



UNIVERSITAT DE
BARCELONA



How to Ensure Trustworthiness in AI for Healthcare: The FUTURE-AI Guideline

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



Part 1 - What is trustworthy AI?

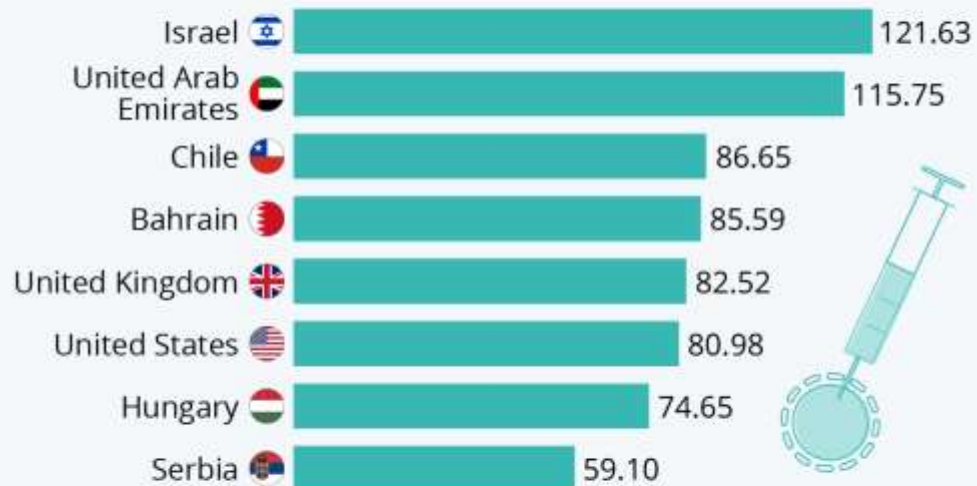
Part 2 - How do we achieve it?



Trustworthiness

The Countries With The Highest Rate Of Covid-19 Vaccination

Covid-19 vaccination doses administered per 100 people (May 17, 2021)*



* Numbers counted as a single dose and may not equal the total number of people vaccinated.

Source: Our World in Data



statista

Novak Djokovic: Tennis star deported after losing Australia visa battle

16 January 2022

B B C





AI in Cardiology

BBC Sign in Home News Sport Reel Worklife

NEWS

Home | War in Ukraine | Climate | Video | World | UK | Business | Tech | Science | Stories

Tech

NHS uses AI scan to detect hidden heart disease

© 29 March 2021 · Comments



GETTY IMAGES

The technology could help save "thousands of lives"

**The
Guardian**

AI eye checks can predict heart disease risk in less than minute, finds study

Breakthrough opens door to a highly effective, non-invasive test that does not need to be done in a clinic



📷 Ophthalmologists may soon be able to carry out cardiovascular screening by checking the retina - without the need for blood tests. Photograph: Zorica Nastasic/Getty Images


Robustness

Robustness of convolutional neural networks to physiological electrocardiogram noise

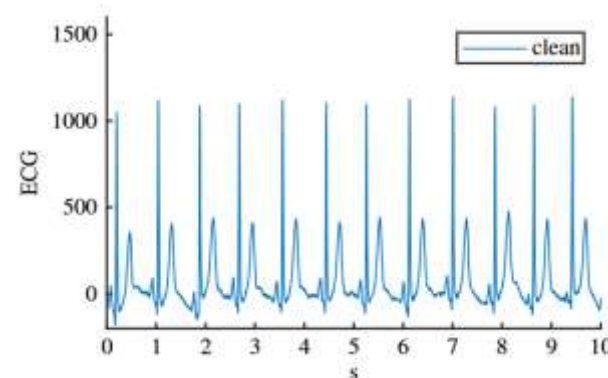
J. Venton¹, P. M. Harris¹, A. Sundar¹, N. A. S. Smith¹
and P. J. Aston^{1,2}

¹Department of Data Science, National Physical Laboratory,
Teddington, UK

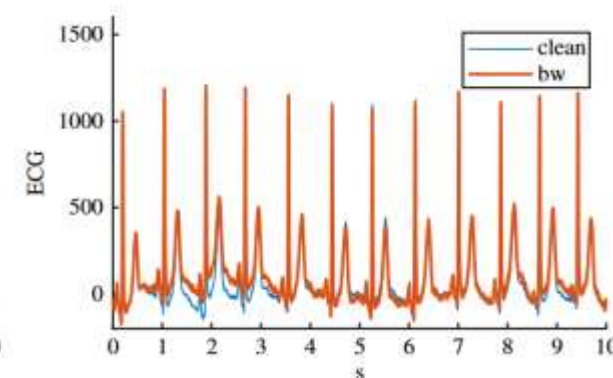
²Department of Mathematics, University of Surrey, Guildford, UK

