews), then authors are expected to explore this potential source of bias. If a meta-analysis was performed, then evaluation of how effects vary by study quality may be undertaken with techniques such as subgroup analysis, meta-regression analysis, or sensitivity analysis. Even if these evaluations are not possible, we recommend that authors provide some commentary as to the possible impact of bias on individual included studies, and the results of the review.

12. Discussing heterogeneity

Heterogeneity in prognostic factor research has many potential causes, of which some are explored as sources of bias in question 8. Accounting for heterogeneity is also briefly covered in question 9. If a meta-analysis was performed, it is generally expected that heterogeneity will be accounted for in the meta-analysis and quantified, for example, using τ^2 (the estimate of variance in prognostic factor effects between studies) and prediction intervals (for a prognostic factor effect in a new study). In situations where only a small number of studies are included, then the estimate of heterogeneity will be very uncertain (and prediction intervals wide), but acknowledgment of this issue should be expected.

Irrespective of whether meta-analysis is undertaken and estimates of heterogeneity are obtained or not, review authors should still consider and discuss potential sources of clinical and methodological heterogeneity, and the possible impact on the results, conclusions, and recommendations of the systematic review.

13. Conflicts of interest

Review authors should document any funding sources or other potential conflicts of interest explicitly in the manuscript. Furthermore, any conflicts of interest in the included studies should also be documented, to help readers better assess the quality of evidence and any potential conflicts in the included studies.

14. Certainty in results

Review authors should deal with the level of certainty around their key findings. The GRADE guidelines⁴⁴ are a commonly used framework for authors, although review authors may use other methods to deal with this issue. At the time of writing, there is not a full range of GRADE guidelines covering all the range of prognostic factor studies, so authors may be required to modify existing resources (eg, existing GRADE guidelines for prognosis^{30 45}) or develop their own methods using similar principles.

Overall confidence in the results of the review

AMSTAR-PF is not designed to give an overall score, because a high score could mask critical failings in one or more key areas. Instead, each domain should be considered potentially vital to the overall confidence in the results of the review, and any errors or oversights noted appraised with this in mind.

We have suggested four categories for overall confidence in the results of the review: high, moderate, low, and critically low. The attached tool and guidance notes provide further detail about our suggested approach to classifying a review; however, we stress that appraisers can, and should, make decisions given their own topic and methodological knowledge.

Application of AMSTAR-PF

Many elements of AMSTAR-PF are open to varying interpretations and may have different levels of importance in different fields or in different reviews. We recommend that team members planning to use the tool meet before undertaking appraisal to ensure consistency in interpretation and application, and to discuss areas of topic or methodological importance that may need clarification or standardisation. We envisage that a common reason for appraising systematic review quality may be as part of an umbrella review. If so, it may be important to clarify areas of quality that are particularly important to the umbrella review question and outcomes sought, particularly if the findings will directly influence policy or practice.

AMSTAR-PF has 19 questions over 14 domains, with some of the questions containing signalling points to assist with answering the question. The response options to each question are yes (Y); probably yes (PY); probably no (PN); no (N); and for some questions, not applicable (N/A). The signalling points have the same response options. For simplicity, all questions and signalling points are worded so that “yes” indicates higher quality, and “no” indicates lower quality. We recommend that authors use these two options when the review has clear evidence for or against the signalling point or question, and that they use the “probably yes” and “probably no” options when the evidence is less clear, or assumptions need to be made when answering. The recommendations for using this answering system are consistent with other quality appraisal and risk of bias tools (eg, ROBINS-I,¹⁸ PROBAST,^{21 22} ROB 2 (version 2 of the Cochrane risk of bias tool for randomised trials),⁴⁶ and ROB-ME

(risk of bias due to missing evidence)⁴⁷). Similar to other appraisal tools,^{18 19 23 46 47} responses of “yes” and “probably yes” have similar implications for the overall appraisal of quality of the review, and as do “no” and “probably no.”

