

The LEADING guideline: Reporting standards for expert panel, best-estimate diagnosis, and longitudinal expert all data (LEAD) methods

Veerle C. Eijsbroek^{a,*}, Katarina Kjell^a, H. Andrew Schwartz^b, Jan R. Boehnke^c, Eiko I. Fried^d, Daniel N. Klein^e, Peik Gustafsson^f, Isabelle Augenstein^g, Patrick M.M. Bossuyt^h, Oscar N.E. Kjell^a

^a Department of Psychology, Lund University, Lund, Sweden

^b Department of Computer Science, Stony Brook University, New York, United States

^c School of Health Sciences, University of Dundee, Dundee, UK

^d Institute of Psychology, Leiden University, Leiden, the Netherlands

^e Department of Psychology, Stony Brook University, New York, United States

^f Faculty of Medicine, Lund University, Lund, Sweden

^g Department of Computer Science, University of Copenhagen, Copenhagen, Denmark

^h Department of Epidemiology and Data Science, Amsterdam University Medical Centers, Amsterdam, the Netherlands

ARTICLE INFO

Keywords:

Reporting
Assessment
Expert panel
Longitudinal Expert All Data
Best-estimate diagnosis
Standard
Psychiatry

ABSTRACT

Accurate assessments of symptoms and illnesses are essential for health research and clinical practice but face many challenges. The absence of a single error-free measure is currently addressed by assessment methods involving experts reviewing several sources of information to achieve a *best-estimate assessment*. This assessment method is called the *Expert Panel* method in medicine, and the *Best-Estimate Diagnosis or Longitudinal Expert All Data (LEAD)* method in psychiatry and psychology. However, due to poor reporting of the assessment methods, the quality of pro-claimed best-estimate assessments is typically difficult to evaluate, and when the method is reported, the reporting quality varies substantially. To tackle this gap, we have developed a reporting guideline following a four-stage approach: 1) drafting reporting standards accompanied by empirical evidence, which were further developed with a patient organization for depression, 2) incorporating expert feedback through a two-round Delphi procedure, 3) refining the guideline based on an expert consensus meeting, and 4) testing the guideline by i) having researchers test it and ii) applying it to previously published studies. The last step also provides evidence for the need for the guideline: 10–63 % (Mean 33 %) of the standards were not reported across thirty randomly selected published studies. The result is the LEADING guideline comprising 20 reporting standards in four groups: the *Longitudinal design*, the *Appropriate data*, the *Evaluation – experts, materials and procedures*, and the *Validity* group. We hope that the LEADING guideline will assist researchers in planning, conducting, reporting, and evaluating research aiming to achieve best-estimate assessments.

1. Introduction

Establishing valid and reliable assessments of symptoms and diagnoses is the foundation of health and clinical sciences. Given that reliable biological markers or specific objective signs for most mental health problems are lacking and many medical conditions only show objective markers in late stages, accurate assessments are difficult [1,2]. Essentially, every single measure of a psychological construct has some potential source of bias (e.g., self-report and recall bias) or can be seen as

fallible in some respect [3,4] – which can result in inaccurate assessments and delayed treatments.

The absence of a single error-free measure can be addressed by involving multiple experts reviewing several sources of information to form a *best-estimate assessment* or a reference standard [5–7]. To understand the quality of such an assessment, it is crucial to understand how it was reached (i.e., the quality of the specific assessment method used). However, the quality of best-estimate assessments is typically very difficult to evaluate due to poor reporting of the assessment

* Corresponding author at: Allhelgona kyrkogata 16a, 223 62 Lund, Sweden.
E-mail address: veerle.eijsbroek@psy.lu.se (V.C. Eijsbroek).

method, and when the method is reported, the reporting quality varies substantially [7]. Here, we tackle this problem by developing a guideline for how to report assessment methods that aim to achieve such best-estimate assessments, i.e., where experts review several sources of (longitudinal) information to achieve a more accurate assessment than a single, error-prone measure.

1.1. Assessment

Assessment includes the evaluation, integration, and interpretation of several sources of information (e.g., outcomes of different measures, tests, or scans) to derive a valid and reliable decision (e.g., a best-estimate diagnosis) [8]. Accurate assessments are crucial. In clinical practice, under- or over-estimation of illnesses can have severe negative impacts on people's lives. In research, inaccurate assessments threaten the validity of scientific results. For policy and implementation development, assessments are the basis for guideline development and the economic and societal evaluations of interventions. Furthermore, obtaining more accurate assessments has become increasingly important considering that high-accuracy assessments are needed in diverse fields such as Biological Psychiatry (e.g., to find reliable biomarkers linked to reference standard assessments [9–11]) and Artificial Intelligence (e.g., to train models to reference standard assessments [12–14]).

1.2. A methodological solution

Here, we connect three bodies of literature that have proposed similar assessment methods: *The Expert Panel* method in medicine [7,15,16] – as well as the *Best-Estimate Diagnosis* [6,17] and the *Longitudinal Expert All Data* (LEAD [5]) methods in clinical psychology and psychiatry. The three methods share the same goal of attaining best-estimate assessments through similar methodological approaches: All three methods use expert panels or consensus teams (e.g., clinical psychologists or medical doctors) to review several sources of information (e.g., clinical questionnaires and medical tests) to establish a more accurate assessment (e.g., a best-estimate diagnosis).

The *Best-Estimate Diagnosis* method was introduced by Leckman et al. [6] as a strategy to set accurate lifetime psychiatric diagnoses. The method focuses on two components, namely 1) using *all data* (e.g., information from medical records and relatives in addition to interviews) that is 2) evaluated by *expert clinicians* (who review all data and then reach a consensus [17]). Consequently, Spitzer [5] proposed the *Longitudinal Expert All Data* (LEAD) method to obtain a criterion or reference standard to validate the Diagnostic Interview Schedule [18] for setting psychiatric diagnoses. LEAD extends the Best-Estimate Diagnosis method and involves three essential components, namely *Longitudinal data collection* (i.e., not limited to a single examination at one point in time), *Expert evaluation* (i.e., the diagnoses are set by expert clinicians), and *All Data* (i.e., the experts have access to multiple data sources). A similar approach was proposed in medicine, where the *Expert Panel* method was developed as a solution for diagnostic accuracy studies with an imperfect or missing reference standard [16,19,20]. Here, a *panel of experts* decides on a medical condition based on *all relevant information*. A review of expert panel studies [7] identified four critical components of the expert panel design, namely 1) the panel constitution, 2) the information presented to the panel, 3) the decision process of the panel, and 4) the validity of the panel diagnosis.

So, all three methods employ a similar approach to obtain *best-estimate assessments* (e.g., for diagnostic purposes or as a reference standard) while accentuating parts of it: The Best-Estimate Diagnosis method accentuates the use of informants and objective tests next to self-reported data [6,17]; the Expert Panel method focuses on the characteristics, constitution, and procedure of the panel [7,15] and only the LEAD method requires a longitudinal design [5], although longitudinal data are also used in some Expert Panel designs ($\approx 27\%$ of studies [7]). Herein, we collectively refer to these three approaches as the *assessment*

methods.

The result of the *assessment methods* is a consensually derived criterion (e.g., a best-estimate assessment) that has been used for many different applications where there is no single error-free measure. It has, for example, been used to i) evaluate the accuracy of a measurement tool or marker through comparison to a best-estimate assessment [21–26]; ii) establish the prevalence of symptoms and disorders [27–29]; iii) establish the temporal stability or development of symptoms and disorders [30–32]; iv) improve (earlier) detection or screening of symptoms or disorders [33–35]; v) study genetics and family history [36–38]; and vi) examine classification systems or diagnostic criteria [39–41]. The applications span diverse fields, including medicine, psychiatry, clinical psychology, public health/epidemiology, and artificial intelligence. Box 1 provides more examples of how the best-estimate assessments have been applied in different types of studies across fields.

1.3. Reporting issues

The assessment methods possess a high potential for achieving best-estimate reference standards in many situations. However, the quality of such proclaimed best-estimate assessments varies substantially and is typically very difficult to evaluate due to poor reporting of the method *how* they were achieved (e.g., see reviews of expert panels [7,15]). A systematic review of assessment methods and reporting of expert panels [7] has demonstrated that the methods used for panel or consensus diagnoses vary substantially across studies and that many aspects of the procedure are often unclear or not reported at all (i.e., in 83 % of the reviewed studies). Many recent studies fail to report central aspects of the assessment procedures, including the quality, structure, or presentation of the data [43], the training and qualifications of the experts [44], the method for avoiding biases and achieving consensus [45], and the time span of the longitudinal design-component [46]. The poor operationalization of the assessment methods jeopardizes the goal of achieving best-estimate assessments – where a vaguely described method makes it difficult to evaluate the research. Referring to an assessment as a *best-estimate* (and sometimes even as a *gold standard*) while vaguely describing or poorly operationalizing the method for achieving the assessment is alarming [47,48].

1.4. The degree of validity

These assessment methods aim to achieve high validity (i.e., the degree to which the assessment captures what it aims to measure). Typically, the assessment methods aim to achieve as high validity as possible (i.e., a “leading” assessment) or, depending on resources, at least more accurate than a single error-prone measure. Despite this central aim, research often fails to clearly describe the degree of validity of the attained assessment. Using these assessment methods does not automatically guarantee high validity – it depends on how well the method is executed.

In addition, the derived assessments are often described with different terms: *reference standard* is often used in medicine, and *criterion standard* or *best-estimate diagnosis* is often used in psychology. We propose that the reporting of these assessment methods benefit from more explicitly describing *what* was measured and *how* well it measures up to different standards – whether and how they relate to a state-of-the-art assessment. Whereas *reference* and *criterion standards* fail to convey an intention of “nearing” a state-of-the-art assessment, the *best-estimate diagnosis* narrowly focuses on the classification of a diagnosis and not on symptom severity. Therefore, we here use the term *best-estimate assessment* in the context of describing a “leading”, state-of-the-art assessment.

