



# Como desmostrar que la IA que uso es robusta?

Prof. Karim Lekadir ICREA Research Professor Universitat de Barcelona Artificial Intelligence in Medicine Lab



## EU Projects



| Project         | Dates     | Role        | Field         |
|-----------------|-----------|-------------|---------------|
| euCanSHare      | 2018-2023 | Coordinator | Cardiology    |
| EarlyCause      | 2020-2024 | Coordinator | Mental health |
| EuCanImage      | 2020-2025 | Coordinator | Oncology      |
| RadioVal        | 2022-2026 | Coordinator | Oncology      |
| AIMIX (ERC)     | 2023-2027 | PI          | Obstetrics    |
| DataTools4Heart | 2022-2026 | Coordinator | Cardiology    |
| AFRICAI-RI      | 2025-2029 | Coordinator | Pulmonology   |



Content



#### 1.3 Objectives

- If you plan to use, develop and/or deploy artificial intelligence (AI) based systems and/or techniques you must demonstrate their technical robustness. AI-based systems or techniques should be, or be developed to become:
  - technically robust, accurate and reproducible, and able to deal with and inform about possible failures, inaccuracies and errors, proportionate to the assessed risk they pose
  - socially robust, in that they duly consider the context and environment in which they operate
  - reliable and function as intended, minimizing unintentional and unexpected harm, preventing unacceptable harm and safeguarding the physical and mental integrity of humans
  - able to provide a suitable explanation of their decision-making processes, whenever they can have a significant impact on people's lives.



#### Terminology



- o Technically & societally robust Al
- o <u>Trustworthy Al</u>
- o Ethical Al
- o Responsible Al
- o Human-centred Al





# Part 1 - What is trustworthy AI? Part 2 - How do we achieve it? Part 3 - Trustworthy AI in an EIC





# Part 1 - What is trustworthy Al? Part 2 - How do we achieve it? Part 3 - Trustworthy Al in an ElC



#### Al in Healthcare





### AI 'outperforms' doctors diagnosing breast cancer



③ 2 January 2020







Tumours missed by 6 clinicians Detected by the Al tool



#### Robustness



#### nature

Explore content V About the journal V Publish with us V

nature > articles > article

Article | Published: 01 January 2020

#### International evaluation of an AI system for breast cancer screening

Scott Mayer McKinney <sup>IM</sup>, Marcin Sieniek, Varun Godbole, Jonathan Godwin, Natasha Antropova, Hutan Ashrafian, Trevor Back, Mary Chesus, Greg S. Corrado, Ara Darzi, Mozziyar Etemadi, Florencia Garcia-Vicente, Fiona J. Gilbert, Mark Halling-Brown, Demis Hassabis, Sunny Jansen, Alan Karthikesalingam, Christopher J. Kelly, Dominic King, Joseph R. Ledsam, David Melnick, Hormuz Mostofi, Lily Peng, Joshua Jay Reicher, Bernardino Romera-Paredes, Richard Sidebottom, Mustafa Suleyman, Daniel Tse <sup>IM</sup>, Kenneth C. Young, Jeffrey De Fauw & Shravya Shetty <sup>IM</sup> → Show fewer authors

<u>Nature</u> 577, 89–94 (2020) <u>Cite this article</u>

Tumours detected by 6 clinicians Missed by the AI tool



#### Transparency



The Machine Making sense of Al Google's breast cancer-predicting AI research is useless without transparency, critics say



#### nature

Explore content Y About the journal Y Publish with us Y

<u>nature</u> > <u>matters arising</u> > article

Matters Arising Published: 14 October 2020

## Transparency and reproducibility in artificial intelligence

Benjamin Haibe-Kains <sup>[C]</sup>, <u>George Alexandru Adam</u>, <u>Ahmed Hosny</u>, <u>Farnoosh Khodakarami</u>, <u>Massive</u> Analysis Quality Control (MAQC) Society Board of Directors, Levi Waldron, Bo Wang, <u>Chris McIntosh</u>, <u>Anna</u> Goldenberg, <u>Anshul Kundaje</u>, <u>Casey S. Greene</u>, <u>Tamara Broderick</u>, <u>Michael M. Hoffman</u>, <u>Jeffrey T. Leek</u>, <u>Keegan Korthauer</u>, <u>Wolfgang Huber</u>, <u>Alvis Brazma</u>, <u>Joelle Pineau</u>, <u>Robert Tibshirani</u>, <u>Trevor Hastie</u>, <u>John P. A.</u> <u>Ioannidis</u>, John Quackenbush & Hugo J. W. L. Aerts

Nature 586, E14–E16 (2020) Cite this article

"The lack of details of the methods and algorithm code undermines its scientific value"



### Universality



#### nature

Explore content Y About the journal Y Publish with us Y

nature > articles > article

Article | Published: 01 January 2020

#### International evaluation of an AI system for breast cancer screening

Scott Mayer McKinney <sup>ICI</sup>, Marcin Sieniek, Varun Godbole, Jonathan Godwin, Natasha Antropova, Hutan Ashrafian, Trevor Back, Mary Chesus, Greg S. Corrado, Ara Darzi, Mozziyar Etemadi, Florencia Garcia-Vicente, Fiona J. Gilbert, Mark Halling-Brown, Demis Hassabis, Sunny Jansen, Alan Karthikesalingam, Christopher J. Kelly, Dominic King, Joseph R. Ledsam, David Melnick, Hormuz Mostofi, Lily Peng, Joshua Jay Reicher, Bernardino Romera-Paredes, Richard Sidebottom, Mustafa Suleyman, Daniel Tse <sup>ICI</sup>, Kenneth C. Young, Jeffrey De Fauw & Shravya Shetty <sup>ICI</sup> → Show fewer authors

Nature 577, 89-94 (2020) Cite this article

#### Test datasets

Image: Problem interpretation25,8563,097InterpretationDouble readingSingle readingScreening interval3 years1 or 2 yearsCancer follow-up39 months27 monthsNumber of cancers414 (1.6%)686 (22.2%)

Al vs. Doctor(s) Equivalent Superior



#### Fairness





medicine

() Check for updates

#### OPEN Underdiagnosis bias of artificial intelligence algorithms applied to chest radiographs in under-served patient populations

Laleh Seyyed-Kalantari<sup>1,2</sup><sup>1,2</sup>, Haoran Zhang<sup>3</sup>, Matthew B. A. McDermott<sup>3</sup>, Irene Y. Chen<sup>3</sup> and Marzyeh Ghassemi<sup>0,2,3</sup>