Deciding which answering option to choose requires a certain amount of judgment, alongside topical knowledge and methodological knowledge. The “probably” option is appropriate where the information available is imperfect and indicates that the reviewers have had to make a judgment because of that. For example, in question 1 (“Did the review clearly define the research question, including the relevant components of PICOTS?”), the first step would be to assess the extent to which the signalling questions have been answered in the systematic review, and mark these accordingly. Using those responses as a guide, appraisers should then consider the question as a whole and to what extent it was answered. If appraisers think that the review has clearly defined all relevant aspects of the review question, then “yes” is an appropriate response. If they think that the question is adequately defined, but missing an element of PICOTS or specificity around a certain signalling point or aspect of the review question, then “probably yes” may be preferred, because it indicates that there is a lack of clarity and a judgment was made. Conversely, we recommend that a review with a poorly defined, or ambiguous research question is given a “no” response, whereas a review with a few well defined elements but still overall lacking clarity and appropriate definitions of the elements may receive a “probably no” response. These response options may vary in different reviews, in different questions, and across different reviewer teams and fields.

The AMSTAR-PF guidance notes (supplementary file 2B) provide more detailed information about each of AMSTAR-PF’s domains and application, and we recommend that users of AMSTAR-PF review these notes and refer to them when using the tool. We stress that given the variety of topic areas, methodologies, review aims, and expertise of appraisers on a team, these notes should be seen as guidance only; teams may wish, or need, to further operationalise their application. Furthermore, teams may find it useful to pre-plan a meeting after appraisers have completed appraisal of 3-6 reviews, in order to confirm their interpretations, compare appraisals, and identify any unexpected inconsistencies in interpretation or application.^{38 48} This approach is especially true if appraisers have different experience levels or knowledge bases, or need more guidance in certain areas.

Discussion

AMSTAR-PF was specifically designed to appraise the quality of systematic reviews of prognostic factors. It is, to our knowledge, the first such tool and fills an important gap for people assessing or undertaking systematic reviews of prognostic factors, which have many unique challenges to consider.⁷ The development

Table 2 | Summary of key differences between AMSTAR 2 and AMSTAR-PF with selected examples

Aspect of tool	Key changes in AMSTAR-PF	Critical examples
Content	Added questions or items for prognostic factor studies and review quality	Question 7c: added question on obtaining prognostic factor effect estimates
		Question 8a: added question about the process of assessing risk of bias
		Question 9a: added question detailing the approach to synthesis
		Question 14: added question on the certainty of findings
	Adapted elements of existing questions to be more relevant to prognostic factors	Question 1: substituted the PICOTS acronym for PICO
		Question 3: stipulated types of prognostic factor studies for inclusion
	Revised signalling points: additions, removals, and modifications	
Question 4: added requirement for the full search strategy to be presented, and recommended 12 months (not 24 months) for the search time frame		
Questions 5, 7a, 8a: added signalling about the plan for resolution of disagreements		
Question 8b: modified the risk of bias assessment to be based on QUIPS		
Question 9a: added specific synthesis guidance specific to prognostic factors		
Structure and style	Changed question answering options	All questions: AMSTAR-PF has yes; probably yes; probably no; no; and for some questions, not applicable; AMSTAR 2 had response options of yes, no, partial yes (for some questions), and equivalent of not applicable
	Changed signalling point answering options	All signalling points: AMSTAR-PF has yes, probably yes, probably no, no, and for some questions, not applicable; AMSTAR 2 had checkboxes only
	Reordered certain questions to assist with a logical flow	Moved excluded studies after the screening, as opposed to after data extraction (questions 5-6 in AMSTAR-PF; questions 5 and 7 in AMSTAR 2) Grouped together questions about meta-analysis and data synthesis (questions 9a, 9b, 10 in AMSTAR-PF; questions 11 and 15 in AMSTAR 2)
	Developed an auto-populating spreadsheet version of the tool, alongside traditional document forms and guidance notes	
Overall rating	Modified method for arriving at an overall judgment of quality for the review; removed the concept of critical domains	
Guidance notes	Developed de novo detailed guidance notes for all aspects of the tool	

AMSTAR 2=a measurement tool to assess systematic reviews, version 2¹⁶; AMSTAR-PF=a measurement tool to assess systematic reviews of prognostic factor studies; PICOTS=population, index prognostic factor, comparator prognostic factors, outcome, timing, setting⁷; PICO=patient, population or problem, intervention, comparison, outcome; QUIPS=quality in prognosis studies.¹⁵

process involved a wide range of researchers, was comprehensive and iterative, and several pilot trials were undertaken before arriving at the final version presented here.