1.5. Reporting standards

Previous well-established guidelines have focused on the complete reporting of specific study designs, such as the *STrengthening the*

Box 1

Overview of applications of the assessment methods across different study designs and fields

The absence of a single error-free measure can be mitigated by involving multiple experts reviewing several sources of information to form a best-estimate assessment. This assessment method is used across clinical fields, such as clinical psychology, psychiatry, medicine, and epidemiology. The use cases below highlight different study designs and fields where the assessment methods are applied. For example, best-estimate assessments are used:

1. To evaluate a measure's accuracy against a reference standard. To understand the accuracy of a measurement tool, there is a need to compare it to a more accurate or best-estimate assessment. For example, it has been used:

- *in psychiatry*, for evaluating MINI-KID diagnoses for children and adolescents [21] and evaluating DSM diagnoses in patients with psychosis [22].
- *in clinical psychology*, for evaluating Major Depression Inventory severity scores [23].
- *in medicine*, for evaluating deep learning models assessing liver cancer [24] and evaluating prediction rules for coronary artery disease [25].
- *in epidemiology*, for evaluating electronic health record algorithms for assessing asthma [26].

2. To establish the prevalence of symptoms or disorders. For example, it has been used:

- *in epidemiology*, for assessing the prevalence of pathological gambling [27].
- *in psychiatry*, for assessing the prevalence of eating disorders in personality disorder patients [28].
- *in medicine*, for assessing the prevalence of clinically relevant findings when diagnosing pulmonary embolism [29].

3. To establish the temporal stability or development of symptoms or disorders. For example, it has been used:

- *in clinical psychology*, for learning about autism spectrum disorder diagnoses during childhood [30] and for learning about the course of bipolar disorder [31].
- *in psychiatry*, for assessing diagnostic stability in individuals with autism spectrum disorder [32].

4. To improve (earlier) detection or screening of symptoms or disorders. For example, it has been used:

- *in psychiatry*, for assessing personality disorders [33].
- *in medicine*, for the early detection of heart failure [35] and for early detection of injuries in physically abused older adults [34].

5. To study genetic history and family heritability. For example, it has been used:

- *in psychiatry*, for learning about genetic risks for ADHD [36].
- *in clinical psychology*, for studying familiarity and heritability of depression subtypes [37] and for studying the familial transmission of mania and depression [38].

6. To examine classification systems or diagnostic criteria. For example, it has been used:

- *in clinical psychology*, for evaluating the DSM criteria for hoarding disorder [39].
- *in psychiatry*, for examining a DSM alternative model for personality disorders [40] and for comparing DSM-IV and – 5 criteria of autism spectrum disorders [41].

Description of an example study including a best-estimate assessment

Best-estimate assessments have, for example, been established to evaluate the validity of the Schedule for Affective Disorders and Schizophrenia for School-Age Children-Present and Lifetime version (K-SADS-PL; a semi-structured diagnostic interview used in child and adolescent psychiatry [42]). The best-estimate assessments were diagnoses of neurodevelopmental and related disorders made by five experienced child psychiatrists based on the DSM-5 criteria. To achieve best-estimate assessments, patients were followed for at least three months, and the psychiatrists had access to all available data (except for the K-SADS-PL diagnoses), including information from medical records, interviews, questionnaires, laboratory tests, as well as information provided by clinical staff, caregivers and teachers. Criterion validity of the K-SADS-PL was established as the agreement of the diagnoses with the best-estimate assessments.

Reporting of *OB*servational studies in *Epidemiology* (STROBE [49]) for observational studies; the *Statement for Reporting for Diagnostic Accuracy* (STARD [50]) for diagnostic accuracy studies; the *Consolidated Standards of Reporting Trials* (CONSORT [51]) for randomized trials, and the *Transparent Reporting of a multivariable prediction model for Individual Prognosis Or Diagnosis* (TRIPOD+AI [52]) for prediction model studies (see the supplementary material [SM] for other relevant guidelines). The STARD guidance is most closely related to the reporting of the assessment methods since best-estimate assessments are often used to evaluate a measure's (diagnostic) accuracy. However, none of the guidelines are sufficient for complete reporting of the assessment methods, where

(multiple) experts review several sources of (longitudinal) information to form a best-estimate assessment. Although an earlier systematic review identified and structured the various choices involved in Expert Panel procedures [7], no attempt was made to develop a formal guideline for the reporting of Expert Panel assessments.

1.6. Aim

Our aim is to develop reporting standards for comprehensive reporting of Expert Panel, Best-Estimate Diagnosis, and LEAD methods – which can help researchers plan and report studies employing these

assessment methods, as well as help readers evaluate them. We call the reporting guideline the LEADING guideline, emphasizing the methodological components and the importance of describing *what* is assessed and *how well* (i.e., how it relates to a “leading” assessment). The individual reporting standards are divided into four groups according to the components of LEAD (*Longitudinal, Expert, All Data* [5]), from which we revised the original meanings to *Longitudinal, Evaluation – experts, materials and procedures, Appropriate Data, and Validity*. In short, the LEADING guideline aims to guide the reporting of assessment method to improve evaluations of the assessment standard.

2. Methods

We developed the reporting guideline over four stages: 1) drafting reporting standards, 2) incorporating expert feedback, 3) refining the final guideline, and 4) testing the guideline. The development method largely followed Moher and colleagues’ guidance for developing reporting guidelines [53] (See Table S1 for elaborations on each recommended step). For organizational purposes, a working group (V.E., K.K., & O.K.) was set up, and a steering group (H.A.S., J.B., E.F., D.K., P.G., I.A., & P.B.) was formed to provide a wide range of expertise. The steering group included seven experts and was selected to cover diverse expertise and fields related to the assessment methods (e.g., psychiatry/clinical psychology, medicine, epidemiology/public health, and Artificial Intelligence). Information regarding ethics is presented after the discussion.

2.1. Drafting reporting standards

The working group, with the support of the steering group members, identified relevant research using or describing the assessment methods, including the three bodies of literature: Expert Panel [7], Best-Estimate Diagnosis [6], and LEAD [5]. In addition, articles using any of the three assessment methods were identified through a literature search using Google Scholar with the following search terms: [“expert panel diagnosis” OR “expert panel assessment” OR “expert panel consensus” OR “expert panel methodology” OR “expert panel standard” OR “expert panel reference”] for Expert Panel studies; [“best-estimate diagnosis” OR “best-estimate diagnostic” OR “best-estimate standard” OR “best-estimate assessment” OR “best-estimate methodology” OR “best-estimate reference”] for Best-Estimate Diagnosis studies; and [“longitudinal expert all data” OR “longitudinal evaluation all data”] for LEAD studies. Articles that clearly stated using one of the three assessment methods were selected. Articles stating another purpose than assessment (e.g. when an expert panel was used to reach a consensus about a treatment strategy) were excluded.

Furthermore, relevant reporting guidelines and systematic reviews were identified, including a review of expert panel applications [7], the STROBE statement [54], and the STARD guidance [50]. Other complementary reporting guidelines and systematic reviews are presented in the SM. The aim was for the reporting standards in the LEADING guideline to complement rather than repeat them (i.e., new standards should extend or complement existing standards rather than repeat them [53]). For example, when reporting a randomized trial that includes best-estimate assessments, one may use CONSORT [55] to report the trial design and main results, the LEADING guideline for describing the specifics for reaching the best-estimate assessments, and the *Consolidated Health Economic Evaluation Reporting Standards* (CHEERS) [56] for reporting the economic evaluations and comparisons.

Potential standards were drafted by the working group with the objective of encompassing a comprehensive reporting of the assessment methods. The reporting standards were grouped into four groups: *Longitudinal design, Appropriate data, Evaluation – experts, materials and procedures, and Validity*. Empirical and theoretical inclusion rationales were stated for the groups and the individual standards (i.e., explanations and elaborations). Lastly, the standards with inclusion rationales

were further developed through a workshop with a patient organization for depression, followed by feedback from the steering group members to receive a wide range of perspectives early in the process.

2.2. Incorporating expert feedback

To systematically collect expert feedback from different perspectives, we used a consensus-building procedure called the *Delphi technique* [57]. We used an iterative process based on two rounds of questionnaires (i.e., Delphi surveys), enabling feedback from round 1 to feed into round 2. Delphi participants received relevant background research, the reporting guideline aims, and the reporting standards with their inclusion rationales. They provided feedback through open- and closed-ended response formats. Through open-ended responses in Round 1, experts could propose new standards and provide feedback on the formulations of existing standards and their inclusion rationales. In addition, two closed-ended questions [58] about standard inclusion (*This item should be included in the reporting checklist*) and perception of study quality (*Whether this information is present or not would influence my perceptions of the quality of a study*) were answered with rating scales ranging from 1 = *Strongly disagree* to 7 = *Strongly agree*. In Round 2, the experts were asked to rate the updated reporting standards using the same two closed-ended questions as in Round 1 and to provide feedback on the clarifications and reformulations through open-ended responses.

To recruit Delphi participants, the first and/or last authors of articles since 2013 using any of the three assessment methods were identified using the search terms described above ($n = 87$ articles, $n = 124$ authors; see the SM for more details). These authors and the seven steering group members were invited via email to participate in the Delphi Round 1 ($n = 131$ participants emailed). In total, 27 participants completed the survey (response rate 21 %). Only participants from Round 1 who provided their contact details were invited to Round 2 ($n = 25$). In total, 20 participants completed the survey (response rate 80 %). All participants provided their informed consent. Fig. 1 presents the research experiences and demographics of the Delphi participants. Participants reported a wide range of academic backgrounds (e.g., Clinical Psychology, Psychiatry, Medicine, Artificial Intelligence, Journal Editors) and an extensive variety of relevant methodological experiences (e.g., Ecological Momentary Assessments, Biological Markers, and Expert Panels; Fig. 1), with an age range of 30–70 years ($M = 51.54$, $SD = 12.40$).

