### UNIVERSITY<sup>OF</sup> BIRMINGHAM

## University of Birmingham Research at Birmingham

# Tackling bias in AI datasets through the STANDING together initiative

Ganapathi, Shaswath; Palmer, Jo; Alderman, Joseph; Calvert, Melanie; Espinoza, Cyrus; Gath, Jacqui; Ghassemi, Marzyeh; Heller, Katherine; Mckay, Francis; Karthikesalingam, Alan; Kuku, Stephanie; Mackintosh, Maxine; Manohar, Sinduja; Mateen, Bilal A; Matin, Rubeta N.; McCradden, Melissa M; Oakden-Rayner, Lauren; Ordish, Johan; Pearson, Russell; Pfohl, Stephen R

DOI.

10.1038/s41591-022-01987-w

License:

Other (please specify with Rights Statement)

Document Version
Peer reviewed version

Citation for published version (Harvard):

Ganapathi, S, Palmer, J, Alderman, J, Calvert, M, Espinoza, C, Gath, J, Ghassemi, M, Heller, K, Mckay, F, Karthikesalingam, A, Kuku, S, Mackintosh, M, Manohar, S, Mateen, BA, Matin, RN, McCradden, MM, Oakden-Rayner, L, Ordish, J, Pearson, R, Pfohl, SR, Rostamzadeh, N, Sapey, E, Sebire, NJ, Sounderajah, V, Summers, C, Treanor, D, Denniston, A & Liu, X 2022, 'Tackling bias in Al datasets through the STANDING together initiative', *Nature Medicine*. https://doi.org/10.1038/s41591-022-01987-w

Link to publication on Research at Birmingham portal

**Publisher Rights Statement:** 

This version of the article has been accepted for publication, after peer review (when applicable) and is subject to Springer Nature's AM terms of use, but is not the Version of Record and does not reflect post-acceptance improvements, or any corrections. The Version of Record is available online at: https://doi.org/10.1038/s41591-022-01987-w

General rights

Unless a licence is specified above, all rights (including copyright and moral rights) in this document are retained by the authors and/or the copyright holders. The express permission of the copyright holder must be obtained for any use of this material other than for purposes permitted by law.

•Users may freely distribute the URL that is used to identify this publication.

- •Users may download and/or print one copy of the publication from the University of Birmingham research portal for the purpose of private study or non-commercial research.
- •User may use extracts from the document in line with the concept of 'fair dealing' under the Copyright, Designs and Patents Act 1988 (?)
  •Users may not further distribute the material nor use it for the purposes of commercial gain.

Where a licence is displayed above, please note the terms and conditions of the licence govern your use of this document.

When citing, please reference the published version.

Take down policy

While the University of Birmingham exercises care and attention in making items available there are rare occasions when an item has been uploaded in error or has been deemed to be commercially or otherwise sensitive.

If you believe that this is the case for this document, please contact UBIRA@lists.bham.ac.uk providing details and we will remove access to the work immediately and investigate.

Download date: 19. Apr. 2024

#### Tackling bias in AI datasets through the STANDING together initiative

Shaswath Ganapathi<sup>1</sup>, Jo Palmer<sup>2</sup>, Joseph E Alderman<sup>1,2</sup>, Melanie Calvert<sup>3,4,5,6,7</sup>, Cyrus Espinoza<sup>8</sup>, Jacqui Gath<sup>8</sup>, Marzyeh Ghassemi<sup>9</sup>, Katherine Heller<sup>10</sup>, Francis Mckay<sup>11</sup>, Alan Karthikesalingam<sup>12</sup>, Stephanie Kuku<sup>13</sup>, Maxine Mackintosh<sup>14</sup>, Sinduja Manohar<sup>15</sup>, Bilal A Mateen<sup>16, 17</sup>, Rubeta Matin<sup>18</sup>, Melissa McCradden<sup>19,20</sup>, Lauren Oakden-Rayner<sup>21</sup>, Johan Ordish<sup>22</sup>, Russell Pearson<sup>22</sup>, Stephen R Pfohl<sup>10</sup>, Negar Rostamzadeh<sup>23</sup>, Elizabeth Sapey<sup>1</sup>, Neil Sebire<sup>15,24</sup>, Viknesh Sounderajah<sup>25, 26</sup>, Charlotte Summers<sup>27</sup>, Darren Treanor<sup>28, 29, 30, 31</sup>, Alastair K Denniston\*<sup>1,2,3,5,15</sup>, Xiaoxuan Liu\*<sup>1,2,3</sup>

