

1 **A Quality Assessment Tool for Artificial Intelligence Centred Diagnostic Test Accuracy Studies: QUADAS-AI**

2
3 Viknesh Sounderajah^{1,2}, Hutan Ashrafian^{1,2}, Sherri Rose⁷, Nigam H. Shah³⁴, Marzyeh Ghassemi³, Robert Golub⁶,
4 Charles E. Kahn, Jr.¹⁹, Andre Esteva¹⁷, Alan Karthikesalingam⁸, Bilal Mateen¹⁸, Dale Webster⁸, Dan Milea²⁰,
5 Daniel Ting²⁰, Darren Treanor^{21,22,23,24}, Dominic Cushman²⁵, Dominic King^{1,26}, Duncan McPherson²⁷, Ben
6 Glocker³⁶, Felix Greaves²⁸, Leanne Harling^{1,2}, Johan Ordish²⁷, Jérémie F. Cohen²⁹, Jon Deeks⁵, Mariska
7 Leeflang¹³, Matthew Diamond³⁰, Matthew D.F. McInnes³¹, Melissa McCradden³², Michael D. Abràmoff³³, Pasha
8 Normahani², Sheraz R. Markar², Stephanie Chang³⁵, Xiaoxuan Liu^{14,15,16}, Susan Mallett⁴, Shravya Shetty⁸,
9 Alastair Denniston^{14,15,16}, Gary S. Collins^{9,10}, David Moher¹¹, Penny Whiting¹², Patrick M. Bossuyt^{*13} and Ara
10 Darzi^{*1,2}

11
12 * Joint senior authorship

13
14 **Author Affiliations:**

15 ¹ Institute of Global Health Innovation, Imperial College London, United Kingdom

16 ² Department of Surgery and Cancer, Imperial College London, United Kingdom

17 ³ Institute for Medical Engineering & Science, Massachusetts Institute of Technology, United States of America

18 ⁴ Centre for Medical Imaging, University College London, United Kingdom

19 ⁵ Institute of Applied Health Research, University of Birmingham, United Kingdom

20 ⁶ JAMA (Journal of the American Medical Association), United States of America

21 ⁷ Center for Health Policy and Center for Primary Care and Outcomes Research, Stanford University, California,
22 United States of America

23 ⁸ Google Health

24 ⁹ Centre for Statistics in Medicine, Nuffield Department of Orthopaedics, Rheumatology and Musculoskeletal
25 Sciences, University of Oxford, Oxford, United Kingdom

26 ¹⁰ NIHR Oxford Biomedical Research Centre, Oxford University Hospitals NHS Foundation Trust, Oxford, United
27 Kingdom

28 ¹¹ Ottawa Hospital Research Institute, Canada

29 ¹² Bristol Medical School, University of Bristol, Bristol, United Kingdom

30 ¹³ Department of Epidemiology and Data Science, Amsterdam University Medical Centres, University of
31 Amsterdam, The Netherlands

32 ¹⁴ Institute of Inflammation and Ageing, College of Medical and Dental Sciences, University of Birmingham,
33 United Kingdom

34 ¹⁵ University Hospitals Birmingham NHS Foundation Trust, Birmingham, United Kingdom

35 ¹⁶ Health Data Research UK, London, United Kingdom

36 ¹⁷ Salesforce Research, San Francisco, United States of America

37 ¹⁸ Wellcome Trust, London, United Kingdom

38 ¹⁹ University of Pennsylvania, Philadelphia, Pennsylvania, United States of America

39 ²⁰ Singapore Eye Research Institute, Singapore National Eye Centre, Singapore

40 ²¹ Leeds Teaching Hospitals NHS Trust, Leeds, United Kingdom

41 ²² University of Leeds, Leeds, UK

42 ²³ Department of Clinical Pathology, and Department of Clinical and Experimental Medicine, Linköping
43 University, Linköping, Sweden

44 ²⁴ Center for Medical Image Science and Visualization (CMIV), Linköping University, Linköping, Sweden

45 ²⁵ NHSX, London, United Kingdom

46 ²⁶ Optum, Paddington, London, United Kingdom

47 ²⁷ Medicines and Healthcare products Regulatory Agency, London, United Kingdom

48 ²⁸ National Institute for Health and Care Excellence, London, United Kingdom

49 ²⁹ Department of Pediatrics and Inserm UMR 1153 (Centre of Research in Epidemiology and Statistics), Necker
 50 - Enfants Malades Hospital, Assistance Publique - Hôpitaux de Paris, Université de Paris, Paris, France
 51 ³⁰ Food and Drug Administration, Maryland, United States of America
 52 ³¹ Departments of Radiology and Epidemiology, University of Ottawa, The Ottawa Hospital Research Institute,
 53 Canada
 54 ³² Department of Bioethics, The Hospital for Sick Kids, Toronto, Canada
 55 ³³ Department of Ophthalmology and Visual Sciences, University of Iowa, Iowa, United States of America
 56 ³⁴ Center for Biomedical Informatics Research, Stanford University, California, United States of America
 57 ³⁵ Annals of Internal Medicine, American College of Physicians, United States of America
 58 ³⁶ Faculty of Engineering, Department of Computing, Imperial College London, United Kingdom