2.2.1. Delphi survey results

The criteria for including a reporting standard was that the median of Delphi expert responses was at least 6 = *Agree* on the question about its inclusion. In Round 1, the mean ratings for the *item inclusion* scale ranged from 5.37 to 6.67 ($M = 6.06$; $SD = 0.31$; Table S2) with a median agreement ranging from 6 = *Agree* to 7 = *Strongly Agree*. No new standards were suggested. The feedback resulted in the removal of one reporting standard and the clarification and reformulation of 20 standards. The standard on *Transparency and replicability* was rated as relevant but removed because it is achieved by reporting the other reporting standards. Standard 4.2 *Validity and Standard* needed a major clarification about the meaning of validity as well as standard. Minor clarifications and reformulations, such as grammar or word changes, were made for 19 standards (see open material). The mean ratings in Round 2 ranged from 5.47 to 6.70 ($M = 6.20$; $SD = 0.37$; Table S3), with open feedback resulting in minor clarifications and reformulations of nine standards.

2.3. Refining the guideline through expert consensus

The authors finalized the guidelines in an expert consensus meeting. The meeting was held online with nine working and steering group members. The content and structure of the consensus meeting were prepared by the working group, and the meeting was led by the last author (O.K.). Participants had access to the guidelines, inclusion

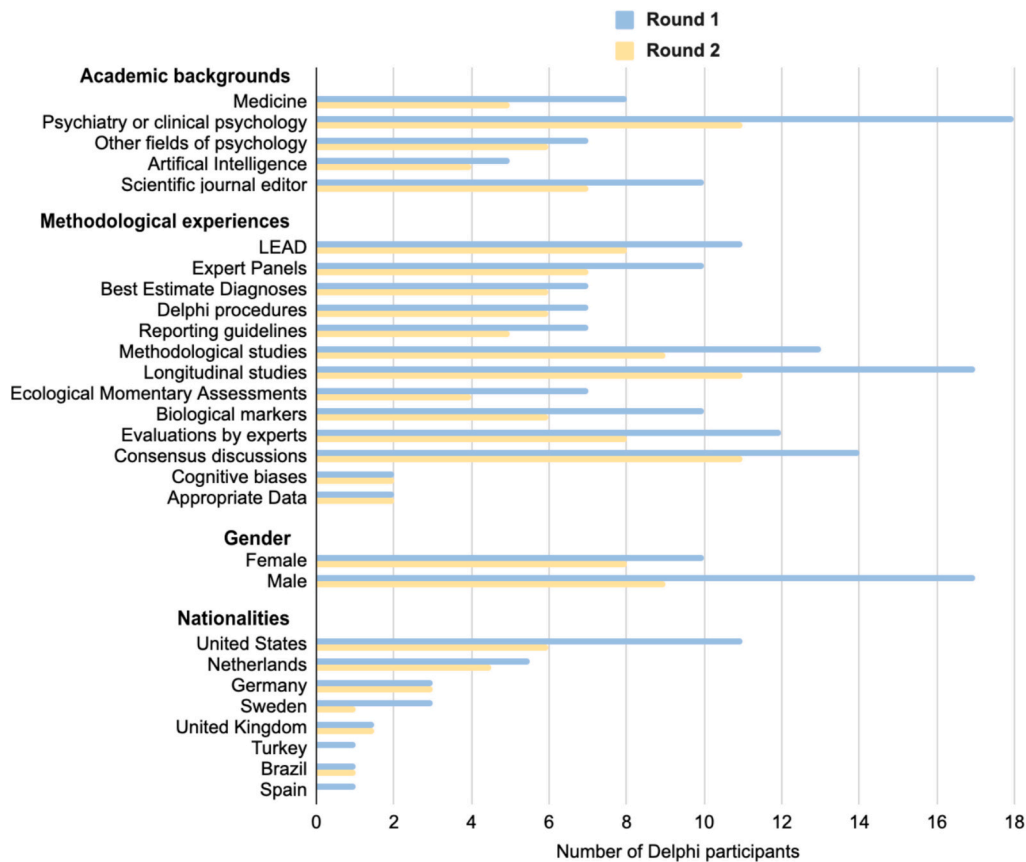


Fig. 1. Research experiences and demographics of the Delphi participants.

Notes. The answer options for Academic backgrounds and Methodological experiences were not mutually exclusive (i.e., multiple backgrounds and experiences could be reported by the participants). In Round 2, the demographics and reported experiences are known for 17 of the 20 participants.

rationales (i.e., elaboration and explanation), and the drafted paper before the meeting, where they also had the option to provide comments and feedback in writing. The meeting included reviewing the Delphi Rounds 1 and 2 findings and discussing the paper draft, including the individual reporting standards and groups. We decided not to carry out another Delphi round since i) the median agreement for each reporting standard in both Delphi Rounds 1 and 2 ranged from 6 = *Agree* to 7 = *Strongly Agree*, ii) no new standards were suggested, and iii) only minor changes were needed after Round 2, which taken together suggest consensus.

2.4. Testing the guideline

To test the applicability of the guideline, the guideline was tested i) by independent researchers with experience of each method piloting the reporting of each standard and ii) by the authors (V.E., K.K.) applying it to published articles. The two test procedures resulted in adding minor clarifications to three standards (2.4 *The access to the index measure*, 3.3 *Blindness and conflict of interest*, and 3.4 *Instructions and training*). Also, a concrete example of how to report the items was added to the guideline instructions.

2.4.1. Incorporating test-user feedback

Two test users (PhD, with experience using the LEAD and Expert Panel method) who had not been involved in the development of the guideline (e.g., in the Delphi procedure) were recruited to pilot the guideline (see the SM for more details). They were asked to report each standard based on a finished, ongoing or planned study using one of the

assessment methods and/or provide feedback about the formulation of the standards.

2.4.2. Applying the standards to published studies

Three separate targeted searches (LEAD, Expert-panel, Best-estimate) were conducted using the search terms described above. The first author (V.E.) examined which standards were reported in 30 randomly selected articles applying the assessment methods in 2022 and 2023 (i.e., five from each method from each year; see the SM for the selection process). Each reporting standard was rated using four categories: standard *not reported*; standard *reported vaguely or insufficiently*; standard *(minimally) sufficiently reported*; or standard *not applicable to the study*. Out of the 30 articles, six were randomly selected (i.e., one from each method from each year) and examined by the second author (K.K.) to get insight into the accuracy of the ratings of the first author. Discussing their disagreements to reach consensus resulted in changing 23 ratings (19 %) of the first author. This testing procedure also provided information about the strengths and shortcomings of contemporary reporting of published articles using the assessment methods (see Results section).

3. Results

The reporting guideline is presented in Table 1 (see Fig. 2 for an overview). It comprises 20 standards for comprehensive reporting of the assessment methods divided into four groups: 1. *The Longitudinal design* group (4 standards), 2. *The Appropriate data* group (4 standards), 3. *The Evaluation – experts, materials, and procedures* group (10 standards), and 4. *The Validity* group (2 standards). The reporting standards encourage

Table 1
The LEADING guideline reporting standards.