 JV, 0000-0003-0547-1226

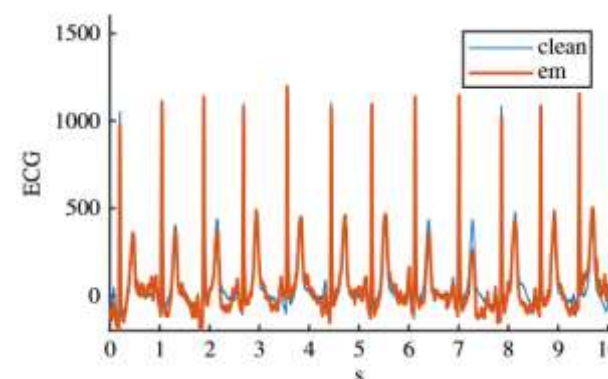
PHILOSOPHICAL
TRANSACTIONS A



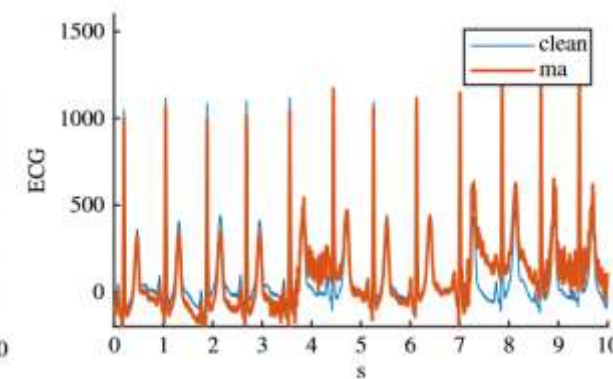
clean



baseline wander



electrode movement



motion artefact



Universality





JACC: Advances
Volume 3, Issue 9, Part 2, September 2024, 101202

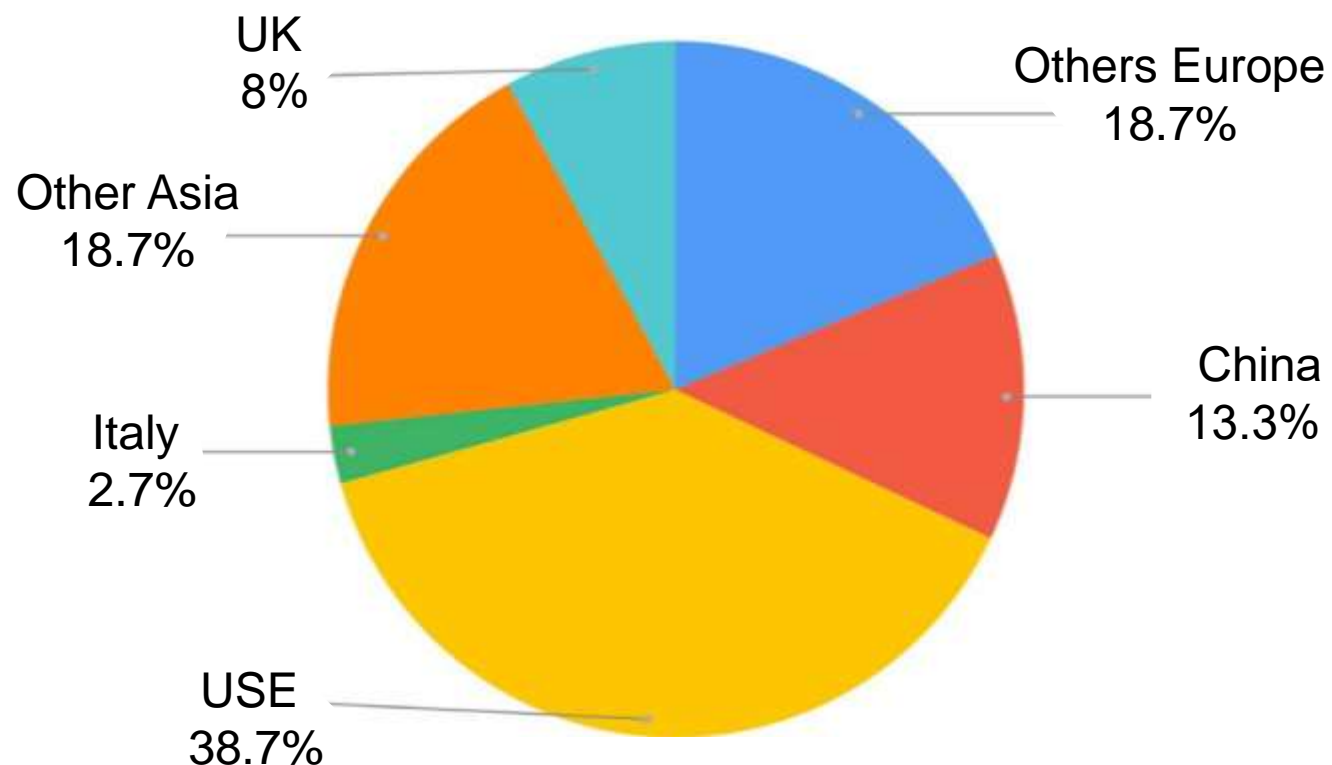


Original Research

Outcomes and Quality

Prospective Human Validation of Artificial Intelligence Interventions in Cardiology: A Scoping Review

Amirhossein Moosavi PhD ^{a b *}, Steven Huang ^{b *}, Maryam Vahabi MSc ^{a b},
Bahar Motamedivafa BSc ^{a b}, Nelly Tian MBAn ^c, Rafid Mahmood PhD ^a, Peter Liu MD ^b,
Christopher L.F. Sun PhD ^{a b}  





Fairness

Circulation: Heart Failure

Volume 17, Issue 1, January 2024; Page e010879






<https://doi.org/10.1161/CIRCHEARTFAILURE.123.010879>

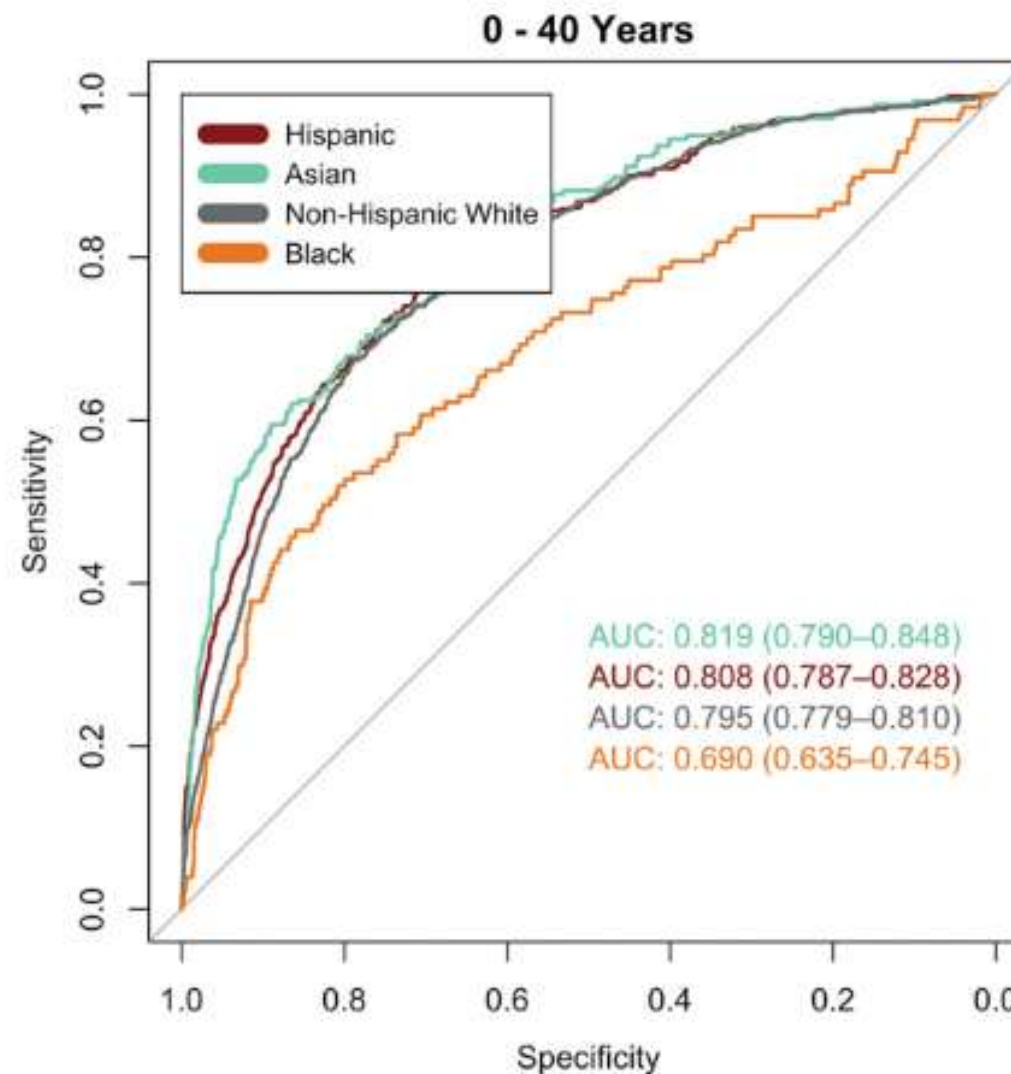


ORIGINAL ARTICLE

Race, Sex, and Age Disparities in the Performance of ECG Deep Learning Models Predicting Heart Failure

See Editorial by [Rosenberg](#)

Dhamanpreet Kaur, BS , J. Weston Hughes, BA, Albert J. Rogers, MD, MBA , Guson Kang, MD, Sanjiv M. Narayan, MD, PhD , Euan A. Ashley, DPhil , and Marco V. Perez, MD 





Traceability

> [JMIRx Med.](#) 2024 Jun 12;5:e45973. doi: 10.2196/45973.