AMSTAR-PF and guidance notes were based on AMSTAR 2 and therefore share commonalities. There are, however, differences in the domains and how they have been presented, with some of AMSTAR-PF's domains divided into questions in order to highlight important elements of that domain. Other modifications were a change in the order of some domains, and changes in guidance to reflect current ideas and recommendations for best practice that have evolved since AMSTAR 2 was published in 2017. For instance, we recommend that it is not enough to state there was a protocol, but rather that it should be publicly available (question 2a). Furthermore, we added a question about the process used in performing risk of bias (question 8a), one about the interpretability of results (question 9a), and one dealing with the certainty around key findings (question 14). Other changes were necessitated to accommodate the different focuses of prognostic factor research; changes were made to the wording of questions, and additional questions (eg, question 7c, around prognostic factor effect estimates) were added.

Categories for the overall rating of the appraised systematic review remain the same in AMSTAR-PF as they are in AMSTAR 2, but answering options for the domains have been changed. AMSTAR-PF includes “probably yes” and “probably no” responses for all questions, whereas a “partial yes” option was only available in some AMSTAR 2 questions, with no “probably no” option in any AMSTAR 2

question. AMSTAR-PF also uses these same options (yes, probably yes, probably no, no, and/or not applicable) for each of the signalling points, which diverges from the checkboxes used in AMSTAR 2. We considered a range of options for answering the questions and signalling points. Some other tools have a “no information” or “unclear” option, or force a binary “yes” or “no” decision for questions and/or signalling points. The change to include the additional “probably” options aims to enhance the usability of the tool, especially considering that judgments will need to be made where there is inherent uncertainty, and the extra options provide added detail that may be useful when two appraisers meet to reach consensus. We acknowledge that such responses can add complexity when coming to a final decision on each domain and on the appraisal as a whole, but we consider that the overall benefit of the added detail makes this change worthwhile, especially given recognised deficits in quality in reporting of prognosis research (www.tripod-statement.org).^{7 8 25 49} This answering structure aligns AMSTAR-PF with many other tools currently being used (eg, ROBINS-I,¹⁸ ROBINS-E),¹⁹ ROB 2,⁴⁶ ROBIS,²⁰ ROB-ME⁴⁷ and PROBAST²¹⁻²³).

AMSTAR 2 uses the concept of critical domains—prespecified domains considered to be of integral importance to a review's quality—to assist in coming to a final judgment on the quality of the review under appraisal. We consider that such critical domains may be less consistent in prognostic factor research than in interventional research, and that it is more beneficial to instead have appraisers approach each and every domain as potentially critical to the outcome and findings of the review. We therefore chose to

diverge from AMSTAR 2 in this respect. Broadly, however, AMSTAR-PF follows a similar format and style to AMSTAR 2. These deliberate similarities may assist users already familiar with AMSTAR 2. Table 2 summarises the main differences between AMSTAR 2 and AMSTAR-PF.

AMSTAR-PF is designed to facilitate a comprehensive appraisal of the quality of systematic review of prognostic factors. Whiting and colleagues³¹ describe three components of quality: internal validity (risk of bias), external validity (applicability/variability), and reporting quality. AMSTAR-PF incorporates considerations of each of these components in the included domains, and all should be considered when coming to a final judgment of the quality of the review in question. Like AMSTAR 2, the highlighting of different elements of quality may make AMSTAR-PF useful as a brief checklist for people conducting systematic reviews of prognostic factors, or as a teaching aid. It is not, however, a substitute for more comprehensive methodology and rationale in conducting systematic reviews, nor a replacement for detailed risk of bias guidance.