Group	#	Reporting standards
Longitudinal Design <i>Report the longitudinal design, by describing:</i>	1.1	The time period. The data collection period covered for each participant (i.e., start and end of the data collection) and to what extent the length is sufficient for capturing the targeted symptoms. <i>For example, the weeks/months a participant is followed and how this matches the criteria for the targeted disease/disorder.</i>
	1.2	The number of time points. Whether and how data were collected on multiple occasions between the start and the end of the time period, the sufficiency of the data collection, and of its frequency and intensity for capturing the target. <i>For example, report the number of check-ins with the participants and the included measures for each assessment.</i>
	1.3	History or lifetime information. Whether and which data from before the start of the data collection were taken into account and how these data are relevant for the assessment of the target. <i>History or lifetime data may include self-report of medical history, childhood memory accounts, or other-than-self information such as from relatives or from medical records.</i>
	1.4	The targeted time point(s) of the experts' assessment. The time point(s) for which the experts provide their assessment, on which time period the data of the assessments are based (i.e., past data, future data, or both), and justifications for the targeted time point(s). <i>For example, the experts can assess the presence of a diagnosis at the start of the study and thus base their assessment on future data from that reference point; or in the middle of the study time period and thus have access to both past and future data from that reference point.</i>
Appropriate data <i>Report the appropriate-ness of the data, by describing:</i>	2.1	The type and quality of the data. The type, quality, and relevance of the data and why these data sources are sufficient and suitable for capturing the target. <i>For example, describe the validity and reliability of the data – and how it relates to capturing the targeted construct.</i>
	2.2	The data triangulation. Whether and why the data come from different methodological approaches and the degree to which these approaches complement each other. <i>For example, how self-reported data is complemented by objective/physical tests and/or other informant data.</i>
	2.3	The data presentation. How the data were structured and presented to the experts for their assessments and why. <i>For example, were the data presented in a case report; and was the information presented with or without any interpretation?</i>
	2.4	The access to the index measure. For an assessment accuracy study, the extent the experts had access to the index measure and why (i.e., an assessment that is being compared to the best-estimate assessment), and how its information was weighted in their assessment. <i>For example, were experts blind to the measure (its outcome and/or its raw data) that is being validated?</i>
Evaluation – experts, materials and procedures <i>Report the evaluation experts, materials and procedures, by describing:</i>	3.1	The expert and panel characteristics. The characteristics of the experts and the panel, as well as how these characteristics are relevant for assessing the target. <i>Relevant characteristics may include clinical and research experiences, professions, education, and demographics.</i>
	3.2	The number of experts and panels. The total number of experts and panels, and how many experts/panels were assessing each case and why. <i>For example, how many and are the same expert(s)/expertise(s) present in every assessment?</i>
	3.3	Blindness and conflicts of interest. Whether and to what extent the experts are blind to the research aims and/or have any conflicts of interest. <i>This may include experts' study authorship or the experts' relationship to the index measure or any other assessment method. If the study examines an index measure (i.e., an assessment that is being compared with the best-estimate assessment), declare the authors' as well as the experts' relationship to it.</i>
	3.4	Instructions and training. The instructions, training, and/or preparation that the experts specifically received for this assessment task and why they did or did not receive this. <i>For example, provide information regarding 1) whether the assessment method and procedure are kept standardized across the individual assessments, 2) the methods to ensure experts' preparedness for the assessment, or 3) any specific measures to limit biases.</i>
	3.5	The assessment procedure. The procedure that the experts followed for their assessment. <i>For example, describe whether there was a standardized procedure and what this procedure included (such as following clear diagnostic criteria).</i>
	3.6	The assessment response format. The response format used by the experts for their individual assessments, what it included, and how it was structured. <i>For example, describe any assessment sheet, including assessment questions and answer options.</i>
	3.7	The data combination method. The method or guidelines for how the data should be weighted, judged, and combined by the individual experts to reach a conclusion in their individual assessment. <i>For example, should any data sources be evaluated first or weighted more strongly; or are the experts asked to assess certain diagnostic criteria/symptoms first, before forming a final diagnosis?</i>
	3.8	Independent expert assessments. Whether and how the experts first evaluated the data individually and made their first individual assessments independently. <i>For example, how it was ensured the experts first reviewed the data individually/independently before discussing their assessment outcome with the other panel members.</i>
	3.9	The inter-rater and inter-panel reliability. The inter-rater/inter-panel reliability, how it was calculated and evaluated, or why it was not possible to calculate it. <i>For example, which reliability metric was used and over how many experts/panels and cases the reliability was calculated.</i>
	3.10	The solution to disagreements. The approach for solving (any) disagreements between the individual expert assessments, the rationale for the chosen approach, and potential problems that may have occurred and how these were assessed. <i>Methods may include reaching a consensus, taking the average, or majority vote. Potential problems may, for example, include power imbalances in the expert panel.</i>
Validity <i>Report what was assessed and how well, by describing:</i>	4.1	The assessment description. Description of what the assessment actually is. <i>For example, is the assessment a diagnosis, symptom severity assessment, course of illness assessment, or treatment response assessment?</i>
	4.2	The validity and standard. Reflect on the degree of validity and describe the standard that the method aims to achieve, how well the assessment method measures up to that degree, and how it compares with current standards. <i>For example, reflect on evidence supporting or against validity aspects such as construct, face, and criterion validity; and state whether the assessment should be seen as a best-estimate assessment standard or an accepted reference standard (see Table S2 for more examples).</i>

Instructions. The LEADING guideline comprises these 20 reporting standards for comprehensive reporting of assessment methods involving expert(s) reviewing several sources of information (over time) to achieve a more accurate assessment (e.g., see Expert Panel, Best-Estimate diagnosis, and Longitudinal Expert All Data methods). The standards aim to help researchers plan and report studies employing these assessment methods, as well as help readers evaluate them. As such, avoid simply answering yes or no to the standards when you instead can (succinctly) describe justifications and courses of action. Ensure the reports of the standards are clear, specific, and justified. To exemplify, standard 1.1 *The time period* could be reported as ‘*The time span was six weeks, which covers more than the two weeks a person should have the symptoms for meeting the criteria for Major Depressive Disorder according to the DSM-5.*’

Not all of the reporting standards will be applicable to all types of studies – however, it is typically better to describe how a standard is not applicable than to leave the information out. Since the guideline covers the reporting of the assessment method, the method section would suit the reporting of most standards in most cases. However, the reporting guideline does *not* standardize where standards should be reported. When standards are considered less relevant or not applicable to a specific study, they can, for example, be described in an Appendix. Since the guideline focuses specifically on the reporting of the assessment method, it is recommended to use a complementary guideline for the reporting of the other study components. Which complementary guideline is dependent on the study type in which the assessment method is employed (e.g., see STARD for diagnostic accuracy studies; STROBE for observational studies; and CONSORT for randomized trials).

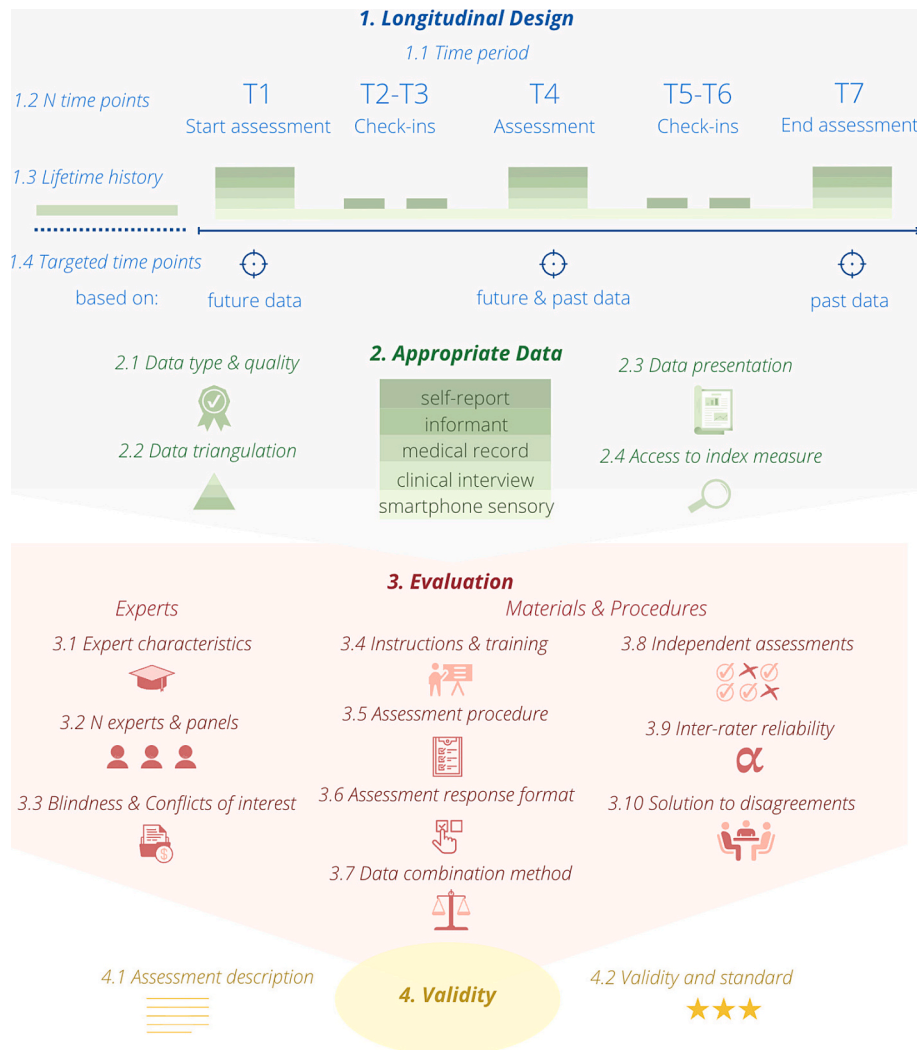


Fig. 2. Overview of the LEADING guideline reporting standards. For more details about each standard, see [Table 1](#).

researchers to elaborate on what was done and why – whilst avoiding normative standards, such as a minimum number of experts. Each standard description in [Table 1](#) is accompanied by an example. Further *Explanations and Elaborations* regarding the individual reporting standards and the four groups are presented in the SM, including Tables S4 and S5. A reporting template for providing an overview of the standard reports can be found on www.leading-guideline.org.

3.1. Applying the LEADING guideline to published studies

Applying the guideline to a random selection of 30 articles indicated severe heterogeneity in *what* of the methods is reported and *how* ([Table 2](#); see the SM for the search strategy). Across the 30 studies, 10 to 63 % (Mean = 33 %) of the standards were *not* reported. Regarding the reporting standards, the type and quality of the data (2.1), the access to the index measure (2.4), the expert and panel characteristics (3.1), the number of experts and panels (3.2), the assessment procedure (3.5), and the assessment description (4.1) were mostly reported (i.e., green in

Table 2
Reports across the LEADING guideline standards in 30 randomly selected articles published in 2022 and 2023

LEADING reporting standards	LEAD*										Expert Panel*										Best-Estimate Diagnosis*											
	2022					2023					2022					2023					2022					2023					% ¹	% ²
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30		
1.1 The time period	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	43	23
1.2 The number of time points	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	40	37
1.3 History or lifetime information	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	43	20
1.4 The targeted time point(s)	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	47	13
2.1 The type and quality of data	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	67	7
2.2 The data triangulation	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	50	10
2.3 The data presentation	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	10	70
2.4 The access to index measure	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	75	0
3.1 The expert and panel characteristics	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	70	3
3.2 The number of experts and panels	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	67	30
3.3 Blindness and conflicts of interest	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	20	23
3.4 Instructions and training	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	13	73
3.5 The assessment procedure	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	70	7
3.6 The assessment response format	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	40	37
3.7 The data combination method	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	0	83
3.8 Independent expert assessments	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	33	47
3.9 The inter-rater and inter-panel reliability	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	31	52
3.10 The solution to disagreements	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	46	36
4.1 The assessment description	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	97	0
4.2 The validity and standard	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	Green	13	63
green %	60	26	16	16	25	37	58	26	53	47	50	70	45	40	70	45	50	45	55	50	42	65	42	40	16	35	37	26	70	37	32 ¹	
red %	30	47	58	53	45	26	22	63	26	32	20	15	30	35	15	15	35	45	20	20	32	18	32	30	37	35	42	53	10	47	33 ²	

