

# IRCAI

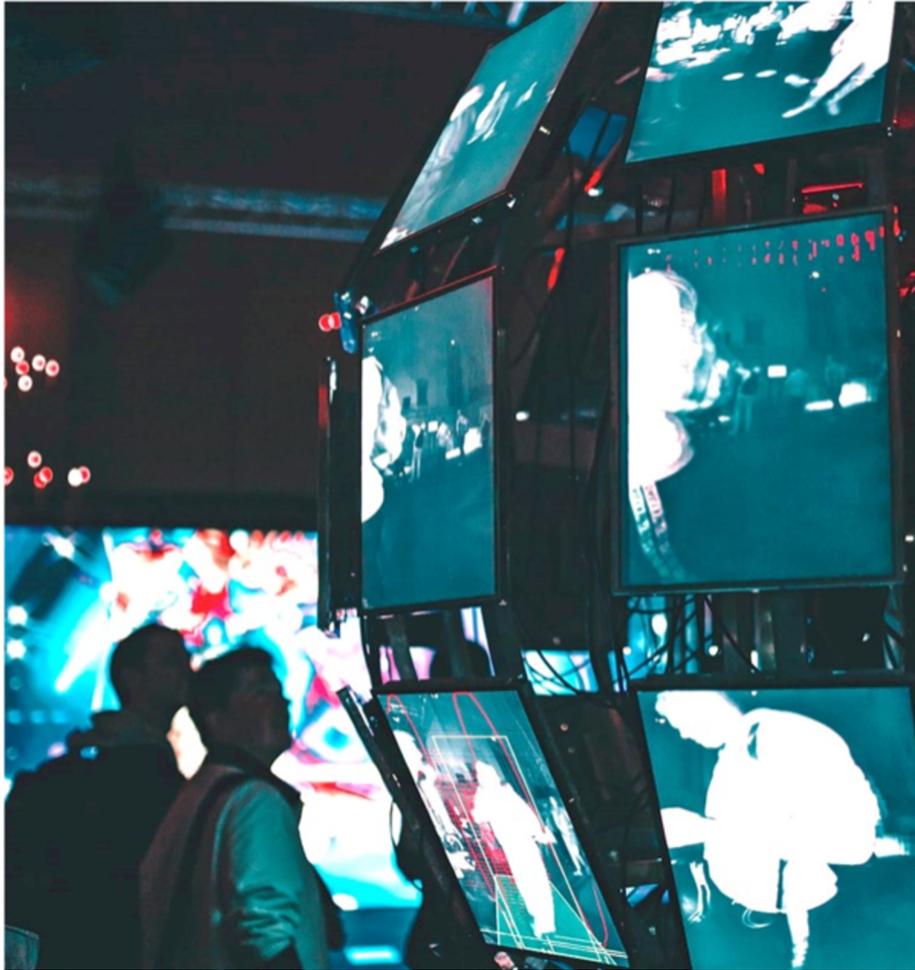
Scientific excellence, global public dialogue and technology for sustainability, inclusion and equality.

## Responsible AIoT and TinyML

School on Applied AI for Sustainable Development | (smr 4210)  
Trieste 24.3.2025



Co-funded by  
the European Union



Potential for new and old threats being enhanced

- Citizen control as envisaged in Brave New World
- Automatic processing of decisions that ignore marginalised groups
- Automatic weapons that can be used to disable legitimate protest or wage destructive wars

## GLOBAL CONFERENCE ON AI AND HUMAN RIGHTS

13 and 14 June 2024  
Faculty of Law,  
University of Ljubljana (Slovenia)



[www.ai-right-to-life.si](http://www.ai-right-to-life.si)



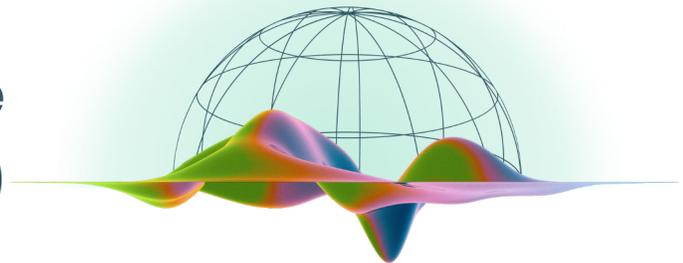
**unesco**



## Recommendation on the Ethics of AI

- After extensive consultations, recommendation adopted at the UNESCO assembly in November 2021
- Human centered: how can individual humans be protected from misuse and enriched by the use of AI
- Closely aligned with European humanistic values and with EU initiatives as well as Council of Europe