AMSTAR-PF was developed for systematic reviews of prognostic factor studies. Prognostic factor research, however, is variably defined; some authors,⁴⁹ for instance, highlight causal prognostic factors as part of this broad subgroup, and differentiate these as “prognostic determinants (causal)” as opposed to “prognostic markers (non-causal).” AMSTAR-PF does not differentiate between causal and non-causal prognostic factors, and it is suitable for either or both. Specific guidance on appraising the quality of the methods used to establish the presence and/or degree of causation is, however, outside the scope of this tool, as are other causal questions. Areas of research outside of prognostic factor reviews may need their own specific guidance, whether from a distinct tool or tailored by modifying existing resources, such as AMSTAR-PF.

The development of AMSTAR-PF and resulting tool had several strengths. We developed AMSTAR-PF using a rigorous, multistage design and validation process, which aligned with recommendations for developing quality assessment tools,³¹ and had input from diverse international experts within and external to the core research group. Responses from the surveys to members of the Cochrane Prognostic Methods Group confirmed a need for a quality appraisal tool for systematic reviews of prognostic factor studies, and feedback was generally positive for the included items and AMSTAR-PF as a whole. With any other feedback, changes were made to improve the tool. Three pilot trials to test agreement and usability, involving 25 independent testers with diverse backgrounds and research expertise, contributed to the final version.

Consideration of the results of our testing of AMSTAR-PF should recognise that the participants in our pilot trials may not be representative of all researchers. In particular, most were not experienced in prognostic factor research, nor were they always

knowledgeable about the topic area of the sample reviews. We consider that this point holds ecological validity—critical appraisal of systematic reviews may often be undertaken by graduate students who are new to both prognostic research and the health condition of interest. The agreement scores obtained in our testing showed slight to substantial interrater and interpair agreement, and fair to almost perfect intrapair agreement, which is similar or better than other published tools for risk of bias or quality appraisal.³⁸ The routinely higher intrapair agreement scores than interrater and interpair agreement scores may point to a benefit of discussing interpretation and understanding of the tool before using it, or after a small number of initial appraisals, in order to better apply the AMSTAR-PF questions within appraisers’ specific areas. Additionally, the guidance notes were modified to provide more thorough guidance following results and feedback from the testing, which may have enhanced the agreement and usability of the final version presented here.

We faced several challenges in developing this tool. We attempted to canvas a wide range of opinions and feedback on AMSTAR-PF, although participation rates of those individuals invited to review it were sometimes low. This proportion is highlighted in the survey to the Cochrane Prognostic Methods Group, where the response rate was around 10% over the four week response period (assuming that all emails were still monitored and were not duplicates). This participation rate is low, although it is within published ranges for emailed surveys.^{50,51} The low response may reflect the time commitment needed to comprehensively review and provide feedback on the document, as well as difficulties with recruiting busy academics, and a lack of incentive for completing it.

Many of the concepts in appraising prognostic factor reviews are complex, and determining whether something has been dealt with adequately is inherently subjective. This challenge is recognised with all such tools, and the final version of AMSTAR-PF represents answering options and decision making guidance that we deemed to balance the potential benefits of prescriptive, standardised scoring with the need to allow appraisers’ judgments and expertise to influence the appraisal process. We have sought to provide additional guidance in coming to judgments, but acknowledge that our guidance is general and may not be detailed enough for all situations; all review teams should include members with appropriate methodological and topical knowledge. In certain situations, appraisers may feel that it is beneficial to modify or amend certain elements to better suit their purposes. This approach comes with risks, however, so we recommend that appraisers carefully consider the potential drawbacks of making alterations, particularly because non-standard tools present difficulties in interpretation and comparability. When modifications are deemed necessary, clear documentation of what changes were made and why they were made will be critical.

Conclusion

To our knowledge, AMSTAR-PF is the first quality appraisal tool specifically developed for systematic reviews of prognostic factors. We undertook a rigorous iterative process with field experts to develop both the tool and the guidance notes. We obtained input from a wide group of experienced prognosis researchers, and tested the tool for agreement, usability, and acceptability in three separate cohorts of largely inexperienced researchers. We present the final tool to the field for immediate application.