Artificial intelligence (AI) systems have increasingly achieved expert-level performance in medical imaging applications. However, there is growing concern that such AI systems may reflect and amplify human bias, and reduce the quality of their performance in historically under-served populations such as female patients, Black patients, or patients of low socioeconomic status. Such biases are especially troubling in the context of underdiagnosis, whereby the AI algorithm would inaccurately label an individual with a disease as healthy, potentially delaying access to care. Here, we examine algorithmic underdiagnosis in chest X-ray pathology classification across three large chest X-ray datasets, as well as one multi-source dataset. We find that classifiers produced using state-of-the-art computer vision techniques consistently and selectively underdiagnosed under-served patient populations and that the underdiagnosis rate was higher for intersectional under-served subpopulations, for example, Hispanic female patients. Deployment of AI systems using medical imaging for disease diagnosis with such biases risks exacerbation of existing care biases and can potentially lead to unequal access to medical treatment, thereby raising ethical concerns for the use of these models in the clinic.

| Attribute | Dataset                  | Avera<br>Gap                            | ge Cross-I<br>Lowest                         | Label Gap<br>Greatest                        | Unfavorable   | Favorable   |
|-----------|--------------------------|---|--|--|---|---|
| Sex       | ALL<br>CXP<br>CXR<br>NIH | <b>0.045</b><br>0.062<br>0.072<br>0.190 | Ef:0.001<br>Ed:0.000<br>Ed:0.011<br>M:0.001  | Pa:0.105<br>Co:0.139<br>EC:0.151<br>Cd:0.393 | Female (4/7)<br>Female (7/13)<br>Female (10/13)<br>Female (8/14)            | Male (4/7)<br>Male (7/13)<br>Male (10/13)<br>Male (8/14)  |
| Age       | ALL<br>CXR<br>CXP<br>NIH | <b>0.215</b><br>0.245<br>0.270<br>0.413 | Ef:0.115<br>SD:0.091<br>SD:0.084<br>In:0.188 | NF:0.444<br>Cd:0.440<br>NF:0.604<br>Em:1.00  | 0-20 (5/7)<br>0-20, 20-40 (7/13)<br>0-20, 20-40, 80- (7/13)<br>60-80 (7/14) | $\begin{array}{c} 40\text{-}60,60\text{-}80(5/7)\\ 60\text{-}80\ (10/13)\\ 40\text{-}60\ (8/13)\\ 20\text{-}40\ (9/14) \end{array}$ |
| Race      | CXR                      | 0.226                                   | NF:0.119                                     | Pa:0.440                                     | Hispanic (9/13)   | White (9/13)  |
| Insurance | CXR                      | 0.100                                   | SD:0.021                                     | PO:0.190                                     | Medicaid (10/13)  | Other (10/13)   |



#### Explainability















Article

https://doi.org/10.1038/s41467-024-46142-w

#### **Empirical data drift detection experiments on real-world medical imaging data**





#### Usability







Figure 2. Current frequency of artificial intelligence use in clinical practice.



## Characteristics Trustworthy AI



- Robustness
- Universality
- Fairness
- Explainability
- Traceability
- Usability

| F    | U         | Τ         | U      | R      | E           |
|------|-----------|-----------|--------|--------|-------------|
| FAIR | UNIVERSAL | TRACEABLE | USABLE | ROBUST | EXPLAINABLE |
|      |           | <u>õ</u>  | Rm     | (JED)  |             |



## Characteristics Trustworthy Al



Based on ethical principles and fundamental rights:







# Part 1 - What is trustworthy Al? Part 2 - How do we achieve it? Part 3 - Trustworthy Al in an ElC











#### RESEARCH METHODS AND REPORTING

#### OPEN ACCESS FUTURE-AI: international consensus guideline for trustworthy and deployable artificial intelligence in healthcare Check for updates

Karim Lekadir, 12 Alejandro F Frangi, 34 Antonio R Porras, 5 Ben Glocker, 6 Cella Cintas, 7 Curtis P Langlotz,<sup>8</sup> Eva Weicken,<sup>9</sup> Folkert W Asselbergs,<sup>10,11</sup> Fred Prior,<sup>12</sup> Gary S Collins,<sup>13</sup> Georgios Kalssis, 14 Glanna Tsakou, 15 Irène Buvat, 16 Jayashree Kalpathy-Cramer, 17 John Mongan, 18 Julia A Schnabel, 19 Kaisar Kushibar, 1 Katrine Riklund, 20 Kostas Marias, 21 Lameck M Amugongo, 22 Lauren A Fromont, 23 Lena Maier-Hein, 24 Leonor Cerdá-Alberich, 25 Luis Martí-Bonmatí, <sup>26</sup> M Jorge Cardoso, <sup>27</sup> Maciej Bobowicz, <sup>28</sup> Mahsa Shabani, <sup>29</sup> Manolis Tsiknakis,<sup>21</sup> Maria A Zuluaga,<sup>30</sup> Marie-Christine Fritzsche,<sup>31</sup> Marina Camacho,<sup>1</sup> Marius George Linguraru, 32 Markus Wenzel, 9 Marleen De Bruijne, 33 Martin G Tolsgaard, 34 Melanie Goisauf, 35 Mónica Cano Abadía, 35 Nikolaos Papanikolaou, 36 Noussair Lazrak, 1 Oriol Pujol,<sup>1</sup> Richard Osuala,<sup>1</sup> Sandy Napel,<sup>37</sup> Sara Colantonio,<sup>38</sup> Smritti Joshi,<sup>1</sup> Stefan Klein,<sup>33</sup> Susanna Aussó, 39 Wendy A Rogers, 40 Zohaib Salahuddin, 41 Martiin P A Starmans 33; on behalf of the FUTURE-AI Consortium

For numbered affiliations see and of the article Correspondence to: K Lekadir kartm.lekadir@ub.edu (ORCID 0000-0002-9456-1612) online only. To view please visit the journal online. http://dx.doi.org/10.1136/ bmi-2024-081554

Accepted: 10 january 2075

Despite major advances in artificial intelligence (AI) research for healthcare, the deployment and Additional material is published adoption of AI technologies remain limited in clinical practice. This paper Of this as BM/2025;388.0081554 describes the FUTURE-AI framework. which provides guidance for the development and deployment of trustworthy AI tools in healthcare. The FUTURE-AI Consortium was founded in 2021 and comprises 117 interdisciplinary experts from 50 countries representing all continents, including Al scientists, clinical researchers, biomedical ethicists, and social scientists. Over a two year period, the FUTURE-Al guideline was

#### SUMMARY POINTS

- Despite major advances in medical artificial intelligence (Al) research, clinical adoption of emerging AI solutions remains challenging owing to limited trust and ethical concerns
- The FUTURE-AI Consortium unites 117 experts from 50 countries to define International guidelines for trustworthy healthcare Al
- The FUTURE AI framework is structured around six guiding principles: fairness, universality, traceability, usability, robustness, and explainability
- The guideline addresses the entire Al lifecycle, from design and development to validation and deployment, ensuring alignment with real world needs and
- ethical requirements The framework includes 30 detailed recommendations for building trustworthy
- and deployable Al systems, emphasising multistakeholder collaboration Continuous risk assessment and mitigation are fundamental, addressing blases, data variations, and evolving challenges during the AI lifecycle FUTURE-AT is designed as a dynamic framework, which will evolve with
- technological advancements and stakeholder feedback

established through consensus based on six guiding principles-fairness. universality, traceability, usability, robustness, and explainability. To operationalise trustworthy AI in healthcare, a set of 30 best practices were defined, addressing technical, clinical, socioethical, and legal dimensions. The recommendations cover the entire lifecycle of healthcare Al, from design, development, and validation to regulation, deployment, and monitoring.