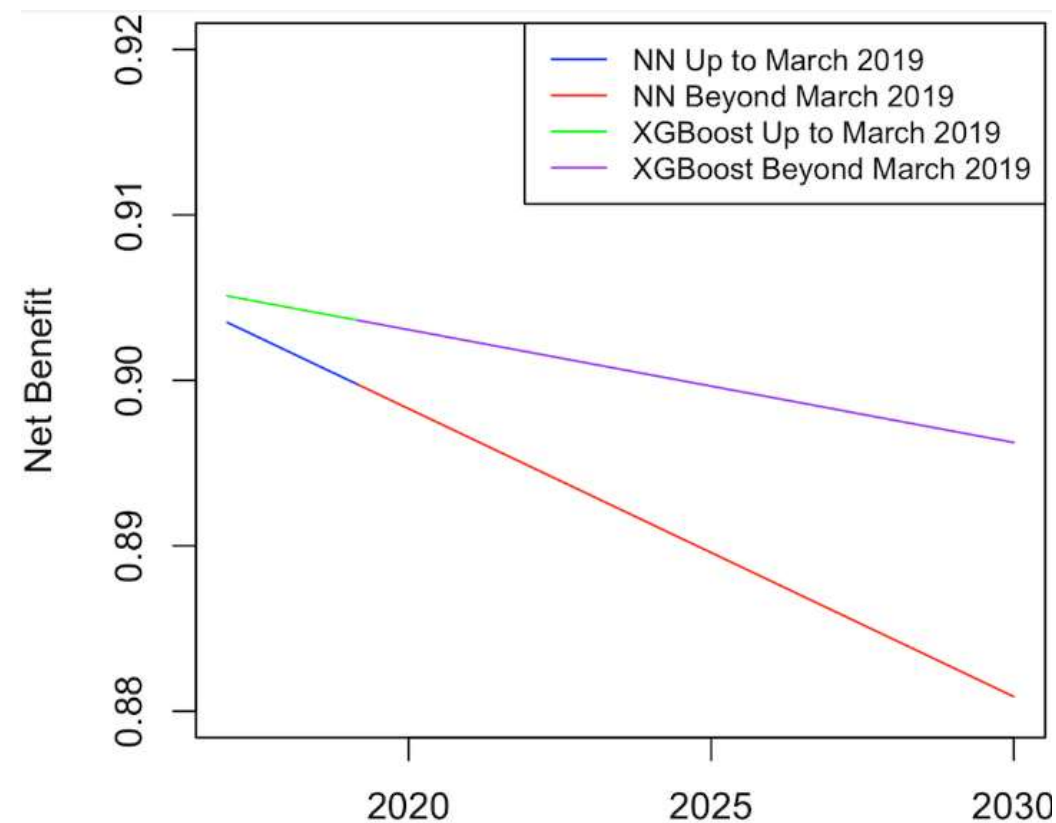
Performance Drift in Machine Learning Models for Cardiac Surgery Risk Prediction: Retrospective Analysis

Tim Dong ¹, Shubhra Sinha ¹, Ben Zhai ², Daniel Fudulu ¹, Jeremy Chan ¹, Pradeep Narayan ³, Andy Judge ¹, Massimo Caputo ¹, Arnaldo Dimagli ¹, Umberto Benedetto ¹, Gianni D Angelini ¹

Affiliations — collapse

Affiliations

- ¹ Bristol Heart Institute, Translational Health Sciences, University of Bristol, Bristol, United Kingdom.
- ² School of Computing Science, Northumbria University, Newcastle upon Tyne, United Kingdom.
- ³ Department of Cardiac Surgery, Rabindranath Tagore International Institute of Cardiac Sciences, West Bengal, India.





Explainability



Canadian Journal of Cardiology 38 (2022) 204–213

Review

Opening the Black Box: The Promise and Limitations of Explainable Machine Learning in Cardiology

Jeremy Petch, PhD, MA, BA(H)^{a,b,c,d}, Shuang Di, MSc, BSc^{a,c,e} and Walter Nelson, BSc(H)^{a,f,g}

^aCentre for Data Science and Digital Health, Hamilton Health Sciences, Hamilton, Ontario, Canada

^bInstitute of Health Policy, Management and Evaluation, University of Toronto, Toronto, Ontario, Canada

^cDivision of Cardiology, Department of Medicine, McMaster University, Hamilton, Ontario, Canada

^dPopulation Health Research Institute, Hamilton, Ontario, Canada

^eDalla Lana School of Public Health, University of Toronto, Toronto, Ontario, Canada

^fDepartment of Statistical Sciences, University of Toronto, Toronto, Ontario, Canada

ABSTRACT

Many clinicians remain wary of machine learning because of long-standing concerns about “black box” models. “Black box” is shorthand for models that are sufficiently complex that they are not straightforwardly interpretable to humans. Lack of interpretability in predictive models can undermine trust in those models, especially in health care, in which so many decisions are—literally—life and death issues. There has been a recent explosion of research in the field of explainable machine learning aimed at addressing these concerns. The promise of explainable machine learning is considerable, but it is

RÉSUMÉ

De nombreux cliniciens restent méfiants envers l'apprentissage automatique en raison de préoccupations de longue date concernant les modèles à « boîte noire ». Le terme « boîte noire » sert à désigner des modèles suffisamment complexes pour échapper à une interprétation simple par un humain. Le manque d'interprétabilité des modèles prédictifs peut miner la confiance en ces modèles, en particulier dans le domaine des soins de santé, où tant de décisions sont littéralement des questions de vie ou de mort. Il y a eu récemment une explosion de la recherche consacrée à l'apprentissage automatique explicable



| Patient Name | Probability of Heart Disease |
|--------------|------------------------------|
| | 94% |
| | 79% |
| | 55% |
| | 49% |
| | 31% |
| | 26% |

Probability of Heart Disease

Probability calculated: 7/21/2021 20:26

79%

High

Factors Contributing to Prediction

| | |
|-------------------------|----------------|
| Thallium Stress Test | Normal |
| Number of Major Vessels | 1 |
| Exercise Induced Angina | No |
| Chest Pain Type | Typical Angina |
| Max Heart Rate Achieved | 174 |
| Sex | Male |