Notes. red = not reported; orange = insufficiently reported; green = (minimally) sufficiently reported; gray = not applicable to report. ¹ = mean of reporting standards; ² = mean of studies. * Ten articles for each assessment method were randomly selected (see SM for selection process): 1 [59]; 2 [60]; 3 [61]; 4 [62]; 5 [63]; 6 [64]; 7 [65]; 8 [66]; 9 [67]; 10 [68]; 11 [69]; 12 [70]; 13 [45]; 14 [71]; 15 [72]; 16 [73]; 17 [74]; 18 [75]; 19 [76]; 20 [77]; 21 [78]; 22 [79]; 23 [80]; 24 [81]; 25 [82]; 26 [83]; 27 [84]; 28 [85]; 29 [86]; 30 [87]. The 30 articles cover different disciplines and fields, including Clinical Psychology or Psychiatry (19 articles; 63%), Medicine (8 articles; 27%), and Artificial Intelligence (3 articles; 10%).

more than 50 % of the studies). However, the data presentation (2.3), the instructions and training (3.4), the data combination method (3.7), the inter-rater and inter-panel reliability (3.9), and the validity and standard (4.2) were *not* reported at all in the majority of the studies (i.e., red in more than 50 % of the studies). Considering that most changes that resulted from discussing disagreements between the raters (V.E. and K.K.) were from green to orange, and that green refers to a (minimally) sufficiently reported and orange to insufficiently reported, this suggests that the table is conservative in regards to the severity of the current state of poor reporting (i.e., potentially showing a more positive picture; for more information see the SM).

4. Discussion

Our objective was to develop a guideline that supports comprehensive reporting of assessment methods collecting longitudinal, appropriate data that experts evaluate to achieve an assessment that is more accurate than a single error-prone measure. This assessment method is known as Expert Panel in medicine, and Best-Estimate Diagnosis or LEAD in psychiatry and clinical psychology. Given that reliable biological markers or specific objective signs are lacking not only in mental health but also in some medical conditions, the assessment approach—and this guideline—have wide applicability across diverse clinical domains (see Box 1). The aim of the LEADING guideline is to help researchers plan, conduct, report, and evaluate the assessment method-

related elements of this study design.

The LEADING reporting standards were established through an open process, incorporating relevant empirical evidence and methodological work, complementary reporting guidelines, and comprehensive iterations of expert feedback and patients' perspectives. As this guideline focuses on the assessment methods, we recommend that researchers also rely on established guidelines for other parts of their research, such as sampling and other epidemiological aspects (e.g., STROBE, CONSORT, and STARD).

4.1. Limitations

We have connected three assessment methods with similar approaches from related fields and drafted applicable reporting standards. We presented the rationale for these three methods and each reporting standard with supporting evidence in the Delphi survey for review, which did not bring up additional methods or reporting standards. As we did not carry out a systematic literature review, we cannot exclude the existence of other assessment methods with similar approaches. We welcome any suggestions about similar methods to which the guideline is applicable.

The Delphi survey participants and the author group had a wide range of experiences and backgrounds; however, geographically, Europe and North America were the most common, whereas several areas were not represented. The Delphi participants were the first or last authors of

studies employing the assessment methods. The quality of the articles and the education or experience of the authors were not taken into account as selection criteria (although it was self-reported, as presented in Fig. 1). The number of Delphi participants (27 in Round 1, 20 in Round 2) is relatively small compared to some other standard developments (e.g., 73 in the development of STARD [88]) but comparable to others (e.g., 24 for the development of the TRIPOD statement [89]). Even though the response rate in Round 1 (21 %) can be considered low, the number of participants was sufficient to cover a broad range of academic backgrounds, methodological experiences, and demographics (Fig. 1). The same limitation is applicable to the size of the steering group ($n = 7$) and the current test-user group ($n = 2$).

4.2. Implementation, adherence, and evolution

4.2.1. Implementation

We encourage implementation of and adherence to the LEADING guideline via www.leading-guideline.org, scientific journals, editorials, commentaries, and the *Enhancing the Quality and Transparency Of health Research* (EQUATOR) Network (www.equator-network.org). The LEADING guideline aims to support authors in writing their research reports, editors and peer reviewers in reviewing submitted reports, and readers in critically evaluating published reports. We encourage editors and publishers to support adherence to the LEADING guideline by referring to it in author guidelines. We recommend that authors submit the guideline as an appendix to their manuscripts (see www.leading-guideline.org for a reporting template). We also encourage dissemination of the guideline via inclusion in research seminars and method courses in clinical studies. Teaching early career researchers about the benefits of comprehensive reporting can promote adoption and adherence, for example, by requiring students to write theses in accordance with the applicable guideline. Further dissemination is encouraged by welcoming translations of the guideline into different languages.

4.2.2. Adherence

Reporting guidelines have become widely available for different study designs, especially in medicine. However, it is important to measure adherence to the guidelines, including identifying barriers and opportunities, and evaluate their impact on reporting quality [90,91]. Potential adherence barriers include prolonged reporting time, especially when multiple guidelines are needed for the report of the complete study design. However, standardized templates, as well as training and repeated practice, can increase efficiency and facilitate adherence [90,91]. We plan to measure adherence to the LEADING guideline via, for example, standardized adherence assessment forms [92] or AI-based tools that are currently being developed for determining reporting guideline compliance [93,94].

4.2.3. Evolution

The LEADING guideline should be regarded as an evolving reporting guideline requiring ongoing evaluation, refinement, and revision. Methodological components of assessment methods evolve: The LEADING guideline will be periodically updated to correspond to the state-of-the-art of *Expert Panel*, *Best-Estimate Diagnosis*, and *LEAD* methods. We encourage readers to provide recommendations for improvements by emailing the corresponding author. Future modifications will be published on the website and aim to reflect feedback and new evidence, ultimately aiming to improve the reporting quality of the assessment methods.

4.3. Conclusions

The LEADING guideline emphasizes the transparent reporting of the methodological components of Expert Panel, Best-Estimate Diagnosis, and LEAD designs and the importance of reporting *what* was assessed and *how* well. Considering the increasing need for high-accuracy

assessments in diverse fields, we hope that the LEADING guideline will be useful in assisting researchers in planning, carrying out, reporting, and evaluating research that aims to achieve accurate assessments.

Ethics approval

Swedish law (2003:460) and the Swedish Ethical Review Authority state that only research that includes i) collecting personal information, and/or ii) that involves “obvious” (uppenbar) risk for physical or psychological harm, or iii) involves manipulating or deceiving individuals, should undergo an external ethical review. Considering that the current research involves asking participants to rate and comment on the reporting recommendations, this research does not fall within these criteria. The participants provided their informed consent and their individual open-ended comments and closed-ended ratings in the open material are anonymised. The closed-ended ratings per reporting standard in the Supplementary Material are presented on group level. The study is deemed exempt from requiring ethical approval according to Swedish Law (see §3–4 of the Act [2003:460] on ethical review of research involving humans in Sweden). Hence, the research should not be reviewed by the Swedish Ethical Review Authority. More information can be found at the Swedish Ethical Review Authority (<https://etikprovningensmyndigheten.se/>).

Open Science

Open data (Delphi rounds 1 and 2), code (analyses), and material (surveys) can be found on the Open Science Framework: <https://osf.io/fkv4b/>.

Public Health Significance

Accurate assessments of symptoms and illnesses are essential for health research and clinical practice but face many challenges. We have developed the LEADING guideline, including reporting standards that assist in planning, reporting, and evaluating methods involving experts reviewing several sources of (longitudinal) information to achieve *best-estimate assessments* in psychology, psychiatry, and medicine.

Patient and public involvement

Prior to the Delphi procedure and test procedures, the reporting standards with inclusion rationales were discussed in an online workshop with a patient organization for depression (the chairman and vice chairman from *Libra Balans Skåne*), followed by receiving feedback from the steering group members to receive a wide range of perspectives early in the process.

Funding sources

V.C. Eijsbroek, K. Kjell, and O. Kjell received funding from FORTE (2022-01022) and H.A. Schwartz from the National Institutes of Health (Grant R01 AA028032-01).

CRediT authorship contribution statement

Veerle C. Eijsbroek: Writing – review & editing, Writing – original draft, Visualization, Resources, Project administration, Methodology, Investigation, Formal analysis, Data curation, Conceptualization. **Katarina Kjell:** Writing – review & editing, Writing – original draft, Resources, Project administration, Methodology, Investigation, Funding acquisition, Conceptualization. **H. Andrew Schwartz:** Writing – review & editing, Resources, Methodology. **Jan R. Boehnke:** Writing – review & editing, Resources, Methodology. **Eiko I. Fried:** Writing – review & editing, Resources, Methodology. **Daniel N. Klein:** Writing – review & editing, Resources, Methodology. **Peik Gustafsson:** Writing – review &

editing, Resources, Methodology. **Isabelle Augenstein**: Writing – review & editing, Resources, Methodology. **Patrick M.M. Bossuyt**: Writing – review & editing, Resources, Methodology. **Oscar N.E. Kjell**: Writing – review & editing, Writing – original draft, Supervision, Resources, Project administration, Methodology, Investigation, Funding acquisition, Formal analysis, Conceptualization.