#### Introduction

In the field of healthcare, artificial intelligence (AI)-that ts, algorithms with the ability to self-learn logic-and data interactions have been increasingly used to develop computer aided models, for example, disease diagnosts, prognosts, prediction of therapy response or survival, and patient stratification.1 Despite major advances, the deployment and adoption of AI technologies remain limited in real world clinical practice. In recent years, concerns have been raised about the technical, clinical, ethical, and societal risks associated with healthcare AL23 In particular, existing research has shown that AI tools in healthcare can be prone to errors and patient harm, btases and increased health inequalities, lack of transparency and accountability, as well as data privacy and security breaches.48

To increase adoption in the real world, it is essential that AI tools are trusted and accepted by patients, clinicians, health organisations, and authorities, However, there is an absence of clear, widely accepted guidelines on how healthcare AI tools should be designed, developed, evaluated, and deployed to be trustworthy-that is, technically robust, clinically safe,





1



### Operationalisation



| FAIR | UNIVERSAL | TRACEABLE | USABLE | ROBUST | EXPLAINABLE |
|------|-----------|-----------|--------|--------|-------------|
|      |           | <u> K</u> | Pm     | CEP .  |             |

| Recommendations   | Research | Deployable |
|---|----------|------------|
| Fairness  |          |            |
| 1. Define any potential sources of bias from an early stage             | ++       | ++         |
| 2. Collect information on individuals' and data attributes              | +        | +          |
| 3. Evaluate potential biases and, when needed, bias correction measures | +        | ++         |

| Recommendations                 | Operations  | Examples   |
|---------------------------------|---|--|
| Define any potential            | Engage relevant stakeholders to define the sources of bias                  | Patients, clinicians, epidemiologists, ethicists, social carers <sup>97 98</sup>                                   |
| sources of bias (fairness<br>1) | Define standard attributes that might affect the AI tool's fairness         | Sex, age, socioeconomic status <sup>99</sup>   |
|                                 | Identify application specific sources of bias beyond<br>standard attributes | Skin colour for skin cancer detection, <sup>100 101</sup> breast density for breast cancer detection <sup>34</sup> |
|                                 | Identify all possible human biases  | Data labelling, data curation <sup>99</sup>  |



## Methodology



- Engage stakeholders (e.g. patients, clinicians, ethicists, ...)
- <u>Understand</u> all potential issues, risks, needs, ethical issues, etc:
  - 1. User requirements (incl. intended use)
  - 2. Intended clinical settings
  - 3. Sources of errors, data variations
  - 4. Sources of bias
  - 5. Explainability options
  - 6. Sources of performance degradation
  - 7. Application-specific ethical risks (e.g. in children)
- Implement suitable appropriate AI methods/mitigation measures
- <u>Evaluate</u> all dimensions of trustworthy AI, then <u>iterate</u> if needed
- <u>Report</u> everything, including benefits and <u>limitations</u>







- Understand intended use and user requirements, preferences, etc.
- Accordingly, implement <u>user-centred solutions</u>
- Examples: User-friendly <u>human-Al interfaces</u>
- Then, evaluate usability measures
- Examples: User satisfaction through <u>questionnaires</u>



## Universality



- Understand intended <u>clinical settings</u> and <u>variations</u>.
- <u>Examples</u>: High-income country, low-income country, <u>big hospital</u>, <u>rural clinic</u>.
- Evaluate applicability across the intended settings.
- Then, apply mitigation measures
- Examples: Train and test with data from <u>multiple clinical centres</u>.



#### Robustness



- Understand potential <u>sources of errors or failures</u>
- Examples: <u>Noise</u>, <u>motion of the patient</u> during scanning, <u>low- quality</u> <u>equipment</u>
- Then evaluate robustness and identify potential issues.
- Then apply mitigation measures, if necessary.
- Examples: Noise removal, data harmonisation, human oversight.







- Understand <u>sources of bias</u>
- Examples: <u>Sex</u>, <u>age</u>, <u>ethnicity</u>, <u>socio-economics</u>
- Then evaluate potential AI biases
- If necessary, apply mitigation measures
- Examples: <u>Data re-sampling</u>, <u>equalised odds techniques</u>



## Explainability



- Understand more suitable explainability options
- Example: Image heatmaps, feature importance, counterfactuals
- Evaluate the explanations with end-users
- If necessary, apply mitigation measures
- Examples: Use a <u>different explainable AI</u> method.



### Traceability



- Understand <u>sources of performance degradation</u> over time
- Examples: <u>Change in imaging protocols</u>
- Accordingly, implement specific mechanisms for monitoring the AI tool over time
- Examples: <u>Quality control</u> of input data, <u>yearly evaluation</u>, <u>model</u> <u>recalibration</u>



## Through AI Lifecycle





#### (4) DEPLOYMENT STAGE

- Un4: Evaluate local clinical validity
- Us3: Provide training materials and activities
- T3: Define mechanisms for quality control
- T4: Implement a periodic auditing system
- T5: Implement a logging system
- G5: Comply with AI regulatory requirements
- T6: Establish mechanisms for Al governance







| G1:  | Engage inter-disciplinary stakeholders          | (1)  |
|------|---|------|
| Us1: | Define intended use and user requirements       | (2)  |
| Un1: | Define clinical settings and related variations | (3)  |
| R1:  | Define all sources of data heterogeneity        | (4)  |
| F1:  | Define all sources of bias                      | (5)  |
| E1:  | Define explainability needs                     | (6)  |
| Un2: | Use community-defined standards                 | (7)  |
| G6:  | Investigate application-specific ethical issues | (8)  |
| G7:  | Investigate social and societal issues          | (9)  |
| T1:  | Define a risk management process                | (10) |