- Institute of Inflammation and Ageing, College of Medical and Dental Sciences, University of Birmingham, UK
- 2. University Hospitals Birmingham NHS Foundation Trust, UK
- 3. Birmingham Health Partners Centre for Regulatory Science and Innovation, University of Birmingham, UK
- 4. Centre for Patient Reported Outcome Research, Institute of Applied Health Research, University of Birmingham, UK
- 5. NIHR Birmingham Biomedical Research Centre, University of Birmingham, UK
- 6. NIHR Surgical Reconstruction and Microbiology Research Centre, University of Birmingham, Birmingham, UK
- 7. NIHR Applied Research Collaborative West Midlands University of Birmingham, Birmingham, UK.
- 8. Patient Partner, UK
- 9. Department of Electrical Engineering and Computer Science; Institute for Medical Engineering and Science, Massachusetts Institute of Technology, USA.
- 10. Google Research, USA
- 11. The Ethox Centre and the Wellcome Centre for Ethics and Humanities, Nuffield Department of Population Health, University of Oxford, UK
- 12. Google Research, UK
- 13. Institute of Women's Health, University College London, UK
- 14. Genomics England, UK
- 15. Health Data Research, UK
- 16. Institute of Health Informatics, University College London, UK
- 17. The Wellcome Trust, UK
- 18. Oxford University Hospitals NHS Foundation Trust, UK
- 19. Department of Bioethics, Hospital for Sick Children, Canada
- 20. Dalla Lana School of Public Health, University of Toronto, Canada
- 21. Australian Institute for Machine Learning, University of Adelaide, Australia
- 22. Medicines and Healthcare Products Regulatory Agency, UK
- 23. Google Research, Canada
- 24. Great Ormond Street Hospital for Children, UK
- 25. Institute of Global Health Innovation, Imperial College London, UK.
- 26. Department of Surgery and Cancer, Imperial College London, UK.
- 27. Wolfson Lung Injury Unit, Heart and Lung Research Institute, University of Cambridge, UK
- 28. Leeds Teaching Hospitals NHS Trust, UK
- 29. University of Leeds, UK
- 30. Department of Clinical Pathology, and Department of Clinical and Experimental Medicine, Linköping University, Sweden

31. Center for Medical Image Science and Visualization (CMIV), Linköping University, Sweden

\*Joint senior author

**Corresponding author:** Dr Xiaoxuan Liu, University Hospitals Birmingham NHS Foundation Trust, Birmingham, UK. x.liu.8@bham.ac.uk

**To the Editor** - As of June 2022, a wide range of artificial intelligence (AI) as a Medical Device (AlaMDs) have received regulatory clearance internationally, with at least 343 devices cleared by the United Stated (US) Food and Drug Administration (FDA). Despite the enormous potential of AlaMDs, their rapid growth in healthcare has been accompanied by concerns that AI models may learn biases engrained in medical practice and exacerbate health inequalities. This has been exemplified through a number of AI systems which have shown the ability of algorithms to systematically misrepresent and exacerbate health problems in minoritised groups. This raises concerns that, without appropriate safeguarding, AI models may perpetuate existing health inequality and mistrust.

Tackling bias in AI requires a multifaceted approach. A recent report by the US National Institute of Standards and Technology on bias in AI emphasised that algorithmic development does not occur through engineering decisions alone, but embeds a myriad of values and behaviours within the data and the humans who interact with them. The report calls for a sociotechnical approach that considers how different biases interact and the social contexts within which AI systems are built and used.<sup>4</sup> Although there is an expanding field of research dedicated to fairness in machine learning, many AlaMD receiving regulatory clearance have not appropriately accounted for biases that disadvantage certain populations. There are also ethical challenges around algorithmic fairness methods (computational techniques seeking to ensure outputs are not unjustifiably influenced by bias), given that these methods are aimed at making *predictions* fair, rather than enabling fair treatment of individuals.<sup>5</sup> Furthermore, current approaches to satisfying regulatory requirements are focused on aggregate-level performance, which can mask stratification across subpopulations.