Probability of Heart Disease

Probability calculated: 7/21/2021 20:26

26%

Low

Factors Contributing to Prediction

| | |
|-------------------------|--------|
| Thallium Stress Test | Normal |
| Max Heart Rate Achieved | 131 |
| Number of Major Vessels | 1 |
| Sex | Male |
| ST Depression | 0.10 |
| Age | 69 |



Usability



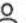

Cardiovascular Digital Health Journal

Volume 4, Issue 3, June 2023, Pages 101-110



Original Article

Artificial intelligence-enabled tools in cardiovascular medicine: A survey of current use, perceptions, and challenges

Alexander Schepart PharmD, MBA ^{*}, Arianna Burton PharmD ^{*}, Larry Durkin MBA [†],
Allison Fuller BA [†], Ellyn Charap MSc [†], Rahul Bhambri PharmD, MBA ^{*},
Faraz S. Ahmad MD, MS [‡]  

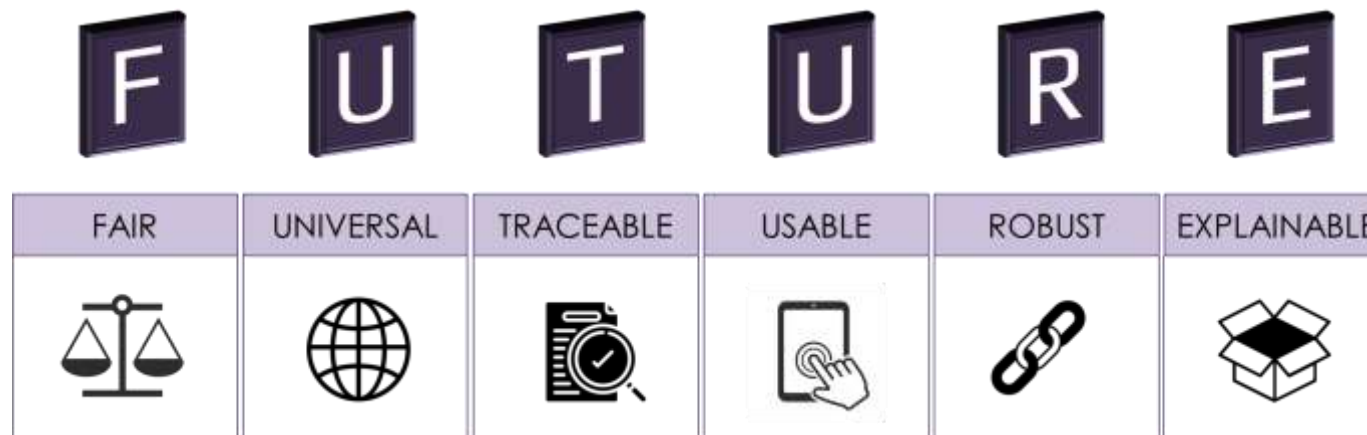
“I think it has to pull from the [health] record without you doing anything...you don’t have to go and put the data into it yourself.” –Cardiologist

“Sometimes when you get all of these risk calculators...and EPIC, it becomes more information, which is not better. We’re pretty busy as clinicians. I think information that’s not actionable or helpful, it just slows us down and we ignore it.” –Cardiologist

“I don’t really know what’s going on as far as generating risk scores for something like [ATTR-CM]. The majority of the time, we look at the detailed data graph and use trends...I know there are scoring systems for drug withdraw though.” –Cardiologist

Characteristics Trustworthy AI

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





Characteristics Trustworthy AI

| | Clusters of requirements | Core principle |
|---|---|------------------------|
| 1 | Diversity, Inclusiveness, Non-discrimination, Bias, Equity | <u>F</u> airness |
| 2 | Generalisability, Adaptability, Interoperability, Applicability | <u>U</u> niversality |
| 3 | Transparency, Monitoring, Auditing, Accountability | <u>T</u> raceability |
| 4 | Human-centred AI, User engagement, Accessibility, Efficiency | <u>U</u> sability |
| 5 | Reliability, Resilience, Safety, Security | <u>R</u> obustness |
| 6 | Interpretability, Understandability, Transparency | <u>E</u> xplainability |

Characteristics Trustworthy AI

Based on ethical principles and fundamental rights:

FAIRNESS



Right to non-discrimination

UNIVERSALITY



Right to equity

TRACEABILITY



Right to accountability

USABILITY



Right to autonomy

ROBUSTNESS



Right to safety

EXPLAINABILITY



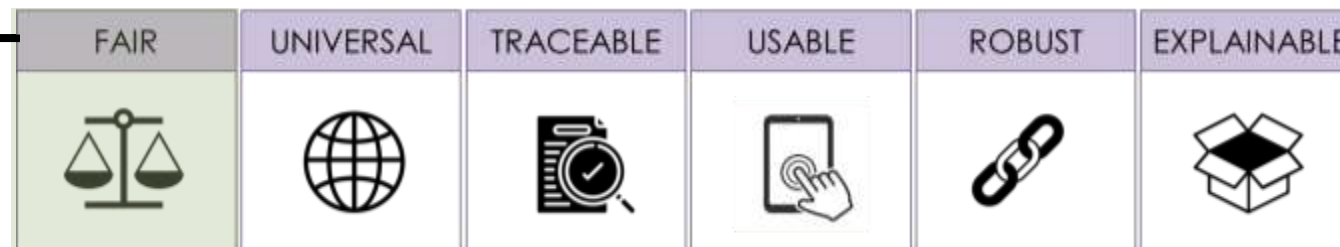
Right to transparency



Part 1 - What is trustworthy AI?

Part 2 - How do we achieve it?