Stakeholder Engagement



| Best practice<br>(What)                                      | Practical steps<br>(How)                           | Examples<br>(References)  |
|--|--|---|
|  | Identify all relevant stakeholders                 | Patients, GPs, nurses, ethicists, data managers<br>(78,79)                              |
| Engage inter-<br>disciplinary<br>stakeholders<br>(General 1) | Provide information on the AI tool<br>and AI       | Educational seminars, training materials,<br>webinars (80)                              |
|  | Set up communication channels<br>with stakeholders | Regular group meetings, one-to-one interviews,<br>virtual platform (81)                 |
|  | Organise co-creation consensus<br>meetings         | One-day co-creation workshop with <i>n</i> =15 multi-<br>disciplinary stakeholders (82) |
|  | Use qualitative methods to gather feedback         | Online surveys, focus groups, narrative interviews<br>(83)                              |

























|  | PARTI   |   | Presi Prilippin  | N   | Prair 8 Degrade  | Pier II Property  | Part & Republishe  |
|--|---|---|--|---|--|---|--|
| Ceretos de lavorsligaciones<br>Services Barrolduce<br>y Neukoardelates | Notes particular de<br>antici dels faits  | Primary care<br>Primary care<br>Invy set image setury set   | Secondary Care   | Tertiary Care<br>Day per lange tradition  | Co<br>Diagnosile<br>Management Angust  | Regular follow-ups  | Hospitalization after diagnosi<br>Wai recent to its server of<br>recenters of the regit  |
| rt I: Patient<br>urney:<br>eart Fallure                                | Staps<br>Will show a device of whether<br>importing it with any   | Patients start their evaluation at<br>primary center-stept (General<br>Modicine Inserval modifier not some<br>allowed), ben withord to accordary (<br>general harpstal certificitigtin<br>according to sept2 evaluation patient<br>is referred to step 3 (tertiary center | At step 2 (General<br>cardiologist) first specific<br>exam 9 not invasive) Specific<br>HF treatment depending on<br>local resources and initial<br>stratification. | Heart Failure cardiologist-<br>invesive evoluation studies-<br>MR-CT-TEE. Patient is<br>stratified according their<br>etiology and risk | stepi-/ General Mediche /Internel mediche<br>not wann alterend/<br>Specific in finationality of the fination<br>features and infast statification<br>investore estation staties MRC-TTEL<br>Politent is statified according their sology<br>ent nik statief medicat strates by device-<br>tionage medication and statistication by device-<br>tionage/millionality.com | stept: 3m-12 months<br>stpe2: 3-24 months<br>step 3: 2 year   | Usually for Acute<br>heart failure   |
| t II: AI4HF In<br>lical practice:<br>en, what, how?                    | Freque trochest<br>Materia di la producta en<br>Materia di la producta en<br>Materia di la producta en<br>Antonio di la producta en<br>An | General medicine<br>Internal medicine   | General cardiologist   | Heart failure<br>cardiologist   | Cardiologist   | cardiologist<br>nurse   | cardiologist   |
|  | When does the day<br>bits pilot<br>the day of the day<br>the day of the day of the day<br>the day of the day<br>of the day of the day of the day of the day<br>of the day of the day of the day of the day<br>of the day of the day of the day of the day of the day<br>of the day of the day of the day of the day of the day<br>of the day of the day<br>of the day of  | [ Withs down all physical<br>locations where this step<br>can take place. E.g.,<br>heathcare canae.<br>httpspite. GP clinic, etc. ]   | [Withe cover all physical<br>locations where this step<br>can take place. Eg.<br>heatboare contact,<br>heapted. (IP clinic, etc. ]                                 | [With down wil physical<br>locations where this step<br>can alike place. E.g.,<br>heathcare center,<br>hospital, GP clinic, etc. ]      | [Withis down all physical<br>locations where this step<br>can take place. E.g.,<br>heathbaire cartillo,<br>heapital, GP clinic, wit: ]   | [Witte down all physical<br>locations where this step<br>can take proce. E.g.,<br>healthcare contex<br>heapter OP sinic, etc. ] | [With down all physical<br>locations where this step<br>can take place. E.g.<br>healthcare center,<br>heapine, GP clinic, etc. ] |
|  | Time duration<br>The needs and the angle<br>angle of the angle of the angle<br>angle of the angle of t  | [ How much time does<br>this step usually taked-<br>until the potent moves<br>to the next step? ]   | [ How much time does<br>this step usually takes<br>until the patient moves<br>to the next step? ]  | [How much time does<br>this step usually takes<br>until the patient moves<br>to the next step? ]  | [ How much time does<br>this stop issuely takes<br>until the patient movies<br>to the next step? ]   | [How much time does<br>this step usually takes<br>until the patient moves<br>. As the next step?]                               | [ How much time does<br>this step usually takes<br>until the patient moves<br>to the next step? ]                                |
|  | Use of Aldsheen<br>prediction back<br>what has now of and prime<br>frame region with and prime<br>frame region with and prime   | annen andre   | E C C  | ······  | ······   | ······································  | ······   |



#### RadioVal FUTURE Al Survey

#### Time required: 45 minutes

#### <u>Goal</u>

We are conducting a survey on trustworthy artificial intelligence (AI) for breast cancer treatment planning. Our aim is to reach healthcare experts across various institutions to get an insight into their day-to-day practice: issues encountered, agreements and more importantly, disagreements regarding clinical practices, definitions, and protocols to set requirements for our tool. We also want to get their perspective on AI tools and ways to make them more reliable and applicable in clinical settings.

3. What sources of bias do you think we should possibly take into account? \*



4. What is your profession? \*

Radiologist

Oncologist

Radiation Oncologist

Surgeon

General practitioner

Pathologist

Other

3. What information should be monitored over time? \*

|             | Definitely | Maybe | Not<br>needed |
|-------------|------------|-------|---------------|
| Usage       | 0          | 0     | 0             |
| Errors      | 0          | 0     | 0             |
| Limitations | 0          | 0     | 0             |





8. Rank the following variables based on their importance for bias estimation



12. Rank the following variables based on their importance for bias estimation

9. Do you have information on ethnicity in your centre/country?







Best Practices – <u>Design</u>



| G1:  | Engage inter-disciplinary stakeholders          | (1)  |
|------|---|------|
| Us1: | Define intended use and user requirements       | (2)  |
| Un1: | Define clinical settings and related variations | (3)  |
| R1:  | Define all sources of data heterogeneity        | (4)  |
| F1:  | Define all sources of bias                      | (5)  |
| E1:  | Define explainability needs                     | (6)  |
| Un2: | Use community-defined standards                 | (7)  |
| G6:  | Investigate application-specific ethical issues | (8)  |
| G7:  | Investigate social and societal issues          | (9)  |
| T1:  | Define a risk management process                | (10) |



#### Intended Use





Secondary risk prevention in heart failure

#### Clinicians

#### Patients

What should the AI tool predict ?