Declaration of competing interest

O. Kjell and K. Kjell have co-founded and hold shares in a start-up using computational language assessments to diagnose mental health problems based on best-estimate assessments and J.R. Boehnke is as editor part of the International Society for Quality of Life Research.

Acknowledgements

Patient organization members

E. Sellberg (Chairman of Libra Balans Skåne, board member of Balans National) and **I. Odenbrand** (Vice Chairman of Libra Balans Skåne).

Participants in Delphi Round 1.

I. Augenstein (Professor, Computer Science, University of Copenhagen), **S. Aydin** (PhD, Developmental and Educational Psychology, Leiden University), **E. Billstedt** (Professor, Neuroscience and Physiology, University of Gothenburg), **J.R. Boehnke** (PhD, School of Health Sciences, University of Dundee), **D.W. Black** (MD, Professor, Medicine, University of Iowa Carver College of Medicine), **B. Cannell** (Associate Professor, Public Health, University of Texas Health Science Center at Houston), **G.A. Carlson** (Professor, Psychiatry, Stony Brook University), **K.A.S. Davis** (Researcher, Psychiatry Psychology and Neuroscience, King's College London), **F. Dereboy** (MD, Psychiatry, Aydin Adnan Menderes University), **E.I. Fried** (Associate Professor, Clinical Psychology, Leiden University), **P. Gustafsson** (Associate Professor, Child and Adolescent Psychiatry, Lund University), **R. Handels** (Assistant Professor, Psychiatry and Neuropsychology, Maastricht University), **K. Jenniskens** (Assistant Professor, Clinical Epidemiology, Utrecht University), **C. Klaiman** (Associate Professor, Pediatrics, Emory University), **D.N. Klein** (Professor, Clinical Psychology, Stony Brook University), **M. McCloskey** (Professor, Cognitive Science and Psychology, Johns Hopkins University), **A.C. Miers** (Associate Professor, Developmental and Educational Psychology, Leiden University), **K.G.M. Moons** (Professor, Clinical Epidemiology, Utrecht University), **L. Mosqueda** (Professor, Medicine, University of Southern California), **H.A. Schwartz** (Associate Professor, Computer Science, Stony Brook University), **M. Stein** (Professor, Psychiatry and Behavioral Sciences, University of Washington), **J.G. Tillman** (PhD, Clinical Psychology, Yale School of Medicine), **Y.P. Wang** (MD, PhD, Medicine, University of Sao Paulo Medical School), and **J. Yonashiro-Cho** (PhD, Medicine, University of Southern California).

Participants in Delphi Round 2.

I. Augenstein (Professor, Computer Science, University of Copenhagen), **S. Aydin** (PhD, Developmental and Educational Psychology, Leiden University), **J.R. Boehnke** (PhD, School of Health Sciences, University of Dundee), **G.A. Carlson** (Professor, Psychiatry, Stony Brook University), **K.A.S. Davis** (Researcher, Psychiatry Psychology and Neuroscience, King's College London), **E.I. Fried** (Associate Professor, Clinical Psychology, Leiden University), **P. Gustafsson** (Associate Professor, Child and Adolescent Psychiatry, Lund University), **R. Handels** (Assistant Professor, Psychiatry and Neuropsychology, Maastricht University), **K. Jenniskens** (Assistant Professor, Clinical Epidemiology, Utrecht University), **C. Klaiman** (Associate Professor, Pediatrics, Emory University), **D.N. Klein** (Professor, Clinical Psychology, Stony Brook University), **M. McCloskey** (Professor, Cognitive Science and Psychology, Johns Hopkins University), **A.C. Miers** (Associate Professor, Developmental and Educational Psychology, Leiden University), **K.G.M. Moons** (Professor, Clinical Epidemiology, Utrecht University), **L.**

Mosqueda (Professor, Medicine, University of Southern California), **Y. P. Wang** (MD, PhD, Medicine, University of Sao Paulo Medical School), and **J. Yonashiro-Cho** (PhD, Medicine, University of Southern California).

Test-users.

T. Ivarsson (Associate Professor, Neuroscience and Physiology, University of Gothenburg) and **P.J. Snelling** (PhD, Medicine, Gold Coast University Hospital and Griffith University).

Appendix A. Supplementary data

The supplementary material includes 1) relevant and complementary reporting guidelines and systematic reviews, 2) elaboration on the steps for developing a health research reporting guideline (Table S1; Moher et al., 2010), 3) the search strategies for identifying articles using the assessments methods and for recruiting Delphi participants and test-users, 4) the closed-ended ratings for each reporting standard from the Delphi surveys (Tables S2 and S3), 5) the explanation and elaboration for the reporting standards, including inclusion rationales and empirical evidence (Table S4), and 6) the procedure for applying the reporting standards to studies published in 2022 and 2023. Supplementary data to this article can be found online at [<https://doi.org/10.1016/j.comppsy.2025.152603>].

References

- [1] Hirschtritt ME, Insel TR. Digital technologies in psychiatry: present and future. *Focus Am Psychiatr Publ* 2018;16(3):251–8. Jul.
- [2] Venkatasubramanian G, Keshavan MS. Biomarkers in psychiatry - a critique. *Ann Neurosci* 2016;23(1):3–5. Mar.
- [3] Cronbach LJ, Meehl PE. Construct validity in psychological tests. *Psychol Bull* 1955;52(4):281–302. Jul.
- [4] Lilienfeld Scott O, Sauvigne Kathryn, Lynn Steven Jay, Litzman Robert D, Cautin Robin, Waldman Irwin D. Fifty psychological and psychiatric terms to avoid: a list of inaccurate, misleading, misused, ambiguous, and logically confused words and phrases. *Front Psychol* 2015;6. Aug 1.
- [5] Spitzer RL. Psychiatric diagnosis: are clinicians still necessary? *Compr Psychiatry* 1983;24(5):399–411. Sep.
- [6] Leckman JF, Sholomskas D, Thompson WD, Belanger A, Weissman MM. Best estimate of lifetime psychiatric diagnosis: a methodological study. *Arch Gen Psychiatry* 1982;39(8):879–83. Aug.
- [7] Bertens LCM, Broekhuizen BDL, Naaktgeboren CA, Rutten FH, Hoes AW, van Mourik Y, et al. Use of expert panels to define the reference standard in diagnostic research: a systematic review of published methods and reporting. *PLoS Med* 2013;10(10). Oct 15. (e1001531).
- [8] Hunsley J, Mash EJ. Evidence-based assessment. *Annu Rev Clin Psychol* 2007;3:29–51.
- [9] Niculescu AB, Le-Niculescu H. Precision medicine in psychiatry: biomarkers to the forefront. *Neuropsychopharmacol Intersect Brain Behav Ther* 2022;47(1):422–3. Jan 1.
- [10] IMAGEN Consortium, Quinlan EB, Banaschewski T, Barker GJ, Bokde ALW, Bromberg U, et al. Identifying biological markers for improved precision medicine in psychiatry. *Mol Psychiatry* 2020;25(2):243–53. Feb.
- [11] García-Gutiérrez María Salud, Navarrete Francisco, Sala Francisco, Gasparyan Ani, Austrich-Olivares Amaya, Manzanares Jorge. Biomarkers in Psychiatry: concept, definition, types and relevance to the clinical reality. *Front Psych* 2020;11. May 1.
- [12] Briganti Giovanni, Le Moine Olivier. Artificial intelligence in medicine: today and tomorrow. *Front Med* 2020;7. Feb 1.
- [13] Rajpurkar P, Chen E, Banerjee O, Topol EJ. AI in health and medicine. *Nat Med* 2022;28(1):31–8. Jan.
- [14] Wiemken TL, Kelley RR. Machine learning in epidemiology and health outcomes research. *Annu Rev Public Health* 2020;41(1):21–36.
- [15] Handels RLH, Wolfs CAG, Aalten P, Bossuyt PMM, Joore MA, Leentjens AFG, et al. Optimizing the use of expert panel reference diagnoses in diagnostic studies of multidimensional syndromes. *BMC Neurol* 2014;14:190. Oct 4.
- [16] Reitsma JB, Rutjes AWS, Khan KS, Coomarasamy A, Bossuyt PM. A review of solutions for diagnostic accuracy studies with an imperfect or missing reference standard. *J Clin Epidemiol* 2009;62(8):797–806. Jan 1.
- [17] Klein DN, Ouimette PC, Kelly HS, Ferro T, Riso LP. Test-retest reliability of team consensus best-estimate diagnoses of axis I and II disorders in a family study. *Am J Psychiatry* 1994;151(7):1043–7. Jul.
- [18] Robins LN, Helzer JE, Croughan J, Ratcliff KS. National Institute of Mental Health diagnostic interview schedule. Its history, characteristics, and validity. *Arch Gen Psychiatry* 1981;38(4):381–9. Apr.
- [19] Gagnon R, Charlin B, Coletti M, Sauvé E, van der Vleuten C. Assessment in the context of uncertainty: how many members are needed on the panel of reference of a script concordance test? *Med Educ* 2005;39(3):284–91. Mar.