AUTHOR AFFILIATIONS

¹IIMPACT in Health, University of South Australia, Karna Country, Adelaide, SA, Australia

²Centre for Health and Wellbeing Across the Lifecourse, Department of Health Sciences, Brunel University London, Uxbridge UB8 3PH, UK

³Department of Applied Health Sciences, School of Health Sciences, College of Medicine and Health, University of Birmingham, Birmingham, UK

⁴National Institute for Health and Care Research (NIHR) Birmingham Biomedical Research Centre, Birmingham, UK

⁵Cochrane Netherlands, University Medical Centre Utrecht, Utrecht University, Utrecht, Netherlands

⁶Julius Centre for Health Sciences and Primary Care, University Medical Centre Utrecht, Utrecht University, Utrecht, Netherlands

⁷School of Epidemiology and Public Health, Faculty of Medicine, University of Ottawa, Ottawa, ON, Canada

⁸Bruyère Health Research Institute, University of Ottawa, Ottawa, ON, Canada

⁹Institute of Public Health, Faculty of Medicine and University Hospital Cologne, University of Cologne, Cologne, Germany

We thank all those individuals who provided feedback and suggestions on earlier iterations of this work, as well as those who piloted the tool, including: Ruth Appiah, Carolyn Berryman, Chloe Blacket, Aidan Cashin, Sophie Crouch, Grace Ferencz, Michael Ferraro, Bethany Gower, Ashley Grant, Katherine Henry, Aleksandra Herman, Emma Karran, Indika Koralegedera, Hayley Leake, Erin MacIntyre, Brendan Mouatt, Karma Phuentsho, Chandan Poddar, Daniel Van Der Laan, Ellana Welsby, Louise Wiles, Erica Wilkinson, Marelle Wilson, Monique Wilson, and members of the Cochrane Prognostic Methods Group. We also thank the writers of AMSTAR and AMSTAR 2, whose comprehensive and user friendly tool and guidance notes formed a basis for the current work.

Contributors: MLH, NEO'C, and GLM conceived the original idea. MLH, NEO'C, RDR, KGMM, BJS, and LH developed the initial versions of the tool and guidance notes. All authors were involved in subsequent revisions. MLH wrote the first draft of the manuscript. All authors critically revised the manuscript and approved the final version. MLH is the guarantor of the article. The corresponding author attests that all listed authors meet authorship criteria and that no others meeting the criteria have been omitted.

Funding: This work received no specific funding. MLH was supported by an Australian Government Research Training Scholarship. RDR was supported by funding from the MRC Better Methods Better Research panel (grant MR/VO38168/1) and by the National Institute for Health and Care Research (NIHR) Birmingham Biomedical Research Centre at the University Hospitals Birmingham NHS Foundation Trust and the University of Birmingham. GLM was supported by a leadership investigator grant from the National Health and Medical Research Council of Australia (ID 1178444). The funders had no role in considering the study design or in the collection, analysis, interpretation of data, writing of the report, or decision to submit the article for publication.

Competing interests: All authors have completed the ICMJE uniform disclosure form at <https://www.icmje.org/disclosure-of-interest/> and declare: no support from any organisation for the submitted work. Between 2020 and 2023, NEO'C was coordinating editor of the Cochrane Pain, Palliative and Supportive Care group, whose activities were funded by the UK NIHR; is chair of the International Association for the Study of Pain (IASP) Methodology, Evidence Synthesis and Implementation Special Interest Group; and has held a grant from the ERA-NET Neuron Co-Fund. RDR is an NIHR senior investigator;

and receives royalties from two textbooks ("Prognosis Research in Healthcare" and "IPD Meta-Analysis") and consultancy fees as a statistical editor for *The BMJ*. GLM has received support, unrelated to the current work, from Reality Health, ConnectHealth UK, Institutes of Health California, AIA Australia, Workers' Compensation Boards, and professional sporting organisations in Australia, Europe, South America, and North America; from professional and scientific bodies for travel costs related to presentation of research on pain and pain education at scientific conferences/symposiums; speaker fees for lectures on pain, pain education, and rehabilitation; and royalties for books on pain and pain education. GLM is non-paid CEO of the non-profit organisation Pain Revolution. There are no disclosures immediately relevant to this work. There are no other relationships or activities that could appear to have influenced the submitted work.