- Change in cardiac function
- Risk of myocardial infarction
- o Risk of mortality

What should the AI tool predict ?

- o Risk of fatigue
- o Risk of backpain
- o Risk of hospital re-admission



Best Practices – <u>Design</u>



| G1:  | Engage inter-disciplinary stakeholders          | (1)  |
|------|---|------|
| Us1: | Define intended use and user requirements       | (2)  |
| Un1: | Define clinical settings and related variations | (3)  |
| R1:  | Define all sources of data heterogeneity        | (4)  |
| F1:  | Define all sources of bias                      | (5)  |
| E1:  | Define explainability needs                     | (6)  |
| Un2: | Use community-defined standards                 | (7)  |
| G6:  | Investigate application-specific ethical issues | (8)  |
| G7:  | Investigate social and societal issues          | (9)  |
| T1:  | Define a risk management process                | (10) |



#### Data Heterogeneity





(1) Canon



(2) General Electric







(4) Siemens



#### electrode movement







Venton et al. 2021. Robustness CNNs to physiological ECG noise.



Best Practices – <u>Design</u>



| G1:  | Engage inter-disciplinary stakeholders          | (1)  |
|------|---|------|
| Us1: | Define intended use and user requirements       | (2)  |
| Unl: | Define clinical settings and related variations | (3)  |
| R1:  | Define all sources of data heterogeneity        | (4)  |
| F1:  | Define all sources of bias                      | (5)  |
| E1:  | Define explainability needs                     | (6)  |
| Un2: | Use community-defined standards                 | (7)  |
| G6:  | Investigate application-specific ethical issues | (8)  |
| G7:  | Investigate social and societal issues          | (9)  |
| T1:  | Define a risk management process                | (10) |



#### Sources of Biases











## Sources of Biases







Lima: Sea level



Arequipa: 2,335 m



Cusco: 3,400 m



Rinconada: 5,100 m



Best Practices – <u>Design</u>



| G1:  | Engage inter-disciplinary stakeholders          | (1)  |
|------|---|------|
| Us1: | Define intended use and user requirements       | (2)  |
| Un1: | Define clinical settings and related variations | (3)  |
| R1:  | Define all sources of data heterogeneity        | (4)  |
| F1:  | Define all sources of bias                      | (5)  |
| E1:  | Define explainability needs                     | (6)  |
| Un2: | Use community-defined standards                 | (7)  |
| G6:  | Investigate application-specific ethical issues | (8)  |
| G7:  | Investigate social and societal issues          | (9)  |
| T1:  | Define a risk management process                | (10) |



#### Explainability











| G2:  | Define measures for privacy and security       | (11) |
|------|--|------|
| F2:  | Collect data on individuals' attributes        | (12) |
| R2:  | Collect representative real-world data         | (13) |
| G3:  | Implement measures against identified AI risks | (14) |
| Us2: | Implement human-Al interaction mechanisms      | (15) |



#### Individuals' Attributes



| Patient 🔽  | Age 🔽 | Sex 🔽  | Ethnicity 🔽 | Neighbourhood 🚽 | Altitude 🔽 | Skin colour  🚽 | Education 🔽 |
|------------|-------|--------|-------------|-----------------|------------|----------------|-------------|
| Patient001 | 20    | Male   | E. Europe   | Gracia          | 0          | White          | High-School |
| Patient002 | 20    | Male   | S. Europe   | Gracia          | 0          | White          | University  |
| Patient003 | 30    | Male   | N. Africa   | Horta           | 500        | White          | University  |
| Patient004 | 30    | Male   | N. Africa   | Horta           | 500        | White          | University  |
| Patient005 | 40    | Male   | S. Africa   | Poblenou        | 1000       | Black          | High-School |
| Patient006 | 40    | Female | S. Africa   | Poblenou        | 1000       | Black          | High-School |
| Patient007 | 50    | Female | S. Asia     | Gervasi         | 1500       | White          | High-School |
| Patient008 | 50    | Female | E. Asia     | Gervasi         | 1500       | White          | University  |
| Patient009 | 60    | Female | L. America  | Eixample        | 3000       | White          | University  |
| Patient010 | 60    | Female | L. America  | Eixample        | 3000       | Black          | University  |



### Training Data







#### High-quality data

#### Heterogeneous data



**Private hospital** 



#### **Public hospital**







| G2:  | Define measures for privacy and security       | (11) |
|------|--|------|
| F2:  | Collect data on individuals' attributes        | (12) |
| R2:  | Collect representative real-world data         | (13) |
| G3:  | Implement measures against identified AI risks | (14) |
| Us2: | Implement human-AI interaction mechanisms      | (15) |



### Universality





Contents lists available at ScienceDirect

Artificial Intelligence In Medicine

journal homepage: www.elsevier.com/locate/artmed



Chieck for

Research paper

Domain generalization in deep learning based mass detection in mammography: A large-scale multi-center study

Lidia Garrucho<sup>\*</sup>, Kaisar Kushibar, Socayna Jouide, Oliver Diaz, Laura Igual, Karim Lekadir Artificial Intelligence in Medicine Lab (BCN-AIM), Faculty of Mathematics and Computer Science, University of Barcelona, Gran Via de les Corts Catalanes 585, Barcelona, 08007, Barcelona, Spain







**Baseline model** 

Domain adaptation







Frontiers | Frontiers in Oncology

TYPE Original Research PUBLISHED 23 January 2023 DOI 10.3389/fonc.2022.1044496

#### High-resolution synthesis of high-density breast mammograms: Application to improved fairness in deep learning based mass detection

Lidia Garrucho<sup>1\*</sup>, Kaisar Kushibar<sup>1</sup>, Richard Osuala<sup>1</sup>, Oliver Diaz<sup>1</sup>, Alessandro Catanese<sup>2</sup>, Javier del Riego<sup>3</sup>, Maciej Bobowicz<sup>4</sup>, Fredrik Strand<sup>5</sup>, Laura Igual<sup>1</sup> and Karim Lekadir<sup>1</sup>

<sup>1</sup>Barcelona Artificial Intelligence in Medicine Lab, Facultat de Matemàtques i Informàtica, Universitat de Barcelona, Barcelona, Spain, <sup>2</sup>Unitat de Diagnóstic per la Imatge de la Mama (UDIM), Hospital Germans Trias i Pujol, Badalona, Spain, <sup>3</sup>Área de Radiología Mamaria y Ginecólogica (UDIAT CD), Parc Taulí Hospital Universitari, Sabadell, Spain, <sup>4</sup>2nd Department of Radiology, Medical University of Gdansk, Gdansk, Poland, <sup>3</sup>Breast Radiology, Karolinska University Hospital and Department of Oncology-Pathology, Karolinska Institutet; Stockholm, Sweden