## Global Forum on the Ethics of AI (GFEAI)



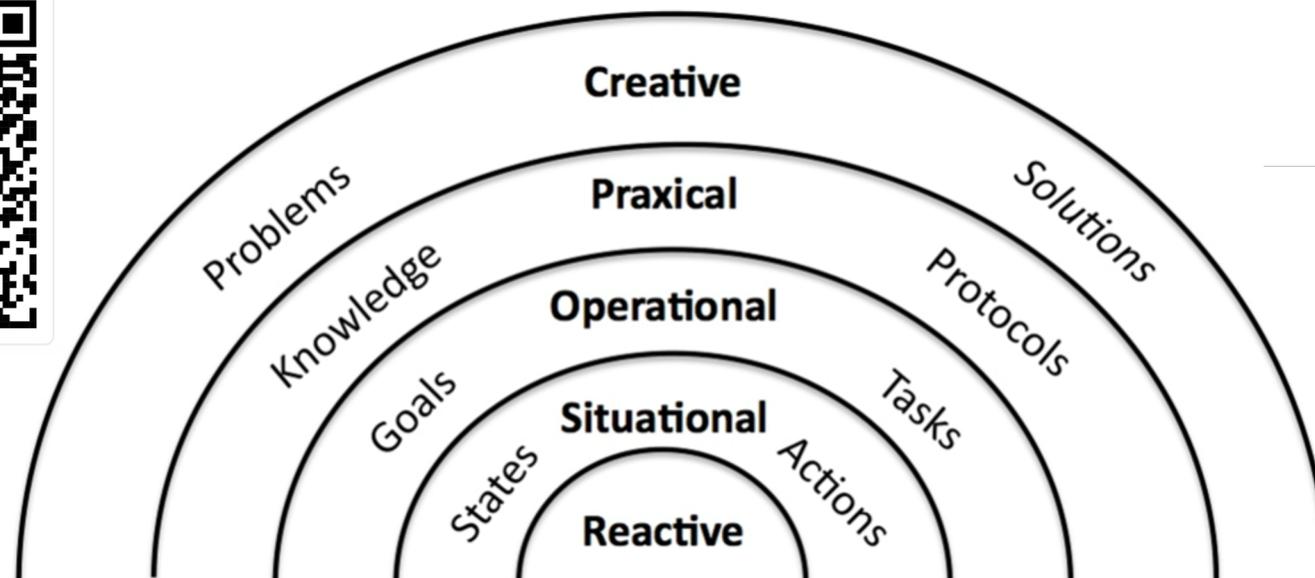


- EU AI Act is first attempt, but risk based so that many systems not covered
- One advantage of the UNESCO recommendation is that it is vague: forcing designers to think through unforeseen consequences:
  - assessing when human rights have been compromised, system acted unfairly, introduced bias, etc.
  - Any formal definition will inevitably be limited and motivate circumvention





- Human-centric AI is at the centre of the European vision of a positive role for AI
- AI that empowers humans to be more effective, more creative, more understanding
- Human-centric AI (HCAI) has been the focus of the Humane AI Network of Excellence
- How does its research agenda envisage HCAI?



Common Ground through Explanation, Instruction, Demonstration, Experience

Humane AI Research agenda highlights the ingredients of Collaborative Intelligent Systems:

- Need to find ‘common ground’ across a range of levels in order to enable effective cooperation/communication
- Levels identified roughly correspond to different styles of collaboration with collaborative systems potentially involving more than one level

HUMANE  AI NET

**HumanE AI Net:  
The HumanE AI Network**

Grant Agreement Number: 952026  
Project Acronym: HumanE AI Net

Project Dates: 2020-09-01 to 2023-08-31  
Project Duration: 36 months

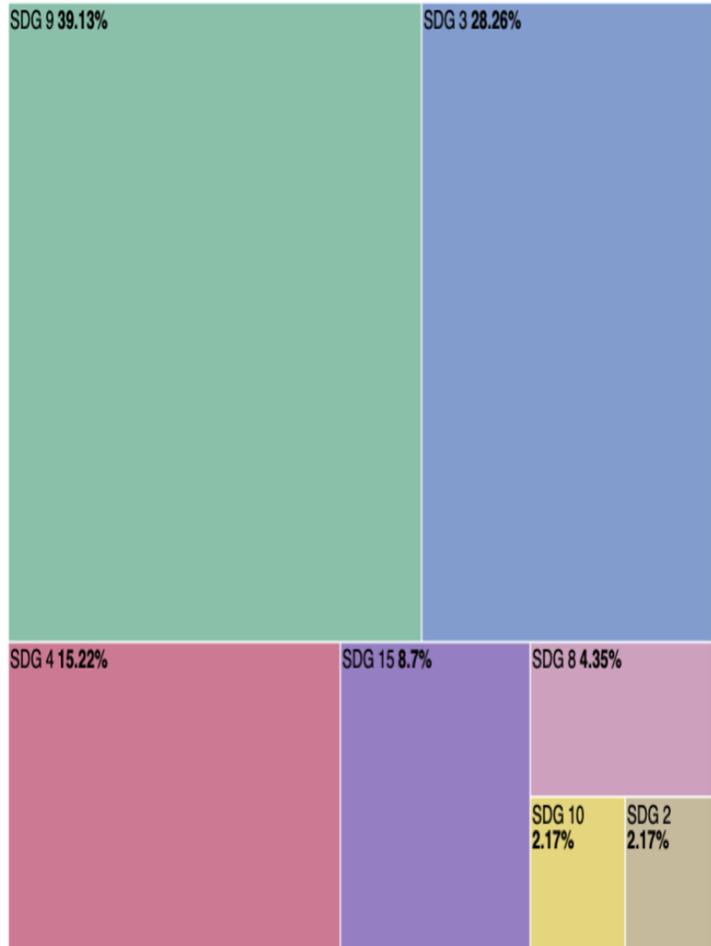
**D6.1 Strategic Research Agenda**

Author(s): Paul Lukowicz  
Contributing partners: John Shawe-Taylor, James Crowley, Antti Oulasvirta, Virginia Dignum, George Kampis.  
Date: Mai 10, 2022  
Approved by: Paul Lukowicz  
Type: Report @  
Status: final  
Contact: [Paul.Lukowicz@dfki.de](mailto:Paul.Lukowicz@dfki.de)