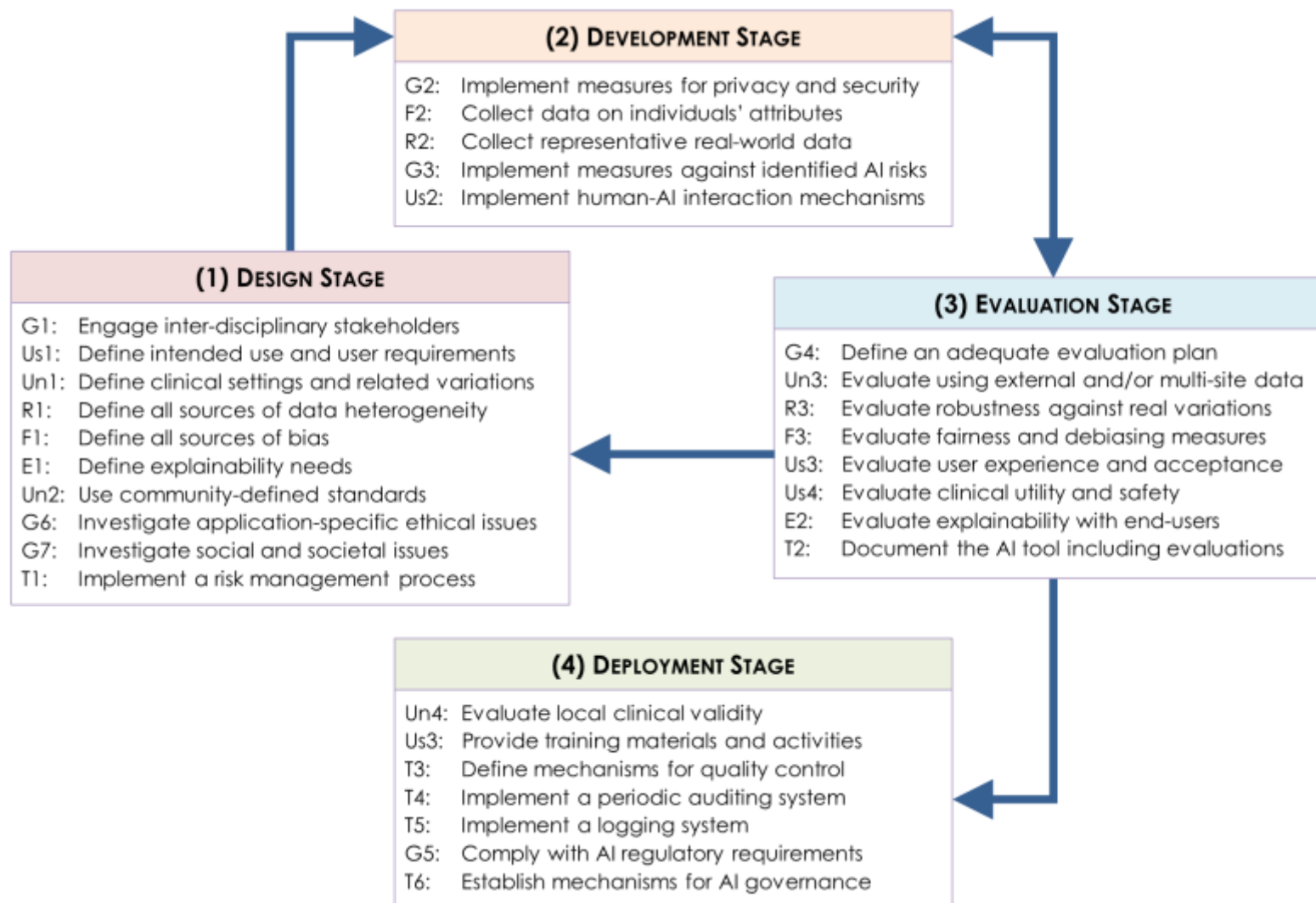
Operationalisation



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

Operationalisation





Best Practices – Design

| | | |
|------|---|------|
| 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 (What)

Practical steps (How)

Examples (References)

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



Stakeholder Engagement





This is the board created for the pilot site **Centro De Investigaciones Tecnológicas, Biomedicas Y Medioambientales (CITBM)**

PART I

| | Phase I: Pre-diagnosis | Phase II: Diagnosis | Phase III: Management | Phase IV: Hospitalization |
|--|---|--|--|---|
| Patient Journey of a patient with Heart Failure | Primary care Entry point through primary care | Secondary Care Entry point through medical specialist | Tertiary Care Entry point through hospitalization | Diagnosis What happens in diagnosis? |
| Steps What does a story of a patient's journey look like? | Patients start their evaluation at primary center (step 1) (General Medicine / Internal medicine not exam allowed), then referred to secondary (general hospital cardiologist) according to step 2 evaluation patient is referred to step 3 (tertiary center) | At step 2 (General cardiologist) first specific exam (not invasive) Specific HF treatment depending on local resources and initial stratification . | Heart Failure cardiologist- invasive evaluation studies- MRI-CT-TEE. Patient is stratified according their etiology and risk | step1- (General Medicine / Internal medicine not exam allowed) Specific HF treatment depending on local resources and initial stratification . Invasive evaluation studies- MRI-CT-TEE. Patient is stratified according their etiology and risk followed medical treatment (device- transplant/ICD-surgery) |
| People involved What does all the work that are involved in this step. If there is more than one clinician involved please specify the type of clinician. | General medicine Internal medicine | General cardiologist | Heart failure cardiologist | Cardiologist |
| Where does this step take place? What does all physical locations where this step can take place. E.g. healthcare center, hospital, GP clinic, etc. | [Write down all physical locations where this step can take place. E.g. healthcare center, hospital, GP clinic, etc.] | [Write down all physical locations where this step can take place. E.g. healthcare center, hospital, GP clinic, etc.] | [Write down all physical locations where this step can take place. E.g. healthcare center, hospital, GP clinic, etc.] | [Write down all physical locations where this step can take place. E.g. healthcare center, hospital, GP clinic, etc.] |
| Time duration How much time does this step usually takes until the patient moves to the next step? | [How much time does this step usually takes until the patient moves to the next step?] | [How much time does this step usually takes until the patient moves to the next step?] | [How much time does this step usually takes until the patient moves to the next step?] | [How much time does this step usually takes until the patient moves to the next step?] |
| Use of Artificial intelligence What are the uses of the artificial intelligence in this step? Please mention which tool you currently use, have used in the past, or plan you want to use. | [Use of Artificial intelligence in this step] | [Use of Artificial intelligence in this step] | [Use of Artificial intelligence in this step] | [Use of Artificial intelligence in this step] |

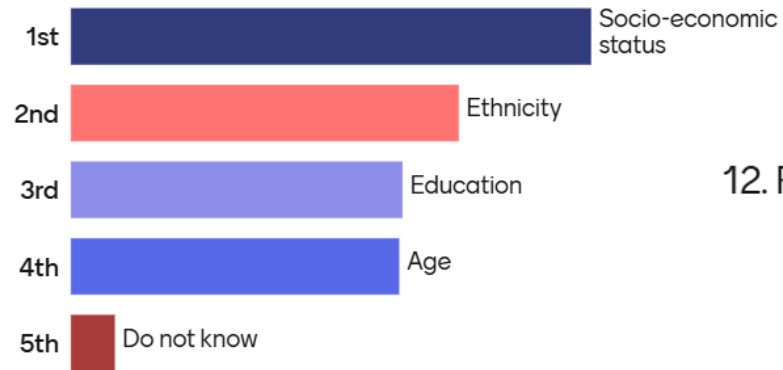
Part I: Patient Journey: Heart Failure

Part II: AI4HF in clinical practice: when, what, how?