#### Fairness



Prontiers | Frontiers in Oncology

TYPE Original Research PUBLISHED 23 January 2023 DOI 10.3389/fonc.2022.1044496

High-resolution synthesis of high-density breast mammograms: Application to improved fairness in deep learning based mass detection

Lidia Garrucho<sup>3\*</sup>, Kaisar Kushibar<sup>1</sup>, Richard Osuala<sup>1</sup>, Oliver Diaz<sup>1</sup>, Alessandro Catanese<sup>2</sup>, Javier del Riego<sup>3</sup>, Maciej Bobowicz<sup>4</sup>, Fredrik Strand<sup>5</sup>, Laura Igual<sup>1</sup> and Karim Lekadir<sup>1</sup>

<sup>1</sup>Barcelona Artificial Intelligence in Medicine Lab, Facultat de Matemàtques i Informàtica, Universitat de Barcelona, Barcelona, Spain, <sup>3</sup>Unitat de Diagnóstic per la Imatge de la Mama (UDIM), Hospital Germans Trias i Pujol, Badalona, Spain, <sup>3</sup>Área de Radiología Mamaria y Ginecólogica (UDIAT CD), Parc Taulí Hospital Universitari, Sabadell, Spain, <sup>4</sup>2nd Department of Radiology, Medical University of Gdansk, Gdansk, Poland, <sup>5</sup>Breast Radiology, Karolinska University Hospital and Department of Oncology-Pathology, Karolinska Institutet, Stockholm, Sweden









| G4:  | Define an adequate evaluation plan             | (16) |
|------|--|------|
| Un3: | Evaluate using external and/or multi-site data | (17) |
| R3:  | Evaluate robustness against real variations    | (18) |
| F3:  | Evaluate fairness and debiasing measures       | (19) |
| Us3: | Evaluate user experience and acceptance        | (20) |
| Us4: | Evaluate clinical utility and safety           | (21) |
| E2:  | Evaluate explainability with end-users         | (22) |
| T2:  | Document the AI tool including evaluations     | (23) |



## Universality Evaluation





Instituto de Investigación Sanitaria La Fe











## **Q**radıoval



#### Fairness Evaluation



#### scientific reports

Check for updates

OPEN Fairness and bias correction in machine learning for depression prediction across four study populations

> Vien Ngoc Dang<sup>1123</sup>, Anna Cascarano<sup>1</sup>, Rosa H. Mulder<sup>2,3</sup>, Charlotte Cecil<sup>2,4,5</sup>, Maria A. Zuluaga<sup>6</sup>, Jerónimo Hernández-González<sup>7</sup> & Karim Lekadir<sup>1,8</sup>

A significant level of stigma and inequality exists in mental healthcare, especially in under-served populations. Inequalities are reflected in the data collected for scientific purposes. When not properly accounted for, machine learning (ML) models learned from data can reinforce these structural inequalities or biases. Here, we present a systematic study of bias in ML models designed to predict depression in four different case studies covering different countries and populations. We find that standard ML approaches regularly present biased behaviors. We also show that mitigation techniques, both standard and our own post-hoc method, can be effective in reducing the level of unfair bias. There is no one best ML model for depression prediction that provides equality of outcomes. This emphasizes the importance of analyzing fairness during model selection and transparent reporting about the impact of debiasing interventions. Finally, we also identify positive habits and open challenges that practitioners could follow to enhance fairness in their models.



#### Standard AI models



## Usability Evaluation



Human evaluators in 5 sites:

- ✓ 2 GPs at each site
- ✓ 2 cardiologists at each site
- ✓ 2 nurses at each site
- $\checkmark$  7 patients for each clinician
- ✓ 2 IT/data manager
- 50% male + 50% female
- 50% early-career, 50% > 5-year experience













| Un4: | Evaluate local clinical validity          | (24) |
|------|---|------|
| Us3: | Provide training materials and activities | (25) |
| T3:  | Define mechanisms for quality control     | (26) |
| T4:  | Implement a periodic auditing system      | (27) |
| T5:  | Implement a logging system                | (28) |
| G5:  | Comply with AI regulatory requirements    | (29) |
| T6:  | Establish mechanisms for Al governance    | (30) |



### Quality Control





Nikiforaki et al. 2024. Image Quality Assessment Tool for Conventional and Dynamic Magnetic Resonance Imaging Acquisitions. Journal of Imaging, 10(5), p.115.





# Part 1 - What is trustworthy AI? Part 2 - How do we achieve it? Part 3 - Trustworthy AI in an EIC



Objectives



Objective 1: Design, develop and demonstrate the first trustworthy AI technology for personalised risk assessment and improved management of HF patients in clinical cardiology.

Objective 2: Implement a human-centred, multi-stakeholder, inclusive approach to improve awareness, acceptance and promotion of trustworthy AI in cardiovascular risk assessment.

Objective 4: Implement the very first international, multi-faceted clinical validation study for AI-powered risk assessment in cardiology in multiple developed and developing countries.

Objective 1: Design, develop and evaluate a trustworthy and ethical AI tool for early cancer detection



## Method Section



- 1. Stakeholder engagement
- 2. Training and testing datasets
- 3. AI methods implementation
- 4. Al validation studies
- 5. Reporting, dissemination, exploitation





| Understand clinical needs, potential barriers    | <ul> <li>1-2 radiologists and 1-2 physicians per SSA site (Total ~20).</li> <li>1 healthcare manager per SSA site (Total 10).</li> </ul> |  |  |  |
|--|--|--|--|--|
| and chinical features of the AI tools.           | • 10 patients / advocates (e.g. patients at local sites, StopTB local groups).   |  |  |  |
|  |  |  |  |  |
| Define nathways towards implementation and       | • 15 local policy makers (e.g. health ministries, regional health agencies).   |  |  |  |
| adoption of the AI tools in real-world practice  | • 7 policy / advocacy organisations (e.g. StopTB, WHO-Africa, National   |  |  |  |
| adoption of the AI tools in real-world practice. | Tuberculosis Control Programs, Digital Health Africa, etc).  |  |  |  |

<u>Participatory action research (PAR)</u>: We will employ various PAR methods, such as qualitative interviews, focus groups and co-creation workshops, to gather insights on the needs, views and concerns of SSA stakeholders. At the start of the project, CISM, KUHES and ISGLOBAL will elaborate a qualitative master protocol with all SSA site, which will be then adapted by each SSA site to ensure relevance to the local contexts. The local social scientists will then collect data through direct interactions with the local stakeholders, employing narrative interviews or focus groups as needed. For instance, each SSA partner will organise focus groups with n=7-10 local patients to discuss ethical concerns and the patient's views on (i) the clinical sites waiving consent on their behalf for the public good, (ii) the application of AI to diagnose their family members or children, (iii) the monetisation of their data to gather extra resources for the healthcare facilities. All qualitative data will be transcribed verbatim then structured through a qualitative data management software (Nvivo).