- [20] Knottnerus JA, Muris JW. Assessment of the accuracy of diagnostic tests: the cross-sectional study. *J Clin Epidemiol* 2003;56(11):1118–28. Nov 1.
- [21] Hogberg C, Billstedt E, Björck C, Björck P, Ehlers S, Gustle L, et al. Diagnostic validity of the MINI-KID disorder classifications in specialized child and adolescent psychiatric outpatient clinics in Sweden. *BMC Psychiatry* 2019;19(1):142. Jan 1.
- [22] Anglin DM, Malaspina D. Ethnicity effects on clinical diagnoses compared to best-estimate research diagnoses in patients with psychosis: a retrospective medical chart review. *J Clin Psychiatry* 2008;69(6):941–5. Jun.
- [23] Bech P, Timmerby N, Martiny K, Lunde M, Soendergaard S. Psychometric evaluation of the major depression inventory (MDI) as depression severity scale using the LEAD (longitudinal expert assessment of all data) as index of validity. *BMC Psychiatry* 2015;15:190. Aug 5.
- [24] Gao Ruitian, Zhao Shuai, Aishanjiang Kedeerya, Cai Hao, Wei Ting, Zhang Yichi, et al. Deep learning for differential diagnosis of malignant hepatic tumors based on multi-phase contrast-enhanced CT and clinical data. *J Hematol Oncol J Hematol Oncol* 2021;14(1):1–7. Sep 1.
- [25] Bösner S, Haasenritter J, Becker A, Heinzel-Gutenbrunner M, Hani MA, Keller H, et al. Ruling out coronary artery disease in primary care: development and validation of a simple prediction rule. *CMAJ. Can Med Assoc J* 2010;182(12):1295–300. Sep 7.
- [26] Cowan KJ, Tandias A, Arndt B, Hanrahan L, Mundt M, Guilbert TW. Defining asthma: validating automated electronic health record algorithm with expert panel diagnosis. *Am J Respir Crit Care Med* 2014;189. Jan 1.
- [27] Black DW, Coryell WH, Crowe RR, McCormick B, Shaw MC, Allen J. A direct, controlled, blind family study of DSM-IV pathological gambling. *J Clin Psychiatry* 2014;75(3):215–21. Mar.
- [28] Reas DL, Rø O, Karterud S, Hummelen B, Pedersen G. Eating disorders in a large clinical sample of men and women with personality disorders. *Int J Eat Disord* 2013;46(8):801–9. Dec 1.
- [29] Hall WB, Truitt SG, Scheunemann LP, Shah SA, Rivera MP, Parker LA, et al. The prevalence of clinically relevant incidental findings on chest computed tomographic angiograms ordered to diagnose pulmonary embolism. *Arch Intern Med* 2009;169(21):1961–5. Nov 23.
- [30] Brian J, Roberts W, Szatmari P, Bryson SE, Smith I M, Roncadin C, et al. Stability and change in autism spectrum disorder diagnosis from age 3 to middle childhood in a high-risk sibling cohort. *Autism* 2016;20(7):888–92. Oct 1.
- [31] Duffy A, Grof P, Goodday S, Keown-Stoneman C. The emergent course of bipolar disorder: observations over two decades from the Canadian high-risk offspring cohort. *Am J Psychiatry* 2019;176(9):720–9. Jan 1.
- [32] Elias R, Lord C. Diagnostic stability in individuals with autism spectrum disorder: insights from a longitudinal follow-up study. *J Child Psychol Psychiatry* 2022;63(9):973–83. Sep 1.
- [33] Pedersen G, Karterud S, Wilberg T, Hummelen B. The impact of extended longitudinal observation on the assessment of personality disorders. *Personal Ment Health* 2013;7(4):277–87. Nov 1.
- [34] Yonashiro-Cho JMF, Gassoumis ZD, Homeier DC, Wilber KH. Improving forensics: characterizing injuries among community-dwelling physically abused older adults. *J Am Geriatr Soc* 2021;69(8):2252–61. Aug 1.
- [35] Oudejans I, Mosterd A, Bloemen JA, Valk MJ, van Velzen E, Wielders JP, et al. Clinical evaluation of geriatric outpatients with suspected heart failure: value of symptoms, signs, and additional tests. *Eur J Heart Fail* 2011;13(5):518–27.
- [36] Mooney MA, Wilmot B, Bhatt P, Nigg JT, Hermosillo RJM, Fair DA, et al. Smaller total brain volume but not subcortical structure volume related to common genetic risk for ADHD. *Psychol Med* 2021;51(8):1279–88. Jun 1.
- [37] Lamers F, Cui L, Hickie IB, Roca C, Machado-Vieira R, Zarate Jr CA, et al. Familial aggregation and heritability of the melancholic and atypical subtypes of depression. *J Affect Disord* 2016;204:241–6. Nov 1.
- [38] Merikangas KR, Cui L, Heaton L, Nakamura E, Roca C, Ding J, et al. Independence of familial transmission of mania and depression: results of the NIMH family study of affective spectrum disorders. *Mol Psychiatry* 2014;19(2):214–9. Feb 1.
- [39] Mataix-Cols D, Billotti D, Fernández De La Cruz L, Nordstletten AE. The London field trial for hoarding disorder. *Psychol Med* 2013;43(4):837–47. Apr.
- [40] Dereboy F, Dereboy Ç, Eskin M. Validation of the DSM-5 alternative model personality disorder diagnoses in Turkey, part 1: LEAD validity and reliability of the personality functioning ratings. *J Pers Assess* 2018;100(6):603–11. Nov.
- [41] Sung M, Goh TJ, Tan BLJ, Chan JS, Liew HSA. Comparison of DSM-IV-TR and DSM-5 criteria in diagnosing autism spectrum disorders in Singapore. *J Autism Dev Disord* 2018;48(10):3273–81. Oct.
- [42] Nishiyama Takeshi, Sumi Satoshi, Watanabe Hiroto, Suzuki Futoshi, Kuru Yukiko, Shiino Tomoko, et al. The kiddie schedule for affective disorders and schizophrenia present and lifetime version (K-SADS-PL) for DSM-5: a validation for neurodevelopmental disorders in Japanese outpatients. *Compr Psychiatry* 2020:96. Jan 1.
- [43] Gerdner A, Kestenberg J, Mattias E. Validity of the Swedish SCID and ADDIS diagnostic interviews for substance use disorders: sensitivity and specificity compared with a LEAD golden standard. *Nord J Psychiatry* 2015;69(1):48–56. Jan 1.
- [44] North CS, Simic Z, Burruss J. Design, implementation, and assessment of a public comprehensive specialty care program for early psychosis. *J Psychiatr Pract* 2019;25(2):91–102. Mar.
- [45] Blackmore R, Gray KM, Melvin GA, Newman L, Boyle JA, Gibson-Helm M. Identifying post-traumatic stress disorder in women of refugee background at a public antenatal clinic. *Arch Womens Ment Health* 2022;25(1):191–8. Feb.
- [46] Osório FL, Loureiro SR, Hallak JEC, Machado-de-Sousa JP, Ushirohira JM, Baes CVW, et al. Clinical validity and intrarater and test-retest reliability of the structured clinical interview for DSM-5 - clinician version (SCID-5-CV). *Psychiatry Clin Neurosci* 2019;73(12):754–60. Dec.
- [47] Flake JK, Fried EI. Measurement schmeasurement: questionable measurement practices and how to avoid them. *Adv Methods Pract Psychol Sci* 2020;3(4):456–65. Dec 1.
- [48] Aguinis H, Ramani RS, Alabduljader N. What you see is what you get? Enhancing methodological transparency in management research. *Acad Manag Ann* 2018;12(1):83–110. Jan.
- [49] von Elm E, Altman DG, Egger M, Pocock SJ, Gøtzsche PC, Vandenbroucke JP, et al. The Strengthening of reporting of observational studies in epidemiology (STROBE) statement: guidelines for reporting observational studies. *Ann Intern Med* 2007;147(8):573–7. Oct 16.
- [50] STARD Group, Bossuyt PM, Reitsma JB, Bruns DE, Gatsonis CA, Glasziou PP, et al. STARD 2015: an updated list of essential items for reporting diagnostic accuracy studies. *BMJ* 2015:351. Oct 28.
- [51] Hopewell S, Chan AW, Collins GS, Hróbjartsson A, Moher D, Schulz KF, et al. CONSORT 2025 statement: updated guideline for reporting randomised trials. *BMJ* 2025;389. Apr 14. (e081123).
- [52] Collins GS, Moons KGM, Dhiman P, Riley RD, Beam AL, Calster BV, et al. TRIPOD+AI statement: updated guidance for reporting clinical prediction models that use regression or machine learning methods. *BMJ* 2024;385. Apr 16. (e078378).
- [53] Moher D, Schulz KF, Simeria I, Altman DG. Guidance for developers of health research reporting guidelines. *PLoS Med* 2010;7(2). Feb 16. (e1000217).
- [54] 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. *Int J Surg* 2014;12(12):1495–9. Dec 1.
- [55] Schulz KF, Altman DG, Moher D, CONSORT Group. CONSORT 2010 statement: updated guidelines for reporting parallel group randomised trials. *BMJ* 2010;340. Mar 23. (c332).
- [56] Husereau D, Drummond M, Augustovski F, de Bekker-Grob E, Briggs AH, Carswell C, et al. Consolidated health economic evaluation reporting standards 2022 (CHEERS 2022) statement: updated reporting guidance for health economic evaluations. *BMC Med* 2022;20(1):23. Jan 12.
- [57] Hsu CC, Sanford BA. The Delphi technique: making sense of consensus. *Pract Assess Res Eval* 2007;12(10):1–8. Jan 1.
- [58] Vernooij Robin WM, Alonso-Coeillo Pablo, Brouwers Melissa, García Laura Martínez, CheckUp Panel. Reporting items for updated clinical guidelines: checklist for the reporting of updated guidelines (CheckUp). *PLoS Med* 2017;14(1). Jan 1. (e1002207–e1002207).
- [59] Morrison EH, Sorkin D, Mosqueda L, Ayutyanont N. Validity and reliability of the scale to report emotional stress signs-multiple sclerosis (STRESS-MS) in assessing abuse and neglect of adults with multiple sclerosis. *Int J MS Care* 2022;1.
- [60] Mackenhauer J, Winslov JH, Holmskov J, Brodsgaard I, Larsen TG, Mainz J. Analysis of suicides reported as adverse events in psychiatry resulted in nine quality improvement initiatives. *Crisis* 2022;43(4):307–14. Jul.
- [61] Hendriks E, Muris P, Meesters C, Houben K. Childhood disorder: dysregulated self-conscious emotions? Psychopathological correlates of implicit and explicit shame and guilt in clinical and non-clinical children and adolescents. *Front Psychol* 2022;13:822725.
- [62] Paap MCS, Heltne A, Pedersen G, Germans Selvik S, Frans N, Wilberg T, et al. More is more: evidence for the incremental value of the SCID-II/SCID-5-PD specific factors over and above a general personality disorder factor. *Personal Disord* 2022;13(2):108–18. Mar.
- [63] Aydin S, Siebelink BM, Crone MR, van Ginkel JR, Numans ME, Vermeiren RRJM, et al. The diagnostic process from primary care to child and adolescent mental healthcare services: the incremental value of information conveyed through referral letters, screening questionnaires and structured multi-informant assessment. *BJP Psychol Open* 2022;8(3). Apr 7. (e81).
- [64] Rocha Neto HG, Lessa JLM, Koiller LM, Pereira AM, de Souza Gomes BM, Veloso Filho CL, et al. Non-standard diagnostic assessment reliability in psychiatry: a study in a Brazilian outpatient setting using Kappa. *Eur Arch Psychiatry Clin Neurosci* 2024 Oct;274(7):1759–70. Epub 2023 Dec 12.
- [65] Spangenberg H, Ramklint M, Ramirez A. A long-term follow-up study of labor market marginalization in psychiatric patients with and without personality disorder. *Ups J Med Sci* 2023;128. <https://doi.org/10.48101/ujms.v128.9014>. Jul 31.
- [66] Muris P, Büttgens L, Koolen M, Manniën C, Scholtes N, van Dooren-Theunissen W. Symptoms of selective mutism in middle childhood: psychopathological and temperament correlates in Non-clinical and clinically referred 6- to 12-year-old children. *Child Psychiatry Hum Dev* 2024 Dec;55(6):1514–25.
- [67] Sven CA, Pedersen G, Hummelen B, Kvarstein EH. Personality disorders: The impact of severity on societal costs. *Eur Arch Psychiatry Clin Neurosci* 2025;275(1):181–92. Epub 2023 Nov 22.
- [68] Pedersen G, Kvarstein EH, Wilberg T, Folmo EJ, Burlingame GM, Lorentzen S. The group questionnaire (GQ)—psychometric properties among outpatients with personality disorders. *Group Dyn Theory Res Pract* 2023;27(2):81–98.
- [69] Sadleir PHM, Clarke RC, Goddard CE, Mickle P, Platt PR. Agreement of a clinical scoring system with allergic anaphylaxis in suspected perioperative hypersensitivity reactions: prospective validation of a new tool. *Br J Anaesth* 2022;129(5):670–8. Nov.
- [70] Khan AM, Ahmed S, Chowdhury NH, Islam MS, McCollum ED, King C, et al. Developing a video expert panel as a reference standard to evaluate respiratory rate counting in paediatric pneumonia diagnosis: protocol for a cross-sectional study. *BMJ Open* 2022;12(11). Nov 15. (e067389).