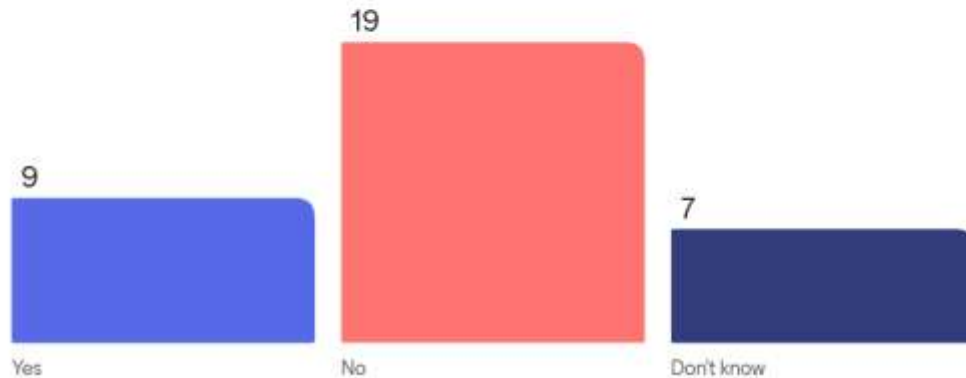


Stakeholder Engagement

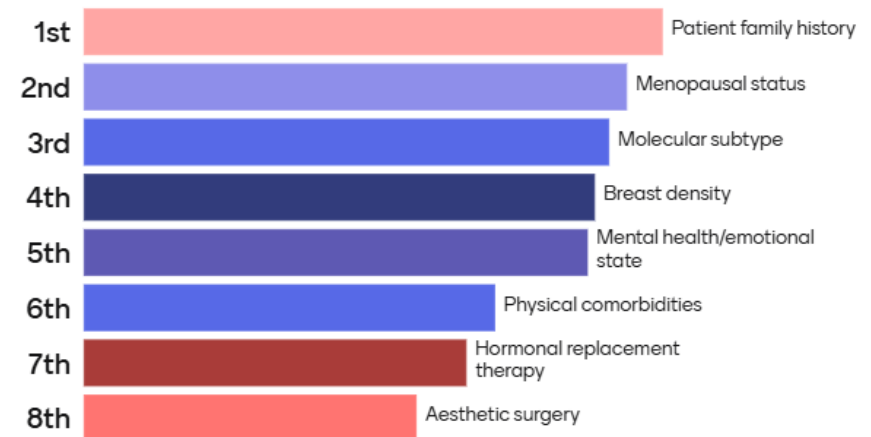
8. Rank the following variables based on their importance for bias estimation



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



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





Best Practices – Design

- | | | |
|------|---|------|
| 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



AI4HF

Secondary risk prevention in heart failure

Clinicians

What should the AI tool predict ?

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

Patients

What should the AI tool predict ?

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



Sources of Biases



AI4HF



Sources of Biases



Japanese



Quechua



European



Mestizo



Sources of Biases



Lima: Sea level



Arequipa: 2,335 m

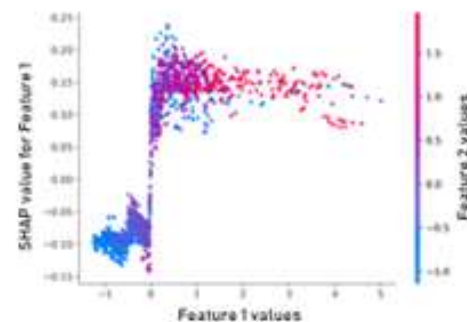
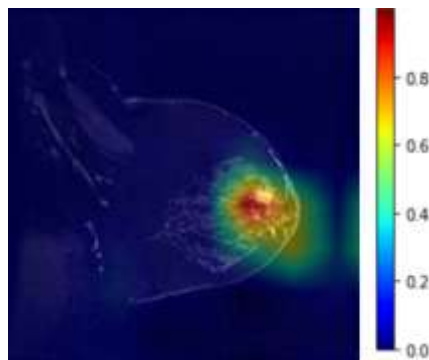


Cusco: 3,400 m

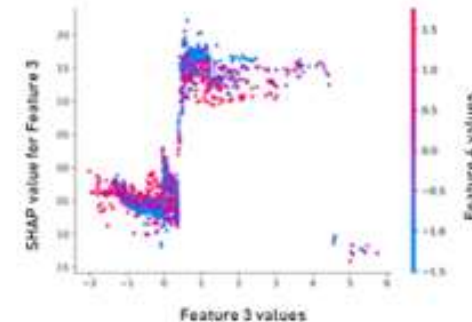


Rinconada: 5,100 m

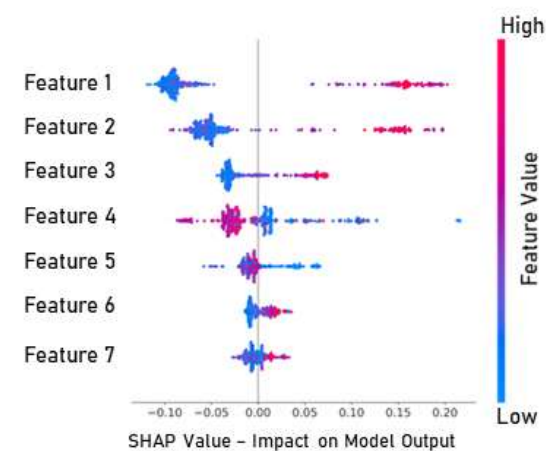
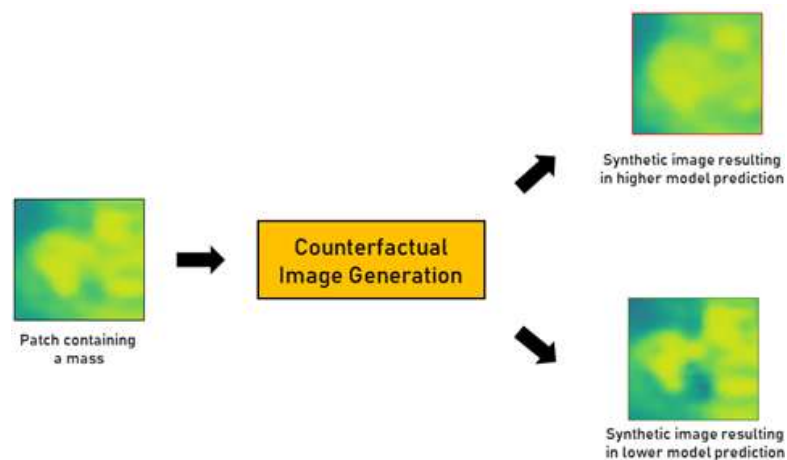
Explainability Options



(a)



(b)





Best Practices – Development

| | | |
|------|--|------|
| 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) |