Datasets



| Name       | Country      | Partner         | Туре     | Number of<br>HF<br>patients | Harmonised | Clinical<br>reports | Cardiac<br>imaging | ECG signals | Cardiac<br>biomarkers | Clinical<br>end-points |
|------------|--------------|-----------------|----------|-----------------------------|------------|---------------------|--------------------|-------------|-----------------------|------------------------|
| Dataset    | 1 (Available | at the start of | AI4HF)   | >715,000                    |            |                     |                    | ,           | 1                     |                        |
| CALIBER    | UK           | NLHI            | EHR      | 502,536                     | Y          | S                   | N                  | N           | Y                     | Y                      |
| UK Biobank | UK           | UOXF            | Cohort   | 8,000                       | Y          | S                   | I+R                | S           | Y                     | Y                      |
| SwedeHF    | SE           | UMCU            | Registry | 150,000                     | Y          | S                   | R                  | R           | Y                     | Y                      |
| ABUCASIS   | ES           | NLHI            | EHR      | 48,000                      | Y          | S                   | R                  | S           | Y                     | Y                      |
| UPOD       | NI.          | UMCU            | EHR      | 8 000                       | Y          | S                   | R                  | S+R         | Y                     | Y                      |

Table 2 – List of selected cardiovascular cohorts to be leveraged in AI4HF



## Al Methods / Validation



*Fairness*: To assess and address AI biases (*e.g.* with respect to sex or ethnicity), we will include dedicated metrics (*e.g.* Statistical Parity, Group Fairness), and bias mitigation techniques (*e.g.* adaptive data re-sampling, equalised odds post-processing).

<u>Universality</u>: To ensure wide applicability of the AI solutions, the toolkit will comprise transfer learning and knowledge distillation techniques to optimise and calibrate pre-existing AI models for each new clinical site across different SSA countries.

<u>Traceability</u>: We will implement statistical tests (*e.g.* population stability index, performance variability) to identify drifts in the input images or AI decisions over time. We will include continuous learning methods to re-calibrate the AI models.

<u>Usability</u>: Existing usability questionnaires, such as the System Usability Scale, will be updated and adapted to local contexts of African clinicians, so they can be utilised to assess acceptance and applicability of newly developed AI solutions.

<u>*Robustness*</u>: Methods for assessing robustness (*e.g.* against variations in data quality or image scanners across SSA sites) will be implemented, together with image harmonisation and domain adaptation techniques to enhance robustness.

*Explainability*: Explainable AI methods will be provided in the open toolkit, such as heatmaps, depending on the needs and preferences of the local clinicians for interpreting visually the decisions of the AI tools on the input images.



#### AI Validation



*In-silico validation of AI-based TB diagnosis from adult CUS*: This study will implement the same logical steps as above across five different SSA partners, in Gabon, Malawi, Gambia, Ghana, and Uganda. We will evaluate the AI tool's performance in adult TB diagnosis with CUS images, with a focus on its adaptability and consistency across various CUS machines and setups that are prevalent in SSA settings. Given the operator-dependent nature of CUS, we will specifically assess how variations in image acquisition by different sonographers influence the AI's diagnostic accuracy. Five radiologists, sonographers, and general practitioners per site will participate to evaluate the AI tool's accessibility and practicality in typical clinical environments. The trial will also explore the effectiveness of explainable AI features in aiding clinicians to identify TB-related abnormalities in CUS such as pleural effusion and echogenic patterns indicative of fibrosis, calcification, or fluid accumulation. Importantly, we will thoroughly compare AI-aided diagnosis of TB in adults with CUS with standard of care based on CXR expert interpretation.



## Gender Dimensions



#### 1.2.15 Gender Dimension

It is established that there are sex differences in respiratory disease, both in terms of pathophysiology and disease manifestation<sup>39</sup>. Hence, in AFRICAI-RI, sex differences will be taken into consideration from day one, as follows:

- First, during the requirements analysis in WP1, our social scientists will ensure a balanced representation in terms of sex and gender in the stakeholder engagement activities, to obtain unbiased requirements.
- During data preparation in WP2-3, we will carefully estimate the distribution of men and women for all data sources, thus ensuring potential bias is reported or corrected (*e.g.* using sampling methods).
- Bias estimation will be applied in WP4-5 to assess if the AI tools are biased or fair across male and female subjects. Correction measures will be systematically applied to correct for any identified biases,
- The selection of patients and clinicians for the multi-site in silico trials in WP6 will ensure a balanced representation and the analyses will be stratified for potential sex differences.
- In WP7, the early career researchers and students will be recruited ensuring a perfect gender balance.
- Finally, findings and recommendations will be disseminated in WP8-9 to promote accessible AI-powered imaging diagnostics of respiratory diseases across all subgroups, including male and female patients.



WPs / Tasks



**T1.2. Qualitative studies (<u>CISM</u>**, ISGLOBAL, ALL) [M1-M48]: At M1-M3, CISM and ISGLOBAL will coordinate elaborate a <u>qualitative master protocol</u> with all SSA sites. At M3-M12, data will be continuously collected through <u>site workshops and individual interviews</u> with stakeholders in all SSA sites (as described in Section 1.2.4). <u>Content and thematic analyses</u> will be used to analyse data (Nvivo software) and extract requirements. The process will be repeated throughout the project depending on emerging questions and needs.

**T1.3. Co-creation activities and consensus requirements (CISM**, ISGLOBAL, ALL) [M9-M48]: Starting from M9, the results will be shared through local workshops to foster feedback. MRCG will organise an in-person co-creation workshop in The Gambia during the first annual meeting to resolve any misalignments (*e.g.* a technical requirement does not align with ethical concerns). Follow-up online co-creation meeting will be convened regularly to refine all requirements. Requirements documents and scientific publications will be elaborated by all sites, coordinated by CISM.