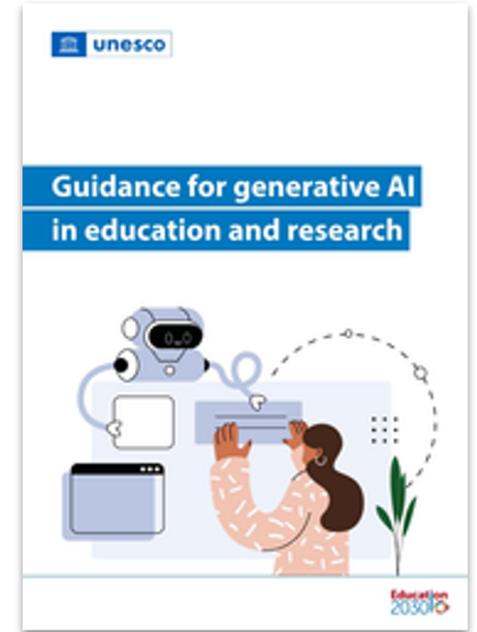
Dissemination Level  
PU | Public

Copyright - This document has been produced under the EC Horizon2020 Grant Agreement 952026-HumanE-AI-Net. This document and its contents remain the property of the beneficiaries of the HumanE AI Net Consortium.

OPERATIONALIZED ETHICAL AND FUNCTIONAL DIMENSIONS WITH PRIORITY AREAS FOR TINYML AND STEAM EDUCATION



Dimension	Description	EdgeAI/TinyML	STEAM Education
<b>Transparency</b>	Clearly communicate AI operations, data handling, and inference logic to participants to build trust and accountability in both community science and educational contexts.	High – enables interpretability of local models and data flows.	High – critical for teaching explainability and responsible AI use.
<b>Inclusivity</b>	Ensure equitable access to AI tools, datasets, and learning resources, addressing disparities in infrastructure, gender representation, and cultural participation.	Medium – supports broader participation in device deployment.	High – foundational for equitable AI literacy and engagement.
<b>Fairness</b>	Identify and mitigate algorithmic bias across data collection, labeling, and model evaluation processes, especially in community or low-resource settings.	High – essential for localized and context-aware model training.	High – core concept for teaching ethical reasoning in AI and data science.
<b>Environmental Responsibility</b>	Assess and reduce lifecycle impacts of AI hardware, including sourcing, energy consumption, and end-of-life management.	High – TinyML’s edge devices directly affect ecological sustainability.	Medium – integrated into STEAM through sustainability modules and project design.
<b>Cultural Diversity</b>	Incorporate local knowledge, indigenous epistemologies, and linguistic diversity into AI datasets, interpretations, and learning activities.	Medium – supports culturally relevant sensor deployment and data interpretation.	High – encourages interdisciplinary, culturally responsive pedagogy.
<b>Operational Assistance</b>	TinyML devices perform localized inference for environmental monitoring, health tracking, or IoT sensing, reducing reliance on cloud infrastructure.	High – primary operational focus of edge ML systems.	Medium – serves as an applied learning example in engineering and computing courses.
<b>Decision Support</b>	Generative AI systems synthesize data insights and summaries to assist collective decision-making and policy formation in citizen science projects.	Medium – complements TinyML data with higher-level synthesis.	High – enhances inquiry-based learning, problem-solving, and critical reflection.
<b>Interpretive / Narrative Collaboration</b>	AI systems co-generate visualizations, reports, and educational materials with human collaborators, promoting reflective learning and knowledge co-creation.	Medium – relevant in projects using embedded visualization or reporting tools.	High – central to creative and narrative integration in STEAM curricula.



One way to try to offset the potential misuse is to regulate greater disclosure

- AI should be educational, broadening our understanding of situations and content
- This could be a requirement that enables users to make up their own minds
- Many of the manipulations rely on suggesting that full disclosure is not available



**Informed Consent:** ensure participants know why their data is being collected, how it will be used, and any associated risks.

**Privacy:** Protect the identities of participants by anonymizing or de-identifying data.

**Transparency:** clear purpose of the data collection, the methods used, and potential benefits and risks.

**Data Sharing:** participants should be aware and have given consent.

**Vulnerable Populations:** extra precautions when collecting data from children, the elderly, and other vulnerable groups.

**Potential Harm:** Assess and minimize potential risks to participants. This includes emotional distress, financial harm, or other adverse effects.