One major source of bias is the data underpinning AI systems. It is often necessary to train models with large quantities of data, which means datasets are often sourced to prioritise sample size. There are concerns that many health datasets do not adequately represent minoritised groups, however the extent of this problem is unknown because many datasets do not provide demographic information, for example on ethnicity and race. Publicly available

datasets for skin cancer and eye imaging have shown inconsistent and incomplete demographic reporting, and are disproportionately collected from a small number of high income countries.<sup>6,7</sup> For skin cancer datasets, reporting of key demographic information, even when clinically relevant (such as ethnicity and skin tone), was only present in 2% of datasets.<sup>7</sup>

Under-representation in datasets can impact the fairness of AI systems through two principal means. During AI development, under-representation within training datasets can negatively impact model performance for under-represented groups.<sup>3</sup> A lack of diversity within the training data risks poor generalisability of model performance post-deployment. During evaluation, under-representation within test datasets increases the uncertainty of performance in that group due to small sample sizes, which reduces the likelihood of detecting underperformance. Therefore, under-representation not only creates models that under-perform within minoritised populations, but also hampers the ability to detect this bias. Furthermore, under-representation in datasets may result in exclusion of populations from the intended use altogether, thereby creating Al systems licensed for only certain groups within society. Even when datasets are inclusive, additional issues can compound bias. Structural inequities can manifest in datasets through the actions of clinical and data curation teams, who are responsible for recording, selecting, labelling and aggregating data, based on assumptions that reflect hegemonic social attitudes. Addressing the consequences of structural biases requires a wider consideration of the dataset: how and why it was created; the setting in which data was collected and by whom; the extent to which the data reflects broader structural biases and axes of injustice; the inclusion/exclusion criteria; and how measurements, observations and labels were constructed. These concerns have motivated calls for better documentation practices and the creation of tools like Datasheets for Datasets and Healthsheets.8,9

The aforementioned problems are becoming increasingly recognised by medical device regulators. In October 2021, The US FDA, Health Canada, and the UK Medicines and Healthcare products Regulatory Agency (MHRA) jointly published 10 guiding principles for Good Machine Learning Practice. This specifically states that data should be representative of the intended population in order to 'manage bias, promote appropriate and generalizable performance across the intended patient population, assess usability and identify circumstances where the model may underperform'. <sup>10</sup> Commitment to identify and mitigate bias by medical regulators is a significant step in the right direction, however, to date, there is a lack of evidence that these principles are adopted by AI device manufacturers. Without specific consensus on

how to assess the appropriateness of datasets, it is unclear what constitutes best practice regarding the use of health data in Al to promote fairness and equity.

To tackle this problem, we are announcing an initiative to develop Standards for Data Diversity, Inclusivity and Generalisability (STANDING Together). STANDING Together is an international, consensus-based initiative that aims to develop recommendations for the composition (who is represented) and reporting (how they are represented) of datasets underpinning medical AI systems. We will engage patients and the public, clinicians, academic researchers across biomedical, computational and social sciences, industry experts, regulators and policy-makers. The standards will represent the culmination of a multiphase evidence generation process, which consists of: dataset mapping reviews to assess limitations in health datasets across different diseases with regard to diversity and inclusivity; interviews with dataset curators to explore the barriers and challenges to ensuring diversity and inclusivity within health datasets; a modified Delphi consensus study to finalise the content that will feature in these recommendations and; an extensive multi-stakeholder piloting phase. The resulting standards will support informed decision-making for those who strive to engineer and implement fair and safe AI systems in healthcare. STANDING Together will be the first in a line of work through which stakeholders can determine what demographic data is collected and how it is represented in datasets. The findings will motivate curators of health datasets to prioritise diversity and inclusiveness as we seek to build and invest in health datasets of the future. We hope this initiative will enable the availability of more inclusive data to promote responsible AI in healthcare, and in the long-term, better health outcomes for all.

We anticipate that the modified Delphi consensus study will begin in late 2022 and the final standards published in 2023. We welcome those with expertise in AI, health data science and health inequalities to participate and encourage expressions of interest through <a href="https://www.datadiversity.org/involvement/participate-in-our-delphi-study">https://www.datadiversity.org/involvement/participate-in-our-delphi-study</a> or by contacting <a href="mailto:contact@datadiversity.org">contact@datadiversity.org</a>.