WPs / Tasks



**T3.3. Robust AI methods for cross-centre generalisability** (<u>UOXF</u>, UB, NLHI, UMCU, VHIR, ICRC, MUHAS, INCOR, BSC, SRDC) [M12-M36]:

- At M12-M18, the baseline models will be assessed with respect to the sources of heterogeneity identified during requirement analysis in WP1 (*e.g.* variations in biomarkers, ECG machines, image scanners, data quality), by performing "leave one centre out" tests based on the 5 cohorts of Dataset 1.
- At M18-M24, UOXF and UB will implement and assess several strategies for enhanced AI robustness across datasets and centres, such as data synthesis, transfer learning, domain adaptation and knowledge distillation.

**T3.4. Bias detection and mitigation methods** (<u>UB</u>, NLHI, UOXF, UMCU, VHIR, ICRC, MUHAS, INCOR, CERTH) [M12-M36]:

- At M12-M18, the baseline models will be tested for their fairness depending on the sources of bias identified in WP1 (*e.g. sex, comorbidity*), by implementing metrics such as Statistical Parity and Group Fairness.
- At M18-M27, based on Dataset 1, UB will implement bias mitigation measures, such as 1) pre-processing approaches (*e.g.* re-sampling or data augmentation), (2) in-processing approaches that explicitly remove discrimination during model training, and (3) post-processing *e.g.* the equalised odds technique.



## Final Point



- Space is limited (20 pages): How much for AI (1 or 10 pages)?
- Reviewers could be <u>AI experts</u> or/and <u>domain experts</u>
- Depends on the EIC grant:
  - Al has a major role: It's an Al-focused EIC grant proposal
  - Al is prominent, but it's not the core of the EIC proposal
  - Al has a minor role: It's just a method amongst others
  - There is no AI at all



#### RESEARCH METHODS AND REPORTING

#### OPEN ACCESS FUTURE-AI: international consensus guideline for trustworthy and deployable artificial intelligence in healthcare Check for updates

Karim Lekadir, 12 Alejandro F Frangi, 34 Antonio R Porras, 5 Ben Glocker, 6 Cella Cintas, 7 Curtis P Langlotz,<sup>8</sup> Eva Weicken,<sup>9</sup> Folkert W Asselbergs,<sup>10,11</sup> Fred Prior,<sup>12</sup> Gary S Collins,<sup>13</sup> Georgios Kalssis, 14 Glanna Tsakou, 15 Irène Buvat, 16 Jayashree Kalpathy-Cramer, 17 John Mongan, 18 Julia A Schnabel, 19 Kaisar Kushibar, 1 Katrine Riklund, 20 Kostas Marias, 21 Lameck M Amugongo, 22 Lauren A Fromont, 23 Lena Maier-Hein, 24 Leonor Cerdá-Alberich, 25 Luis Martí-Bonmatí, <sup>26</sup> M Jorge Cardoso, <sup>27</sup> Maciej Bobowicz, <sup>28</sup> Mahsa Shabani, <sup>29</sup> Manolis Tsiknakis,<sup>21</sup> Maria A Zuluaga,<sup>30</sup> Marie-Christine Fritzsche,<sup>31</sup> Marina Camacho,<sup>1</sup> Marius George Linguraru, 32 Markus Wenzel, 9 Marleen De Bruijne, 33 Martin G Tolsgaard, 34 Melanie Goisauf, 35 Mónica Cano Abadía, 35 Nikolaos Papanikolaou, 36 Noussair Lazrak, 1 Oriol Pujol,<sup>1</sup> Richard Osuala,<sup>1</sup> Sandy Napel,<sup>37</sup> Sara Colantonio,<sup>38</sup> Smritti Joshi,<sup>1</sup> Stefan Klein,<sup>33</sup> Susanna Aussó, 39 Wendy A Rogers, 40 Zohaib Salahuddin, 41 Martiin P A Starmans 33; on behalf of the FUTURE-AI Consortium

For numbered affiliations see and of the article Correspondence to: K Lekadir kartm.lekadir@ub.edu (ORCID 0000-0002-9456-1612) online only. To view please visit the journal online. http://dx.doi.org/10.1136/ bmi-2024-081554

Accepted: 10 january 2075

Despite major advances in artificial intelligence (AI) research for healthcare, the deployment and Additional material is published adoption of AI technologies remain limited in clinical practice. This paper Of this as BM/2025;388.0081554 describes the FUTURE-AI framework. which provides guidance for the development and deployment of trustworthy AI tools in healthcare. The FUTURE-AI Consortium was founded in 2021 and comprises 117 interdisciplinary experts from 50 countries representing all continents, including Al scientists, clinical researchers, biomedical ethicists, and social scientists. Over a two year period, the FUTURE-Al guideline was

#### SUMMARY POINTS

- Despite major advances in medical artificial intelligence (Al) research, clinical adoption of emerging AI solutions remains challenging owing to limited trust and ethical concerns
- The FUTURE-AI Consortium unites 117 experts from 50 countries to define International guidelines for trustworthy healthcare Al
- The FUTURE AI framework is structured around six guiding principles: fairness, universality, traceability, usability, robustness, and explainability
- The guideline addresses the entire Al lifecycle, from design and development to validation and deployment, ensuring alignment with real world needs and
- ethical requirements The framework includes 30 detailed recommendations for building trustworthy
- and deployable Al systems, emphasising multistakeholder collaboration Continuous risk assessment and mitigation are fundamental, addressing blases, data variations, and evolving challenges during the AI lifecycle FUTURE-AT is designed as a dynamic framework, which will evolve with
- technological advancements and stakeholder feedback

established through consensus based on six guiding principles-fairness. universality, traceability, usability, robustness, and explainability. To operationalise trustworthy AI in healthcare, a set of 30 best practices were defined, addressing technical, clinical, socioethical, and legal dimensions. The recommendations cover the entire lifecycle of healthcare Al, from design, development, and validation to regulation, deployment, and monitoring.

#### Introduction

In the field of healthcare, artificial intelligence (AI)-that ts, algorithms with the ability to self-learn logic-and data interactions have been increasingly used to develop computer aided models, for example, disease diagnosts, prognosts, prediction of therapy response or survival, and patient stratification.1 Despite major advances, the deployment and adoption of AI technologies remain limited in real world clinical practice. In recent years, concerns have been raised about the technical, clinical, ethical, and societal risks associated with healthcare AL23 In particular, existing research has shown that AI tools in healthcare can be prone to errors and patient harm, btases and increased health inequalities, lack of transparency and accountability, as well as data privacy and security breaches.48

To increase adoption in the real world, it is essential that AI tools are trusted and accepted by patients, clinicians, health organisations, and authorities, However, there is an absence of clear, widely accepted guidelines on how healthcare AI tools should be designed, developed, evaluated, and deployed to be trustworthy-that is, technically robust, clinically safe,





1



#### Many Thanks!



















## **Q**radioval



karim.lekadir@ub.edu