**Purpose Change:** seek fresh consent from participants.

**Limitation:** avoid over-collection.

**Accuracy:** data is accurate and represents the truth.

**Beneficence:** data collection should be for good.

**Models vs. Data:** storing models, forgetting the data

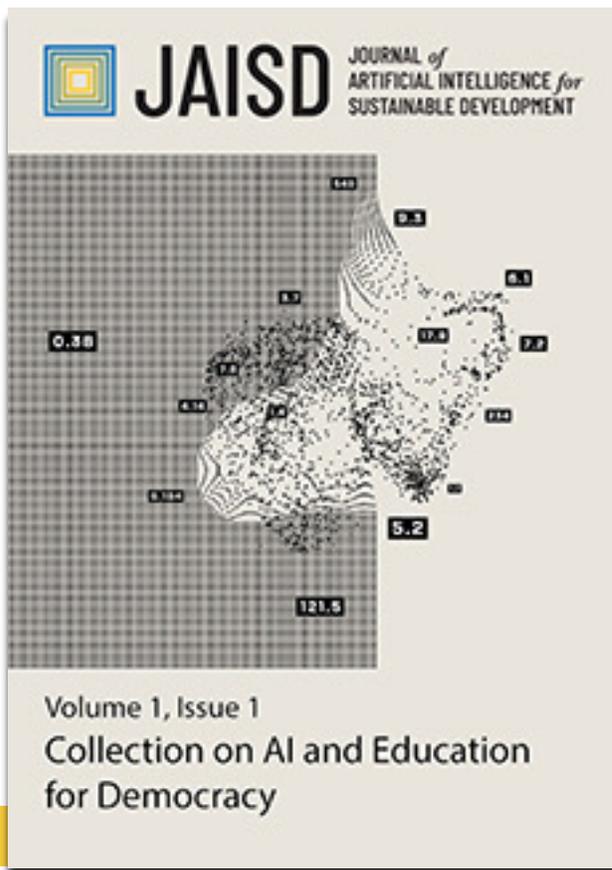
**Access:** who has access to the data.

**Long-term Storage:** ensure secure storage methods.

**Feedback:** offer participants feedback or results from the study if they express interest.

**Regulation and Legislation:** adhere to local, national, and international data protection laws and regulations.

**Cultural Sensitivity:** understand and respect cultural norms and traditions when collecting data



**Transparency:** AI algorithms should be transparent in terms of how they work and how decisions are made – XAI.

**Bias and Fairness:** It's crucial to identify, reduce, and disclose biases to ensure fairness in decision-making.

**Accountability:** determine who is responsible for AI decisions.

**Privacy:** protect the privacy of individuals when processing their data.

**Data Security:** safeguard this data against breaches and unauthorized access.

**Autonomy:** People should always have a choice, especially concerning decisions that significantly affect their lives.

**Informed Consent:** Users should be informed when their data is being processed by AI and have a clear understanding of how the AI system uses their data.

**Explainability:** AI decisions should be interpretable and explainable.

**Generalizability and Robustness:** AI models are not overfitted to their training data and can generalize well to real-world situations.

**Economic and Social Impacts:** Consider the broader societal and economic effects of AI for good

**Long-term Considerations:** longer-term implications of AI, including the potential for systems to evolve or be used in ways not originally intended.

**Environmental Impact:** consider the environmental footprint and strive for sustainable AI research.

**Continual Monitoring:** AI systems, especially those deployed in dynamic environments, should be continuously monitored to ensure they are behaving as expected.

**Stakeholder Participation:** Involve relevant stakeholders, including those affected by AI systems, in the design, development, and deployment processes.

**Regulation and Standards:** Align AI practices with existing regulations and contribute to the development of ethical standards for AI.



# Continuous Compliance with Legislations for Privacy, IRP & Data Collection

## AI Ethics & Principles by Design Framework

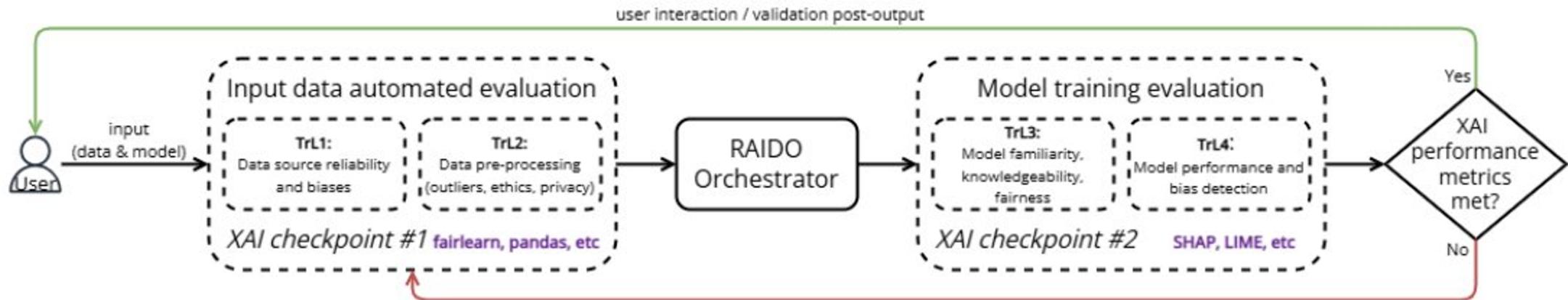
## Reinforced Benchmarking & Feedback-based Progress Monitoring