59
 60 **Corresponding authors**

61 Ara Darzi.
 62 Institute of Global Health Innovation & Department of Surgery and Cancer, Imperial College London, United
 63 Kingdom. a.darzi@imperial.ac.uk

64
 65 Patrick M. Bossuyt
 66 Department of Epidemiology and Data Science, Amsterdam University Medical Centres, University of
 67 Amsterdam, The Netherlands. p.m.bossuyt@amsterdamumc.nl

68
 69 **To the Editor** - Over the next decade, it is projected that artificial intelligence (AI), particularly machine
 70 learning, centred systems will become key components of several workflows within the health sector. Medical
 71 diagnosis is particularly seen as one of the first areas that would harbour the adoption of AI innovations.
 72 Indeed, over 90% of health-related AI systems that have reached regulatory approval by the U.S. Food and
 73 Drug Administration belong to the field of diagnostics ¹.

74
 75 In the current paradigm, the majority of diagnostic investigations require interpretation from a clinician to
 76 identify the presence of a target condition; a crucial step in determining subsequent treatment strategies.
 77 Despite being an essential step in the provision of patient care, many health systems find it increasingly
 78 difficult to meet the demand for diagnostic test interpretation. To address this issue, diagnostic AI systems
 79 have been characterised as medical devices which may alleviate the burden placed upon diagnosticians:
 80 serving as case triage tools, enhancing diagnostic accuracy, and stepping in as a second reader when
 81 necessary. As AI centred diagnostic test accuracy (AI DTA) studies emerge, there has been a concurrent rise in
 82 systematic reviews which amalgamate the findings of comparable studies.

83 Strikingly, of these published AI DTA systematic reviews, 94% have been conducted in the absence of an AI
 84 specific quality assessment tool ². The most commonly used instrument is the QUADAS-2 (Quality Assessment
 85 of Diagnostic Accuracy Studies) tool ³. QUADAS-2 is a risk of bias and applicability tool whose use is encouraged
 86 by PRISMA 2020 guidance ⁴. QUADAS-2 does not, however, accommodate for niche terminology encountered
 87 in AI DTA studies nor does it signal researchers to the sources of bias found within this class of studies.
 88 Examples of such biases, when framed against the established domains of QUADAS-2 (Patient Selection; Index
 89 Test; Reference Standard; and Flow and Timing) are listed in table 1.

Domain	Description	Biases
Patient Selection	A description of included patients detailing prior	In AI DTA studies, eligible patients are often excluded on account of competing input data entry requirements (e.g. image quality) which, themselves, are variably reported. As highlighted by the CONSORT-AI

	<p>testing, presentation, setting and the intended use of the index test.</p>	<p>guidelines⁵, there is a need to accurately characterise the source, size and quality of input data alongside clear patient eligibility criteria.</p> <p>Data source issues can negatively impact the performance and overall applicability of an index test. For example, to minimise research costs, there has been increasing usage of datasets sourced from open-source repositories. Whilst this offers a pragmatic option, many open-source datasets have been found to house the inadvertent duplication of data across repositories, erroneous labelling, and incomplete patient demographic data.</p> <p>Manuscripts reporting both the development and validation of an index test rarely present the rationale and breakdown of its training, validation, and test sets. Small datasets, particularly those that lack complexity and balance, can result in overfitting, whereby the final index test resembles the training data too closely and is unable to reliably fit additional data. The clinical manifestation of this issue is the inability to accurately diagnose instances of a pathology if its clinical presentation does not closely resemble the training cases that the index test had previously encountered.</p> <p>There are various points within the data curation pipeline where quality may be compromised. For example, image pre-processing, a practice whereby image formats and resolutions are homogenised for the purpose of training, is an essential step in AI workflows. However, either down- or up-scaling resolution may impact the ability of certain index tests to identify diagnostic features effectively. Moreover, the lack of image metadata can also preclude the ability to explore an index test's dependence on specific data acquisition parameters, for example, the model of scanner used to acquire imaging data.</p>
<p>Index Test</p>	<p>The diagnostic test being evaluated and how it has been conducted and interpreted within the context of the study.</p>	<p>Only a limited number of published studies have undertaken adequate external evaluation when presenting the development and evaluation of their diagnostic tests. Reliance upon data from the same dataset that is used to train the diagnostic test (internal holdout set) can overestimate diagnostic performance.</p>
<p>Reference Standard</p>	<p>The choice of reference standard and how it has been conducted and interpreted within the context of the study.</p>	<p>There are multiple instances, as highlighted by Harris et al.⁶, in which studies have reported the development of index tests against inappropriate reference standards, as opposed to more appropriate tests that provide higher sensitivity and specificity. For example, a clinician using a chest X-ray to diagnose pulmonary tuberculosis rather than the more accurate use of sputum culture. Studies with inappropriate reference standards are poorly reflective of real-world clinical practice in which reference standards consist of the amalgamation of clinical, radiological and laboratory data.</p>