Patient and public involvement: We did not include patients or the public in this methodologically focused project.

Transparency: The lead author (MLH) affirms that the manuscript is an honest, accurate, and transparent account of the study being reported; that no important aspects of the study have been omitted; and that any discrepancies from the study as planned (and registered) have been explained.

Dissemination to participants and related patient and public communities: There are no specific patient and public dissemination plans for this paper. The work is methodologically focused and primarily intended for the research and clinical community, and did not involve patients or public during its development.

Provenance and peer review: Not commissioned; externally peer reviewed.

Data sharing: Data related to agreement testing can be found in the referenced paper, or on reasonable request from the corresponding author at neiloconnell@brunel.ac.uk. The draft versions of the tool and guidance notes generated during the development process are available on reasonable request, from the corresponding author.

- 1 Moons KG, Royston P, Vergouwe Y, Grobbee DE, Altman DG. Prognosis and prognostic research: what, why, and how? *BMJ* 2009;338:b375.
- 2 Moons KG, Hooft L, Williams K, Hayden JA, Damen JA, Riley RD. Implementing systematic reviews of prognosis studies in Cochrane. *Cochrane Database Syst Rev* 2018;10:ED000129.
- 3 Hemingway H, Croft P, Perel P, et al, PROGRESS Group. Prognosis research strategy (PROGRESS) 1: a framework for researching clinical outcomes. *BMJ* 2013;346:e5595.
- 4 Riley RD, Hayden JA, Steyerberg EW, et al, PROGRESS Group. Prognosis Research Strategy (PROGRESS) 2: prognostic factor research. *PLoS Med* 2013;10:e1001380.
- 5 Steyerberg EW, Moons KG, van der Windt DA, et al, PROGRESS Group. Prognosis Research Strategy (PROGRESS) 3: prognostic model research. *PLoS Med* 2013;10:e1001381.
- 6 Hingorani AD, Windt DA, Riley RD, et al, PROGRESS Group. Prognosis research strategy (PROGRESS) 4: stratified medicine research. *BMJ* 2013;346:e5793.
- 7 Riley RD, Moons KGM, Snell KIE, et al. A guide to systematic review and meta-analysis of prognostic factor studies. *BMJ* 2019;364:k4597.
- 8 Riley RD, van der Windt DA, Croft P, Moons KGM. *Prognosis research in healthcare: concepts, methods, and impact*. Oxford University Press, 2019.
- 9 Debray TPA, de Jong VMT, Moons KGM, Riley RD. Evidence synthesis in prognosis research. *Diagn Progn Res* 2019;3:13.
- 10 Hoffmann F, Allers K, Rombey T, et al. Nearly 80 systematic reviews were published each day: Observational study on trends in epidemiology and reporting over the years 2000-2019. *J Clin Epidemiol* 2021;138:1-11.
- 11 Ioannidis JPA. The mass production of redundant, misleading, and conflicted systematic reviews and meta-analyses. *Milbank Q* 2016;94:485-514.
- 12 Page MJ, Shamseer L, Altman DG, et al. Epidemiology and reporting characteristics of systematic reviews of biomedical research: a cross-sectional study. *PLoS Med* 2016;13:e1002028.
- 13 Almeida MO, Yamato TP, Parreira PDCS, Costa LOP, Kamper S, Saragiotto BT. Overall confidence in the results of systematic reviews on exercise therapy for chronic low back pain: a cross-sectional analysis using the Assessing the Methodological Quality of Systematic Reviews (AMSTAR) 2 tool. *Braz J Phys Ther* 2020;24:103-17.
- 14 Hayden JA, Côté P, Bombardier C. Evaluation of the quality of prognosis studies in systematic reviews. *Ann Intern Med* 2006;144:427-37.
- 15 Hayden JA, van der Windt DA, Cartwright JL, Côté P, Bombardier C. Assessing bias in studies of prognostic factors. *Ann Intern Med* 2013;158:280-6.