Acceptance of AI use in life and job/industry		
20	I am interested in exploring new technological developments for my field of work/in my daily life.	<ul style="list-style-type: none"> <li>Strongly Disagree</li> <li>Disagree</li> <li>Neutral/Neither</li> <li>Agree nor Disagree</li> <li>Agree</li> <li>Strongly Agree</li> </ul>
21	I could imagine including AI in my current work or in my daily life.	<ul style="list-style-type: none"> <li>Strongly Disagree</li> <li>Disagree</li> <li>Neutral/Neither</li> <li>Agree nor Disagree</li> <li>Agree</li> <li>Strongly Agree</li> </ul>
22	I see the benefit of using AI tools in my current work or in my daily life.	<ul style="list-style-type: none"> <li>Strongly Disagree</li> <li>Disagree</li> <li>Neutral/Neither</li> <li>Agree nor Disagree</li> <li>Agree</li> <li>Strongly Agree</li> </ul>
23	Polymakers strongly support AI.	<ul style="list-style-type: none"> <li>Strongly Disagree</li> <li>Disagree</li> <li>Neutral/Neither</li> <li>Agree nor Disagree</li> <li>Agree</li> <li>Strongly Agree</li> </ul>
24	(For respondents in employment) The use of AI in my field aligns with what society considers appropriate for the industry.	<ul style="list-style-type: none"> <li>Strongly Disagree</li> <li>Disagree</li> <li>Neutral/Neither</li> <li>Agree nor Disagree</li> <li>Agree</li> <li>Strongly Agree</li> </ul>
25	(For respondents in employment) I am the person who decides on the use of AI at my job. (For people in employment)	<ul style="list-style-type: none"> <li>Strongly Disagree</li> <li>Disagree</li> <li>Neutral/Neither</li> <li>Agree nor Disagree</li> <li>Agree</li> <li>Strongly Agree</li> </ul>
26	(For respondents in employment) I feel that my organisation would support me in the adoption of AI. (New question for people in employment).	<ul style="list-style-type: none"> <li>Strongly Disagree</li> <li>Disagree</li> <li>Neutral/Neither</li> <li>Agree nor Disagree</li> <li>Agree</li> </ul>



VANESSA NUROCK

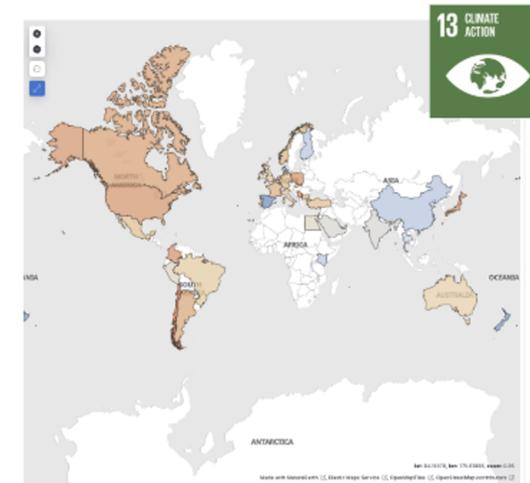
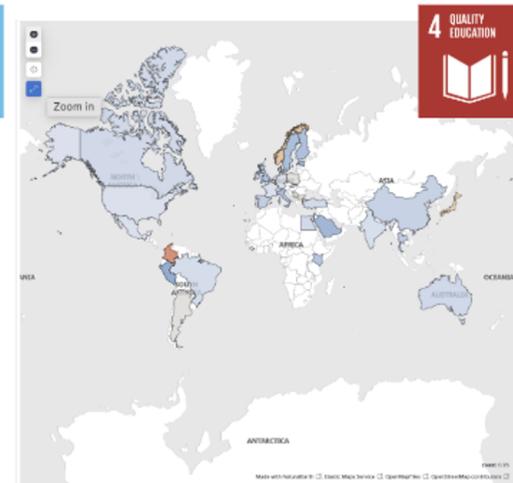
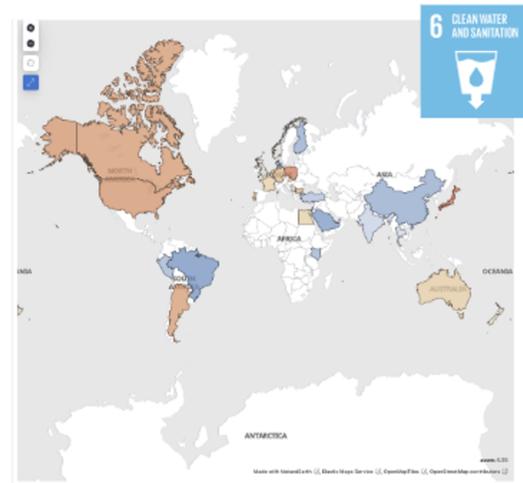
Professor of Philosophy at the Université de Côte d'Azur (France) and UNESCO EVA Chair (Ethique du Vivant et de l'Artificiel / Ethics of the Living and the Artificial)



## Framework for AI Explainability, Transparency & Trustworthiness

- Compute sentiment analysis (using VADER) on title + text of OECD AI policies translated
- Compare geographically the neutral sentiment on the AI policies to identify tendencies
- Differentiate between SDGs to explore the variability across the 17 targets and their topics