Flow and Timing	The time interval and the use of any interventions between the application of the index test and the reference standard	The timing between index test and reference standard is often poorly reported. As highlighted in a recent systematic review ⁷ , studies which reported the performance of index tests to diagnose SARS-CoV-2 from chest X-rays did not routinely note the timing of the confirmatory RT-PCR test in relation to the imaging data. It is well understood that RT-PCR is a time sensitive assay and failing to report this relationship significantly hinders the overall clinical validity of the study results.
-----------------	---	--

90

91 **Table 1: Examples of bias within AI DTA studies**

92 In order to tackle the sources of bias described above, as well as AI specific examples such as algorithmic bias,
 93 we propose an AI-specific extension to QUADAS-2 and QUADAS-C ⁸, a risk of bias tool developed for
 94 comparative accuracy studies. This new tool, QUADAS-AI, will provide researchers and policy makers with a
 95 specific framework to evaluate the risk of bias and applicability when conducting reviews evaluating AI DTA
 96 and reviews of comparative accuracy studies evaluating at least one AI centred index test.

97 QUADAS-AI will be complementary to ongoing reporting guideline tool initiatives, such as STARD-AI ⁹ and
 98 TRIPOD-AI ¹⁰. QUADAS-AI is being coordinated by a global Project Team and Steering Committee consisting of
 99 clinician scientists, computer scientists, epidemiologists, statisticians, journal editors, EQUATOR Network
 100 representatives, regulatory leaders, industry leaders, funders, health policy makers and bioethicists. Given the
 101 reach of AI technologies, we view that connecting global stakeholders is of the utmost importance for this
 102 initiative. In turn, we would welcome contact from any new potential collaborators.

103 **References**

104

- 105 1. Benjamins, S., Dhunoo, P. & Meskó, B. *npj Digit. Med.* **3**, 118 (2020).
 106 2. Jayakumar, S. *et al.* *What are the Quality Assessment Standards used in Artificial Intelligence Diagnostic*
 107 *Accuracy Systematic Reviews?* <https://www.researchsquare.com> (2021) doi:10.21203/RS.3.RS-
 108 329433/V1.
 109 3. Whiting, P. F. QUADAS-2: *Ann. Intern. Med.* **155**, 529 (2011).
 110 4. Page, M. J. *et al.* The PRISMA 2020 statement: *The BMJ* vol. 372 (2021).
 111 5. Liu, X. & Rivera, S. C. *Nat. Med.* **2020 269 26**, 1364–1374 (2020).
 112 6. Harris, M. *et al.* *PLoS One* **14**, (2019).
 113 7. Roberts, M. *et al.* *Nat. Mach. Intell.* **3**, 199–217 (2021).
 114 8. Yang, B. *et al.* (2018) doi:10.17605/OSF.IO/HQ8MF.
 115 9. Sounderajah, V. *et al.* *Nature Medicine* vol. 26 807–808 (2020).
 116 10. Collins, G. & Moons, K. *Lancet* **393**, 1577–1579 (2019).

117

118 Acknowledgments

119 Infrastructure support for this research was provided by the NIHR Imperial Biomedical Research Centre (BRC).
 120 GSC is supported by the NIHR Biomedical Research Centre, Oxford, and Cancer Research UK (programme
 121 grant: C49297/A27294). DT is funded by National Pathology Imaging Co-operative, NPIC (Project no. 104687),
 122 supported by a £50m investment from the Data to Early Diagnosis and Precision Medicine strand of the
 123 government’s Industrial Strategy Challenge Fund, managed and delivered by UK Research and Innovation
 124 (UKRI). FG is supported by the National Institute for Health Research Applied Research Collaboration
 125 Northwest London. The views and opinions expressed herein are those of the authors and do not necessarily
 126 reflect the views of their employers or funders.

127

128 Author contributions

129 VS, SR, NS, MG, RG, CEK, XL, GSC, DW, AE, HA, DM (Dan Milea), DM (Duncan McPherson), JO, DT, JFC, ML, MM,
130 MDFM, MDA, SM, PW and PMB prepared the first draft of the manuscript. Critical edits and feedback have
131 been attained from all co-authors. The study described in the manuscript has been conceptualised, discussed,
132 and agreed upon between all co-authors.

133

134 Competing interests

135 AK, SS and DW are employees at Google. AD and HA are employees at Flagship Pioneering UK Ltd. AE is an
136 employee at Salesforce. DK is an employee at Optum. None of the other authors have any competing interests.