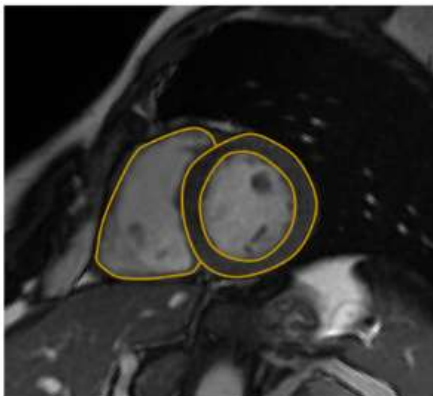
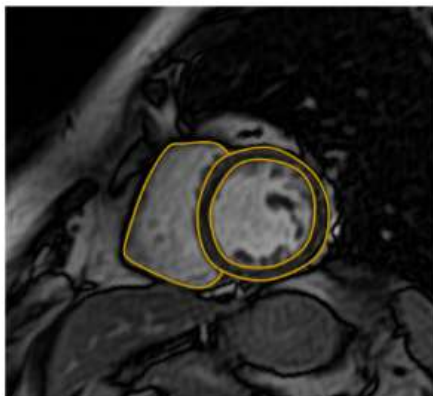
Data Collection

| Patient ▼ | Age ▼ | Sex ▼ | Ethnicity ▼ | Neighbourhood ▼ | Altitude ▼ | Skin colour ▼ | Education ▼ | Menopause ▼ |
|------------|-------|--------|-------------|-----------------|------------|---------------|-------------|-------------|
| Patient001 | 20 | Male | E. Europe | Gracia | 0 | White | High-School | No |
| Patient002 | 20 | Male | S. Europe | Gracia | 0 | White | University | No |
| Patient003 | 30 | Male | N. Africa | Horta | 500 | White | University | No |
| Patient004 | 30 | Male | N. Africa | Horta | 500 | White | University | No |
| Patient005 | 40 | Male | S. Africa | Poblenou | 1000 | Black | High-School | No |
| Patient006 | 40 | Female | S. Africa | Poblenou | 1000 | Black | High-School | Yes |
| Patient007 | 50 | Female | S. Asia | Gervasi | 1500 | White | High-School | Yes |
| Patient008 | 50 | Female | E. Asia | Gervasi | 1500 | White | University | Yes |
| Patient009 | 60 | Female | L. America | Eixample | 3000 | White | University | No |
| Patient010 | 60 | Female | L. America | Eixample | 3000 | Black | University | No |

Data Representativeness

Multi-Centre, Multi-Vendor and Multi-Disease Cardiac Segmentation: The M&Ms Challenge

Victor M. Campello¹, Polyxeni Gkontra², Cristian Izquierdo, Carlos Martín-Isla, Alireza Sojoudi, Peter M. Full³, Klaus Maier-Hein, Yao Zhang⁴, Zhiqiang He, Jun Ma⁵, Mario Parreño⁶, Alberto Albiol⁷, Fanwei Kong, Shawn C. Shadden⁸, Jorge Corral Acero⁹, Vaanathi Sundaresan¹⁰, Mina Saber, Mustafa Elattar¹¹, Hongwei Li¹², Bjoern Menze, Firas Khader, Christoph Hauburger, Cian M. Scannell¹³, Mitko Veta¹⁴, Adam Carscadden, Kumaradevan Punithakumar¹⁵, *Senior Member, IEEE*, Xiao Liu, Sotirios A. Tsaftaris¹⁶, Xiaoqiong Huang, Xin Yang¹⁷, Lei Li, Xiahai Zhuang¹⁸, David Viladés¹⁹, Martín L. Descalzo²⁰, Andrea Guala²¹, Lucia La Mura²², Matthias G. Friedrich, Ria Garg²³, Julie Lebel, Filipe Henriques, Mahir Karakas, Ersin Çavuş, Steffen E. Petersen²⁴, Sergio Escalera²⁵, Santi Seguí²⁶, José F. Rodríguez-Palomares²⁷, and Karim Lekadir²⁸



| Centre | Vendor | Model | Field strength (T) |
|--------|---------|-----------------|--------------------|
| 1 | Siemens | MAGNETOM Avanto | 1.5 |
| 2 | Philips | Achieva | 1.5 |
| 3 | Philips | Achieva | 1.5 |
| 4 | GE | Signa Excite | 1.5 |
| 5 | Canon | Vantage Orian | 1.5 |
| 6 | Siemens | MAGNETOM Skyra | 3.0 |



Best Practices – Validation

| | | |
|------|--|------|
| 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



AI4HF



Spain



Netherlands



Amsterdam UMC
University Medical Centers



UMC Utrecht



Czech Republic



ST. ANNE'S UNIVERSITY HOSPITAL BRNO
INTERNATIONAL CLINICAL RESEARCH CENTER



Peru



**Instituto Nacional
Cardiovascular**



Tanzania



**Muhimbili University of Health and
Allied Sciences**



Usability Evaluation

Human evaluators in 5 sites:

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

The System
Usability Scale (SUS)
helps measure:



Efficiency:

How fast someone can use a product



Intuitiveness:

How effortlessly someone can understand a product



Ease:

How simple a product is to use



Satisfaction:

How much a user subjectively likes or dislikes using a product



Explainability Evaluation



AI Explainability Score:

Questions to assess explainability with clinicians:

- Did you find the AI explanations clear?
- Did you find the AI consistent between cases?
- Did the AI explanations increase confidence in the decisions?
- Were the AI visualisations easy to use?