- 16 Shea BJ, Reeves BC, Wells G, et al. AMSTAR 2: a critical appraisal tool for systematic reviews that include randomised or non-randomised studies of healthcare interventions, or both. *BMJ* 2017;358:j4008.
- 17 Shea BJ, Grimshaw JM, Wells GA, et al. Development of AMSTAR: a measurement tool to assess the methodological quality of systematic reviews. *BMC Med Res Methodol* 2007;7:10.
- 18 Sterne JA, Hernán MA, Reeves BC, et al. ROBINS-I: a tool for assessing risk of bias in non-randomised studies of interventions. *BMJ* 2016;355:i4919.
- 19 Higgins JPT, Morgan RL, Rooney AA, et al. A tool to assess risk of bias in non-randomized follow-up studies of exposure effects (ROBINS-E). *Environ Int* 2024;186:108602.
- 20 Whiting P, Savović J, Higgins JP, et al, ROBIS group. ROBIS: A new tool to assess risk of bias in systematic reviews was developed. *J Clin Epidemiol* 2016;69:225-34.
- 21 Wolff RF, Moons KGM, Riley RD, et al, PROBAST Group. PROBAST: A Tool to Assess the Risk of Bias and Applicability of Prediction Model Studies. *Ann Intern Med* 2019;170:51-8.
- 22 Moons KGM, Wolff RF, Riley RD, et al. PROBAST: A Tool to Assess Risk of Bias and Applicability of Prediction Model Studies: Explanation and Elaboration. *Ann Intern Med* 2019;170:W1-33.
- 23 Moons KGM, Damen JAA, Kaul T, et al. PROBAST+AI: an updated quality, risk of bias, and applicability assessment tool for prediction models using regression or artificial intelligence methods. *BMJ* 2025;388:e082505.
- 24 Moons KG, Altman DG, Reitsma JB, et al. Transparent Reporting of a multivariable prediction model for Individual Prognosis or Diagnosis (TRIPOD): explanation and elaboration. *Ann Intern Med* 2015;162:W1-73.
- 25 Collins GS, Reitsma JB, Altman DG, Moons KGM. Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (TRIPOD): the TRIPOD Statement. *BMC Med* 2015;13:1-10.
- 26 Collins GS, Moons KGM, Dhiman P, et al. TRIPOD+AI statement: updated guidance for reporting clinical prediction models that use regression or machine learning methods. *BMJ* 2024;385:e078378.
- 27 Moons KG, de Groot JA, Bouwmeester W, et al. Critical appraisal and data extraction for systematic reviews of prediction modelling studies: the CHARMS checklist. *PLoS Med* 2014;11:e1001744.
- 28 Whiting PF, Rutjes AW, Westwood ME, et al, QUADAS-2 Group. QUADAS-2: a revised tool for the quality assessment of diagnostic accuracy studies. *Ann Intern Med* 2011;155:529-36.
- 29 Iorio A, Spencer FA, Falavigna M, et al. Use of GRADE for assessment of evidence about prognosis: rating confidence in estimates of event rates in broad categories of patients. *BMJ* 2015;350:h870.
- 30 Foroutan F, Guyatt G, Zuk V, et al. GRADE Guidelines 28: Use of GRADE for the assessment of evidence about prognostic factors: rating certainty in identification of groups of patients with different absolute risks. *J Clin Epidemiol* 2020;121:62-70.
- 31 Whiting P, Wolff R, Mallett S, Simeria I, Savović J. A proposed framework for developing quality assessment tools. *Syst Rev* 2017;6:204.
- 32 Gates M, Gates A, Duarte G, et al. Quality and risk of bias appraisals of systematic reviews are inconsistent across reviewers and centers. *J Clin Epidemiol* 2020;125:9-15.
- 33 Lee SWH. What tool do undergraduate pharmacy students prefer when grading systematic review evidence: AMSTAR-2 or ROBIS? *Cochrane Evid Synth Methods* 2023;1:e12023.
- 34 Pieper D, Puljak L, González-Lorenzo M, Minozzi S. Minor differences were found between AMSTAR 2 and ROBIS in the assessment of systematic reviews including both randomized and nonrandomized studies. *J Clin Epidemiol* 2019;108:26-33.