**Sentiment Bias Aware.**

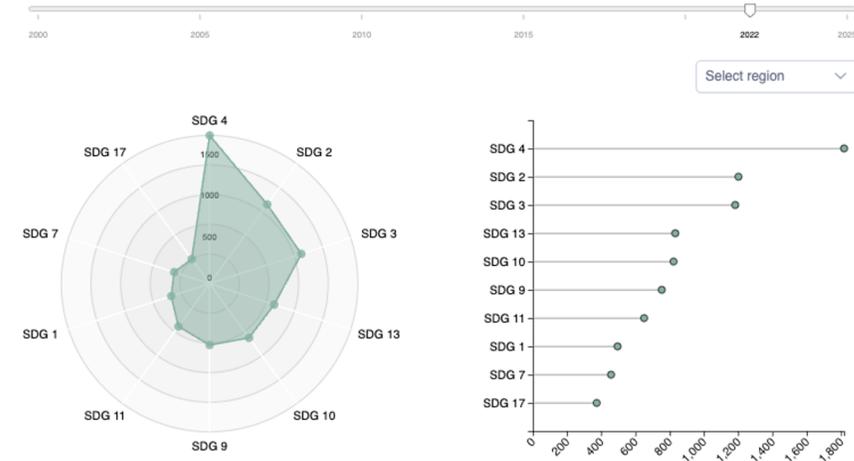


ai4gov-project.eu



101094905

- Evaluate the amount of SDG-classified topics across time
- Compare by geographic world regions to explore prioritisation allowing for comparison
- Encode the SDG topics based on wikidata concepts identified in the textual documents independently of their language



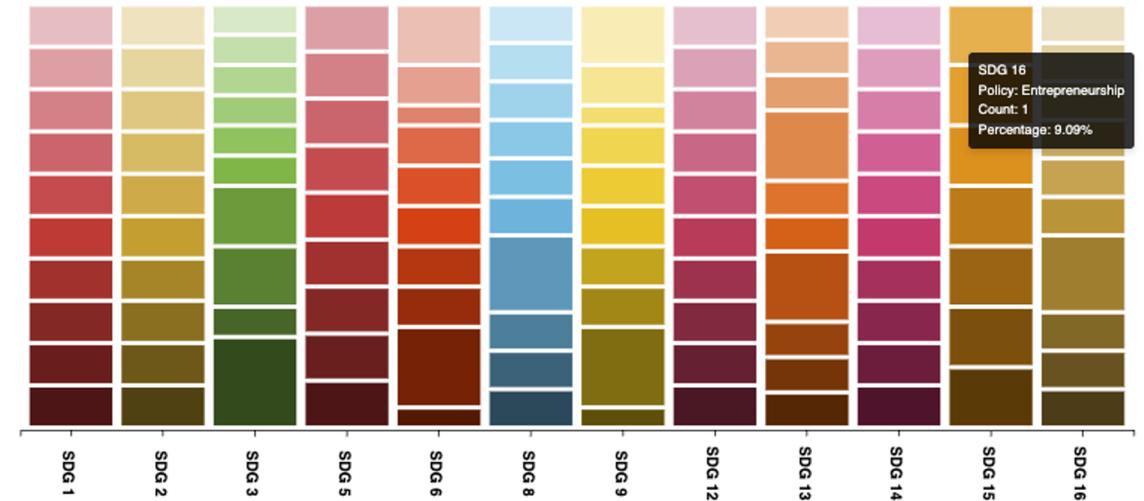
Coverage Bias Aware.



raido-project.eu



101135800



- Integration of cognitive science with computational social science.
- Empirical framework for bias detection in AI discourse.
- Use of semantic linking to align media content with AI taxonomy.
- Identification of cognitive mechanisms shaping AI adoption narratives.
- First large-scale framework connecting cognitive biases to AI innovation adoption.