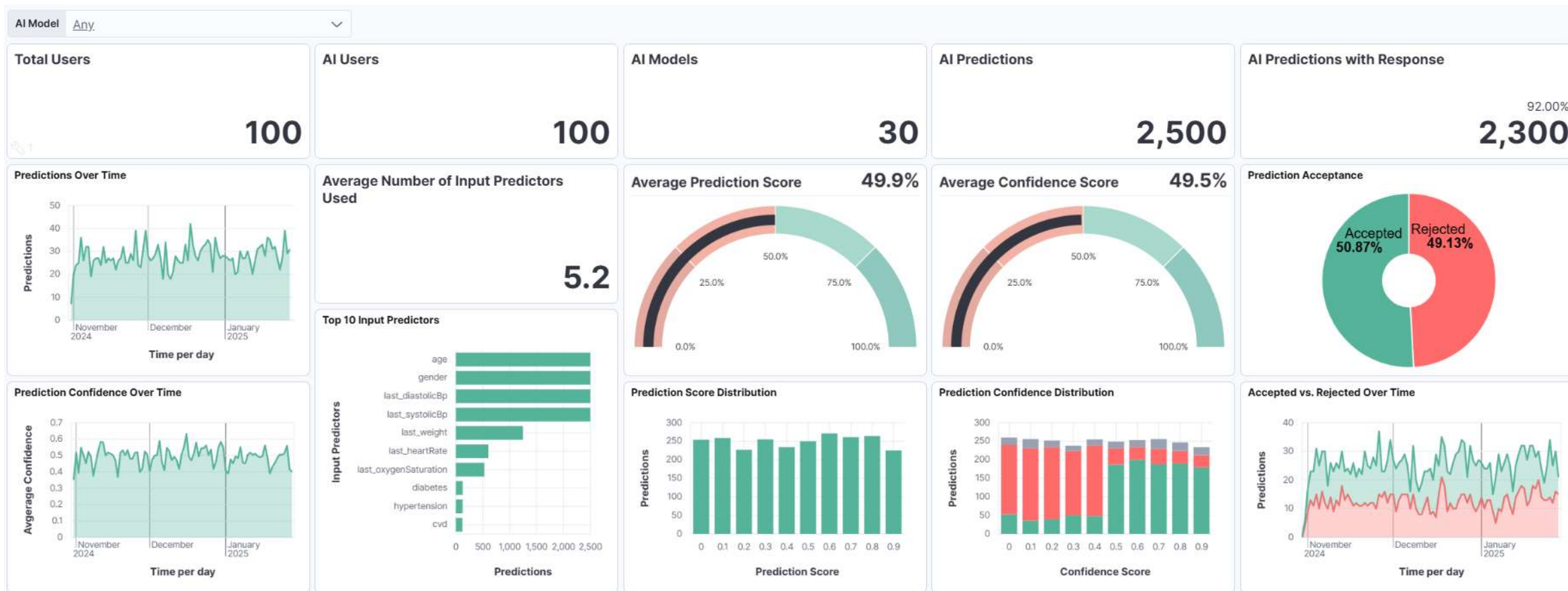


Best Practices – Deployment

| | | |
|------|---|------|
| 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 AI governance | (30) |



Periodic Auditing



AI4HF

RESEARCH METHODS AND REPORTING

OPEN ACCESS

Check for updates

FUTURE-AI: international consensus guideline for trustworthy and deployable artificial intelligence in healthcare

Karim Lekadir,^{1,2} Alejandro F Frangi,^{3,4} Antonio R Porras,⁵ Ben Glocker,⁶ Celia Cintas,⁷ Curtis P Langlotz,⁸ Eva Weicken,⁹ Folkert W Asselbergs,^{10,11} Fred Prior,¹² Gary S Collins,¹³ Georgios Kaltsis,¹⁴ Gianna Tsakou,¹⁵ Irène Buvat,¹⁶ Jayashree Kalpathy-Cramer,¹⁷ John Mongan,¹⁸ Julia A Schnabel,¹⁹ Kaisar Kushiab,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-Bonmati,²⁶ M Jorge Cardoso,²⁷ Maciej Bobowicz,²⁸ Mahsa Shabani,²⁹ Manolis Tsiknakis,³¹ Maria A Zuluaga,³⁰ Marie-Christine Fritzsche,³¹ Marina Camacho,¹ Marius George Lingurar,³² Markus Wenzel,⁹ Marleen De Bruijne,³³ Martin G Tolsgaard,³⁴ Melanie Golsauf,³⁵ Mónica Cano Abadía,³⁵ Nikolaos Papanikolaou,²⁸ Noussair Lazrak,¹ Oriol Pujol,¹ Richard Osuala,¹ Sandy Napel,³⁷ Sara Colantonio,³⁸ Smriti Joshi,¹ Stefan Klein,³⁹ Susanna Aussó,³⁹ Wendy A Rogers,⁴⁰ Zohaib Salahuddin,⁴¹ Martijn P A Starmans³³; on behalf of the FUTURE-AI Consortium

For numbered affiliations see end of the article.
Correspondence to: K Lekadir (karim.lekadir@lup.edu; ORCID: 0000-0002-9456-1612)
Additional material is published online only. To view please visit the journal online.
Cite this as: *BMJ* 2025;388:e081954
http://dx.doi.org/10.1136/bmj.2024-081954
Accepted: 10 January 2025

Despite major advances in artificial intelligence (AI) research for healthcare, the deployment and adoption of AI technologies remain limited in clinical practice. This paper 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 AI scientists, clinical researchers, biomedical ethicists, and social scientists. Over a two year period, the FUTURE-AI guideline was

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 AI, from design, development, and validation to regulation, deployment, and monitoring.

Introduction

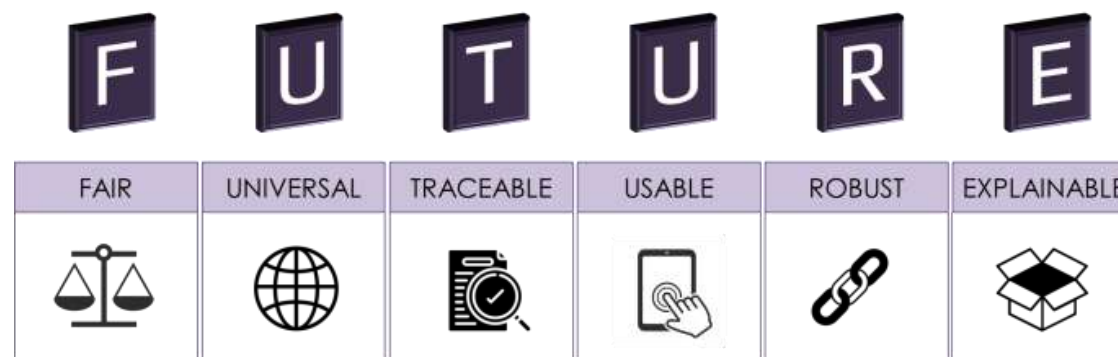
In the field of healthcare, artificial intelligence (AI)—that is, algorithms with the ability to self-learn logic—and data interactions have been increasingly used to develop computer aided models, for example, disease diagnosis, prognosis, prediction of therapy response or survival, and patient stratification.¹ 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 AI.^{2,3} In particular, existing research has shown that AI tools in healthcare can be prone to errors and patient harm, biases and increased health inequalities, lack of transparency and accountability, as well as data privacy and security breaches.^{4,5}

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,



SUMMARY POINTS

Despite major advances in medical artificial intelligence (AI) 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 AI.
The FUTURE-AI framework is structured around six guiding principles: fairness, universality, traceability, usability, robustness, and explainability.
The guideline addresses the entire AI 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 AI systems, emphasising multistakeholder collaboration.
Continuous risk assessment and mitigation are fundamental, addressing biases, data variations, and evolving challenges during the AI lifecycle.
FUTURE-AI is designed as a dynamic framework, which will evolve with technological advancements and stakeholder feedback.





Next Steps

- Feedback gathering from over 10 EU projects
- Feedback gathering from external actors (e.g. ESC members)
- Next papers:
 - FUTURE-AI guideline: The patient perspective
 - FUTURE-AI guideline: Clinical validation methods
 - FUTURE-AI guideline: Implications for AI regulations
 - FUTURE-AI guideline: Adaptations for large language models



Join Us!

[FUTURE-AI BEST PRACTICES](#) ▾[NEWS AND EVENTS](#)[CURRENT PROJECTS](#) ▾[CONSORTIUM](#)[CONTACT US](#)[JOIN US](#)

FUTURE-AI: Best practices for trustworthy AI in medicine

FUTURE-AI is an international, multi-stakeholder initiative for defining and maintaining concrete guidelines that will facilitate the design, development, validation and deployment of trustworthy AI solutions in medicine and healthcare based on six guiding principles: Fairness, Universality, Traceability, Usability, Robustness and Explainability.

karim.lekadir@ub.edu



Many Thanks!



karim.lekadir@ub.edu