- [71] Loots FJ, Smits M, Hopstaken RM, Jenniskens K, Schroeten FH, van den Bruel A, et al. New clinical prediction model for early recognition of sepsis in adult primary care patients: a prospective diagnostic cohort study of development and external validation. *Br J Gen Pract J R Coll Gen Pract* 2022;72(719):e437–45. Jun.
- [72] Leroux A, Frey KP, Crainiceanu CM, Obremskey WT, Stinner DJ, Bosse MJ, et al. Defining incidence of acute compartment syndrome in the research setting: a proposed method from the PACS study. *J Orthop Trauma* 2022;36(Suppl. 1): S26–32. Jan 1.
- [73] Kocks JWH, Cao H, Holzhauer B, Kaplan A, FitzGerald JM, Kostikas K, et al. Diagnostic performance of a machine learning algorithm (asthma/chronic obstructive pulmonary disease [COPD] differentiation classification) tool versus primary care physicians and pulmonologists in asthma, COPD, and asthma/COPD overlap. *J Allergy Clin Immunol Pract* 2023;11(5):1463–1474.e3. May 1.
- [74] Lacroix L, Papis S, Mardegan C, Luterbacher F, L'Huillier A, Sahyoun C, et al. Host biomarkers and combinatorial scores for the detection of serious and invasive bacterial infection in pediatric patients with fever without source. *PLoS One* 2023; 18(11):e0294032.
- [75] Nienhuis PH, van Nieuwland M, van Praagh GD, Markusiewicz K, Colin EM, van der Geest KSM, et al. Comparing diagnostic performance of short and long [18F] FDG-PET acquisition times in giant cell arteritis. *Diagn Basel Switz* 2023;14(1):62. Dec 27.
- [76] Himmelreich JCL, Harskamp RE. Diagnostic accuracy of the PMcardio smartphone application for artificial intelligence–based interpretation of electrocardiograms in primary care (AMSTELHEART-1). *Cardiovasc Digit Health J* 2023;4(3):80–90. Jun 1.
- [77] de la Matta M, Alonso-González M, García-Santigosa M, Arance-García M, Sánchez-Peña J, Castro-Liñán LM, et al. Accuracy and comprehensiveness in recording information of a web-based application for preoperative assessment: a prospective observational study. *J Perianesth Nurs* 2023;38(3):440–7. Jun 1.
- [78] Peterson BS, Kaur T, Baez MA, Whiteman RC, Sawardekar S, Sanchez-Peña J, et al. Morphological biomarkers in the amygdala and Hippocampus of children and adults at high familial risk for depression. *Diagnostics* 2022;12(5):1218. May.
- [79] Reiersen AM, Noel JS, Doty T, Sinkre RA, Narayanan A, Hershey T. Psychiatric diagnoses and medications in Wolfram syndrome. *Scand J Child Adolesc Psychiatry Psychol* 2022;10(1):163–74. Jan.
- [80] Bradshaw J, Shi D, Hendrix CL, Saulnier C, Klaiman C. Neonatal neurobehavior in infants with autism spectrum disorder. *Dev Med Child Neurol* 2022;64(5):600–7. May.
- [81] Hesam-Shariati S, Overs BJ, Roberts G, Toma C, Watkeys OJ, Green MJ, et al. Epigenetic signatures relating to disease-associated genotypic burden in familial risk of bipolar disorder. *Transl Psychiatry* 2022;12(1):310. Aug 3.
- [82] Shima C, Lee R, Coccaro EF. Associations of aggression and use of caffeine, alcohol and nicotine in healthy and aggressive individuals. *J Psychiatr Res* 2022;146:21–7. Feb.
- [83] Kvig EI, Nilssen S. Does method matter? Assessing the validity and clinical utility of structured diagnostic interviews among a clinical sample of first-admitted patients with psychosis: a replication study. *Front Psych* 2023;14:1076299.
- [84] Detera-Wadleigh SD, Kassem L, Besancon E, Lopes F, Akula N, Sung H, et al. A resource of induced pluripotent stem cell (iPSC) lines including clinical, genomic, and cellular data from genetically isolated families with mood and psychotic disorders. *Transl Psychiatry* 2023;13(1):397. Dec 16.
- [85] Hill SY, Hostyk J. A whole exome sequencing study to identify rare variants in multiplex families with alcohol use disorder. *Front Psychiatry*. 2023;14:1216493. Oct 17.
- [86] Jones W, Klaiman C, Richardson S, Lambha M, Reid M, Hamner T, et al. Development and replication of objective measurements of social visual engagement to aid in early diagnosis and assessment of autism. *JAMA Netw Open* 2023;6(9). Sep 5. (e2330145).
- [87] Dimian AF, Estes AM, Dager S, Piven J, Wolff JJ. Network for the I. Predicting self-injurious behavior at age three among infant siblings of children with autism. *Autism Res* 2023;16(9):1670–80.
- [88] Korevaar DA, Cohen JF, Reitsma JB, Bruns DE, Gatsonis CA, Glasziou PP, et al. Updating standards for reporting diagnostic accuracy: the development of STARD 2015. *Res Integr Peer Rev* 2016;1(1):7. Jun 7.
- [89] Collins GS, Reitsma JB, Altman DG, Moons KGM. Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (TRIPOD): the TRIPOD statement. *BMJ* 2015;350. Jan 7. (g7594).
- [90] Blanco D, Altman D, Moher D, Boutron I, Kirkham JJ, Cobo E. Scoping review on interventions to improve adherence to reporting guidelines in health research. *BMJ Open* 2019;9(5). May 9. (e026589).
- [91] Samaan Z, Mbuagbaw L, Kosa D, Debono VB, Dillenburg R, Zhang S, et al. A systematic scoping review of adherence to reporting guidelines in health care literature. *J Multidiscip Healthc* 2013;6(6):169–88. May.
- [92] Heus P, Damen JAAG, Pajouheshnia R, Scholten RJPM, Reitsma JB, Collins GS, et al. Uniformity in measuring adherence to reporting guidelines: the example of TRIPOD for assessing completeness of reporting of prediction model studies. *BMJ Open* 2019;9(4). Apr 1. (e025611).
- [93] Li Z, Luo X, Yang Z, Zhang H, Wang B, Ge L, et al. RAPID: Reliable and efficient Automatic generation of submission rePorting checklists with Large language models. *bioRxiv* 2025. p. 2025.02.13.638015.
- [94] Wrightson JG, Blazey P, Moher D, Khan KM, Ardern CL. GPT for RCTs?: Using AI to determine adherence to reporting guidelines. *medRxiv* 2024. p. 2023.12.14.23299971.