- 35 Perry R, Whitmarsh A, Leach V, Davies P. A comparison of two assessment tools used in overviews of systematic reviews: ROBIS versus AMSTAR-2. *Syst Rev* 2021;10:273.
- 36 Harris PA, Taylor R, Minor BL, et al, REDCap Consortium. The REDCap consortium: Building an international community of software platform partners. *J Biomed Inform* 2019;95:103208.
- 37 Harris PA, Taylor R, Thielke R, Payne J, Gonzalez N, Conde JG. Research electronic data capture (REDCap)—a metadata-driven methodology and workflow process for providing translational research informatics support. *J Biomed Inform* 2009;42:377-81.
- 38 Henry M, O'Connell N, Riley R, et al. Agreeability testing of AMSTAR-PF, a tool for quality appraisal of systematic reviews of prognostic factor studies. *medRxiv* [Preprint]. 2025:2025.04.10.25325555. doi: 10.1101/2025.04.10.25325555.
- 39 Gwet KL. Computing inter-rater reliability and its variance in the presence of high agreement. *Br J Math Stat Psychol* 2008;61:29-48.
- 40 Gwet KL. *Handbook of Inter-Rater Reliability: The Definitive Guide to Measuring the Extent of Agreement Among Raters*. 4th ed. Advanced Analytics, LLC, 2014.
- 41 Landis JR, Koch GG. The measurement of observer agreement for categorical data. *Biometrics* 1977;33:159-74.
- 42 Geersing GJ, Bouwmeester W, Zuihthoff P, Spijker R, Leeflang M, Moons KG. Search filters for finding prognostic and diagnostic prediction studies in Medline to enhance systematic reviews. *PLoS One* 2012;7:e32844.
- 43 Stallings E, Gaetano-Gil A, Alvarez-Diaz N, et al. Development and evaluation of a search filter to identify prognostic factor studies in Ovid MEDLINE. *BMC Med Res Methodol* 2022;22:107.
- 44 Schünemann H, Brozek J, Guyatt G, Oxman A, eds. GRADE handbook for grading quality of evidence and strength of recommendations. Updated October 2013. The GRADE Working Group, 2013. <https://gdt.grade.pro.org/app/handbook/handbook.html>
- 45 Huguet A, Hayden JA, Stinson J, et al. Judging the quality of evidence in reviews of prognostic factor research: adapting the GRADE framework. *Syst Rev* 2013;2:71.
- 46 Sterne JAC, Savović J, Page MJ, et al. RoB 2: a revised tool for assessing risk of bias in randomised trials. *BMJ* 2019;366:14898.
- 47 Page MJ, Sterne JAC, Boutron I, et al. ROB-ME: a tool for assessing risk of bias due to missing evidence in systematic reviews with meta-analysis. *BMJ* 2023;383:e076754.
- 48 Boutron I, Page MJ, Higgins JPT, Altman DG, Lundh A, Hróbjartsson A. Chapter 7: Considering bias and conflicts of interest among the included studies [last updated August 2022]. In: Higgins JPT, Thomas J, Chandler J, et al, editors. *Cochrane Handbook for Systematic Reviews of Interventions*. Version 6.5. Cochrane; 2024. www.training.cochrane.org/handbook
- 49 Kent P, Cancelliere C, Boyle E, Cassidy JD, Kongsted A. A conceptual framework for prognostic research. *BMC Med Res Methodol* 2020;20:172.
- 50 Shih T-H, Fan X. Comparing response rates in e-mail and paper surveys: A meta-analysis. *Educ Res Rev* 2009;4:26-40.
- 51 Wu M-J, Zhao K, Fils-Aime F. Response rates of online surveys in published research: A meta-analysis. *Comput Hum Behav Rep* 2022;7:100206.

Supplementary file 1A-1D: AMSTAR-PF survey; Appraiser demographics; List of references used during the piloting process; Deviations from pilot testing protocols

Supplementary file 2A-2C: AMSTAR-PF tool; AMSTAR-PF guidance notes; AMSTAR Excel form