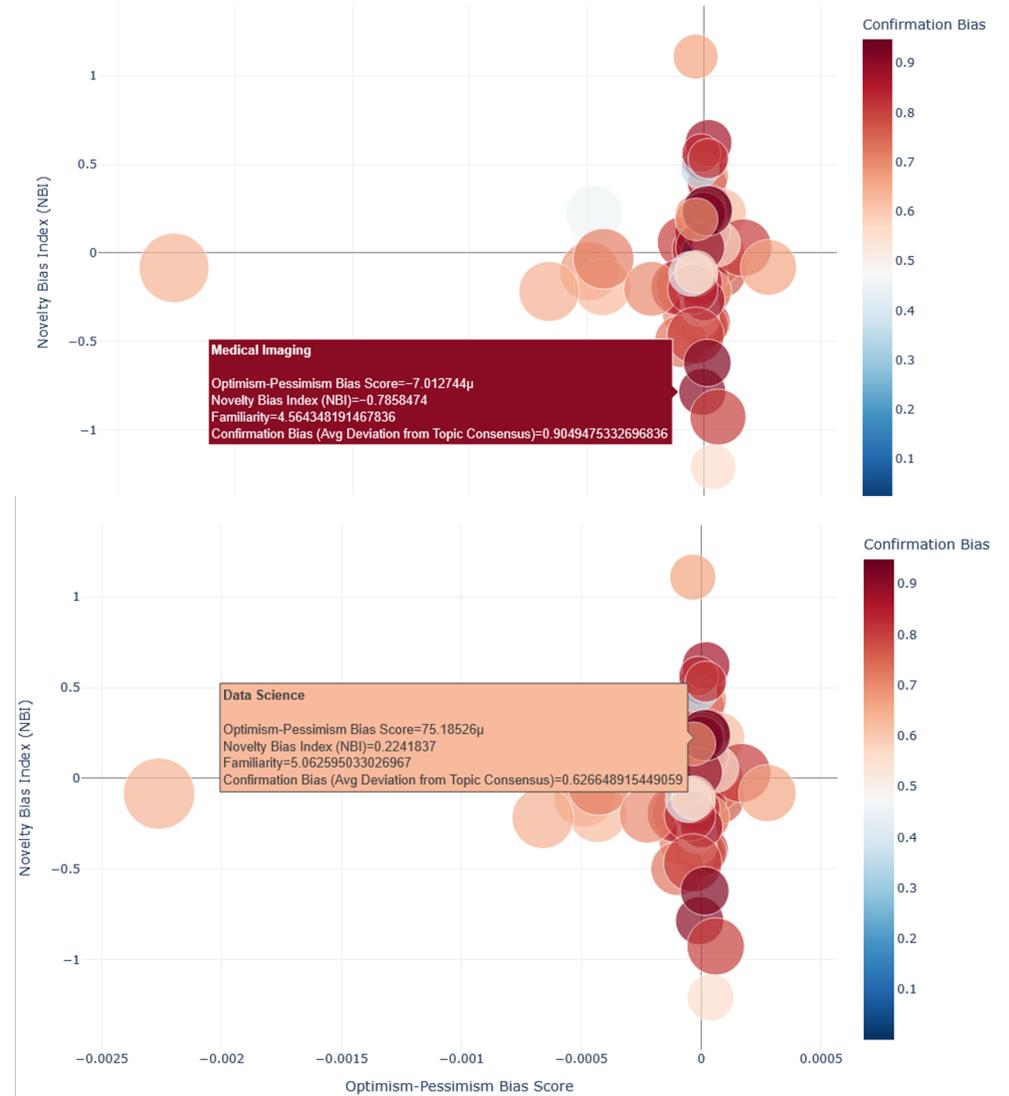
## Cognitive Bias Aware.



elias-ai.eu



101120237



OECD AI Policy Observatory | GPAI

Home > Tools & metrics > Tools

## Catalogue of Tools & Metrics for Trustworthy AI

These tools and metrics are designed to help AI actors develop and use trustworthy AI systems and applications that respect human rights and are fair, transparent, explainable, robust, secure and safe.

Overview | Tools | Metrics | About the catalogue | [Contribute to the catalogue](#)

Show tools | Show use cases

SEARCH @

Search by name...

Publication date

TYPE

APPROACH

- Technical
- Educational
- Procedural

TOOL TYPE

Filter by...

OBJECTIVE

Filter by...

USAGE RIGHTS

ORIGIN

STAKEHOLDER GROUP

COUNTRY OF ORIGIN

ORGANISATION

Filter by...

SCOPE

LIFECYCLE STAGE(S)

TARGET GROUP(S)

TARGET USER(S)

TARGET SECTOR(S)

IMPACTED STAKEHOLDERS

PURPOSE(S)

TOOL READINESS

List of tools (921)

[Eticas Bias](#)

Technical | United States | Uploaded on Mar 24, 2025

An open-source Python library designed for developers to calculate fairness metrics and assess bias in machine learning models. This library provides a comprehensive set of tools to ensure transparency, accountability, and ethical AI development.

Objectives: Fairness, Robustness

Related lifecycle stage(s): Operate & monitor, Build & interpret model, Collect & process data

[Behavior Elicitation Tool](#)

Technical, Procedural | France, European Union | Uploaded on Mar 24, 2025

Behavior Elicitation Tool (BET) is a complex-AI system that systematically probes and elicits specific behaviors from cutting-edge LLMs. Whether for red-teaming or targeted behavioral analysis, this automated solution is Dynamic Optimized and Adversarial (DAO) and can be configured to test the robustness precisely and help to have a better control of the AI system.

Objectives: Robustness, Safety

Related lifecycle stage(s): Deploy, Verify & validate, Build & interpret model

[AIRO \(AI Risk Ontology\)](#)

Educational | Ireland | Uploaded on Jan 29, 2025

The AI Risk Ontology (AIRO) is an open-source formal ontology that provides a minimal set of concepts and relations for modelling AI use cases and their associated risks. AIRO has been developed according to the requirements of the EU AI Act and international standards, including ISO/IEC 23894 on AI risk management and ISO 31000 family of standards.

Objectives: Transparency

Related lifecycle stage(s): Operate & monitor, Verify & validate

[COMPL-AI](#)

Technical | Switzerland, European Union | Uploaded on Jan 24, 2025

COMPL-AI is an open-source compliance-centered evaluation framework for Generative AI models

Objectives:

Related lifecycle stage(s):

Overview | Tools | Metrics | About the catalogue | [Contribute to the catalogue](#)

## CarefulAI: Prompt-LLM Improvement Method (PLIM)

Website

Educational | United Kingdom | Uploaded on Dec 9, 2024

When working with large language models (LLMs), accuracy is important. However, there is a lack of understanding of the co-dependency between LLM outputs and prompts. Existing LLM benchmarks do not specify this; they allude to historical accuracy scores on LLM benchmarks that may not be relevant to the end user. In addition, LLMs are usually dynamic in practice. Their behaviour is not static, but changes over time, and often cannot be explained by LLM providers. Users, therefore, can only partially depend upon LLM benchmarks. In practice, to make LLMs fit for purpose and safe, users are required to constantly test Prompt-LLM outputs for specific cases. This can be time-consuming.