#### **Acknowledgements:**

This project is funded by The NHS AI Lab at the NHS Transformation Directorate and the Health Foundation and managed by the National Institute for Health and Care Research (AI\_HI200014). The views expressed in this publication are those of the author(s) and not necessarily those of the NHS Transformation Directorate, the Health Foundation or the National Institute for Health and Care Research.

#### **Author contributions**

XL, AKD, JEA, JP: project conception. XL, SG, AKD, JEA, JP: manuscript drafting. RM, CE, JG, MG, MMa, SM, BAM, MMc, LOR, JO, RP, NS, KH, VS, AK, NR, SRP, SK, DT, FmcK, MJC, ES, CS: manuscript review. All authors reviewed and approved the final version of the manuscript. All authors are involved in the wider conduct and direction of the STANDING Together Programme.

#### Conflicts of interest:

SG, BAM, RM, CE, JG, MG, MMa, SM, MMc, LOR, JO, RP, NS, VS: none

KH, AK, NR and SRP are employees of Google. SK is a consultant for Hardian Health.

DT and FMcK are funded by National Pathology Imaging Co-operative (NPIC, Project no. 104687) which is supported by a £50m investment from the Data to Early Diagnosis and Precision Medicine strand of the government's Industrial Strategy Challenge Fund, managed and delivered by UK Research and Innovation (UKRI).

XL, AKD, JEA and JP are funded by NIHR, the NHS Transformation Directorate and the Health Foundation (Al HI200014).

MJC is Director of the Birmingham Health Partners Centre for Regulatory Science and Innovation, Director of the Centre for the Centre for Patient Reported Outcomes Research and is a National Institute for Health and Care Research (NIHR) Senior Investigator. MJC receives funding from the NIHR, UK Research and Innovation (UKRI), NIHR Birmingham Biomedical Research Centre, the NIHR Surgical Reconstruction and Microbiology Research Centre, NIHR ARC West Midlands, UK SPINE, European Regional Development Fund – Demand Hub and Health Data Research UK at the University of Birmingham and University Hospitals Birmingham NHS Foundation Trust, Innovate UK (part of UKRI), Macmillan Cancer Support, UCB Pharma, Janssen, GSK and Gilead. MC has received personal fees from Astellas, Aparito Ltd, CIS Oncology, Takeda, Merck, Daiichi Sankyo, Glaukos, GSK and the Patient-Centered Outcomes Research Institute (PCORI) outside the submitted work. In addition, a family member owns shares in GSK.

ES receives research funding from UKRI [MR/V033654/1 and MR/S002782/1], the British Lung Foundation, and Alpha 1 Foundation and NIHR.

CS receives research funding from the National Institute for Health and Care Research [NIHR133788], UKRI [MR/P502091/1 and MR/X005070/1], the Wellcome Trust, and the NIHR Cambridge Biomedical Research Centre [BRC1215-20014].

#### References

- Center for Devices & Radiological Health. https://www.fda.gov/medical-devices/software-medical-device-samd/artificial-intelligence-and-machine-learning-aiml-enabled-medical-devices (2021).
- 2. Obermeyer, Z., et al. Science 366, 447–453 (2019).
- 3. Seyyed-Kalantari, L., et al. Nat. Med. 27, 2176–2182 (2021).
- 4. Schwartz, R. et al. Towards a standard for identifying and managing bias in artificial intelligence. (2022) doi:10.6028/nist.sp.1270.
- 5. McCradden, M. D., et al. Lancet Digit Health 2, e221–e223 (2020).
- 6. Khan, S. M. et al. Lancet Digit Health 3, e51–e66 (2021).
- 7. Wen, D. et al. Lancet Digit Health 4, e64-e74 (2022).
- 8. Rostamzadeh, N. et al. arXiv [cs.Al] (2022). doi:10.48550/arXiv.2202.13028
- Gebru, T. et al. Datasheets for Datasets. arXiv [cs.DB] (2018).
   doi:10.48550/arXiv.1803.09010
- 10. Medicines and Healthcare products Regulatory Agency. Good Machine Learning Practice for Medical Device Development: Guiding Principles. Gov.uk https://www.gov.uk/government/publications/good-machine-learning-practice-for-medical-device-development-guiding-principles (2021).