CarefulAI's approach to this is based on the discovery that by serving a model with a standard set of end-user-specific examples of questions and answers—validated by the end-user community (with each prompt validated by a minimum of 3 subject matter experts/end users), the time taken to get acceptable answers is significantly reduced (tenfold). In addition to getting Prompt-LLM combinations that are deemed safe, the approach enables sector/subject matter prompt benchmarking against multiple models.

PLIM is designed to make benchmarking and continuous monitoring of LLMs safer and more fit for purpose.

Website | GitHub | Hugging Face

About the tool

You can click on the links to see the associated tools

Developing organisation(s): CarefulAI

Objective(s): Performance

Impacted stakeholder(s): Consumers, Regulators

Country of origin: United Kingdom

Type of approach: Educational

Maturity: Implemented in multiple projects

Usage rights: Fee-based

Overview | Tools | Metrics | About the catalogue | [Contribute to the catalogue](#)

Show metrics | Show use cases

SEARCH @

Search by name...

Relevance

OBJECTIVE

Filter by...

SCOPE

RISK MANAGEMENT STAGES

Filter by...

Assess

Assess risks & impacts

Define

Govern

Treat

Treat: Cease risks & impacts

PURPOSE(S)

Filter by...

SUBMIT A METRIC

If you have a tool that you think should be featured in the Catalogue of AI Tools & Metrics, we would love to hear from you!

SUBMIT

List of technical metrics (130)

This page includes technical metrics and methodologies for measuring and evaluating AI trustworthiness and AI risks. These metrics are often represented through mathematical formulas that assess the technical requirements for achieving trustworthy AI in a particular context. They can help to ensure that a system is fair, accurate, explainable, transparent, robust, safe, or secure.

[Accuracy](#) 168 related use cases

Accuracy is the proportion of correct predictions among the total number of cases processed. It can be computed with:

Accuracy = (TP + TN) / (TP + TN + FP + FN), where:

TP: True positive

TN: True negative

FP: False positive

FN...

Objectives: Performance, Robustness

[Mean Intersection over Union \(IoU\)](#) 35 related use cases

Mean Intersection over Union (IoU) is the area of overlap between the predicted segmentation and the ground truth divided by the area of union between the predicted segmentation and the ground truth.

For binary (two classes) or multi-class segmentation...

Objectives: Performance, Robustness

[Anonymity Set Size](#) 32 related use cases

Website | GitHub | Hugging Face

Apache 2.0

Data scientist

Developer

Other

Academia

Business

Government

Always up to date

International

ai responsible performance benchmarking ilm



<https://oecd.ai/en/catalogue/tools>

# Bias Detector Catalogue

The Bias Detector Catalogue stands as a pioneering tool, met expansive repository represents a concerted effort by innovat tailored to diverse stages of the training process. From data c Catalogue is a testament to the collective determination to m



**Bias in Automated Speaker Recognition** Accuracy: MODERATE Cost: MODERATE

<b>AIF360: AI Fairness 360 toolkit</b>	Accuracy: HIGH	Cost: LOW
<b>FairMLHealth</b>	Accuracy: UNKNOWN	Cost: LOW
<p>Source: <a href="https://github.com/KenSciResearch/fairMLHealth">https://github.com/KenSciResearch/fairMLHealth</a> Type: MITIGATION Programming Language: PYTHON</p> <p><b>Description:</b> FairMLHealth is a healthcare-specific tool for bias analysis. It provides machine-learning fairness, healthcare applications, and variation analysis.</p> <p><b>Applicability:</b> HEALTHCARE</p> <p><b>Limitations:</b> The 'fair' range to be used for these metrics requires judgement on the part of the analyst.</p> <p><b>References:</b> Ahmad et al., (2020). Fairness in Machine Learning for Healthcare. Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery &amp; Data Mining. <a href="https://doi.org/10.1145/3394486.3406461">https://doi.org/10.1145/3394486.3406461</a>.</p>		
<b>Mitigating Unwanted Biases with Adversarial Learning</b>	Accuracy: HIGH	Cost: LOW
<b>Themis-ml: A Fairness-Aware Machine Learning Interface for End-To-End Discrimination Discovery and Mitigation</b>	Accuracy: UNKNOWN	Cost: UNKNOWN
<b>Bias in Automated Speaker Recognition</b>	Accuracy: MODERATE	Cost: MODERATE
<b>Bias Assessment Metrics and Measures</b>	Accuracy: UNKNOWN	Cost: UNKNOWN
<b>Biaslyze</b>	Accuracy: UNKNOWN	Cost: LOW



<https://cluster-ai4gov.euprojects.net/>



**International Research Centre on Artificial Intelligence  
under the auspices of UNESCO (Category II)**  
Jožef Stefan Institute, Jamova cesta 39, SI-1000 Ljubljana  
E: [info@ircai.org](mailto:info@ircai.org) | W: <https://ircai.org/>

**Joao Pita Costa**  
[in/joaopitacosta/](#)

