



Centre  
Under the auspices  
of UNESCO

**IRC AI**

International Research Centre  
on Artificial Intelligence  
under the auspices of UNESCO

**GITEX**

**ASIA**

*Singapore*



EVERYTHING  
-SINGAPORE-



IN AFRICA

# BUILDING RESPONSIBLE AI ECOSYSTEMS: FROM THEORY TO ACTION FOR PUBLIC GOOD

**GITEX ASIA**

Singapore, 9.4.2026



**NAIXUS**

Network for Artificial Intelligence, Knowledge and Sustainable Development -  
arena and central meeting point between AI and SDGs

**IRC AI**

International Research Centre  
of Artificial Intelligence  
under the auspices of UNESCO



IN AFRICA

**GITEX**

**ASIA**

*Singapore*



EVERYTHING  
-SINGAPORE-

# IRCAI: Who are we and what do we do?



- LAUNCH publicly in March 29 and 30 2021 with **1083 registered participants** from **123 countries** as an ecosystem of research institutions, companies, start-ups and non-profit players.
- Partially founded by the Government of the Republic of Slovenia and has its headquarters at the **Jožef Stefan Institute** in **Ljubljana, Slovenia**.
- Connections to **160 private companies**, **76 public entities** (Foundation/ non-profit), **21 public (University)**; **6 public (Government)**, and **6 Public (NGO)**, helped/ing preparing applications for IRCAIs in other regions (**Morocco, Saudi Arabia, Ghana, New Zealand**)
- IRCAI's director **John Shawe-Taylor** is Chairholder of the **UNESCO Chair in Artificial Intelligence** within the UCL doing work in AI and acts as the trustee of the Knowledge 4 All foundation doing work in the Global South.

IRCAI VIRTUAL LAUNCH  
<https://ircai.org/launch-2021/>



# IRCAI INTERNATIONAL RESEARCH CENTRE FOR AI UNDER THE AUSPICES OF UNESCO

## Policies Innovation

(UNESCO, Council of Europe, OECD, EU, GPAI, National policies, D4D)

## Research on consequences of AI

(Legal frameworks, Ethics in AI, Added Value Models, Jobs,...)

## Projects addressing global challenges with AI

(Global Ground Truth, Open Education)

## Large Open Analytics Infrastructures

(SDG Observatories, Open Libraries)

## Verification programme

(clearing house, mediation verification programme, validation of the ethicality of algorithms and data management)

## Financing programme

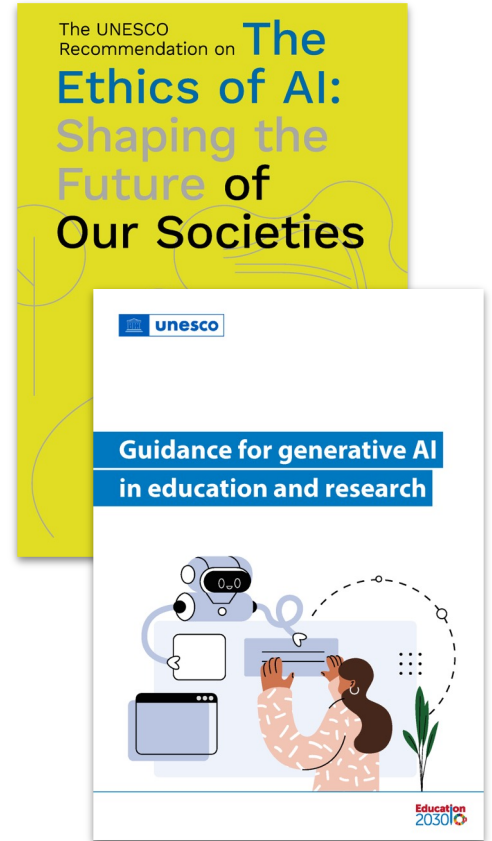
(Social Impact Bonds, venture capital to finance research and acceleration of start-ups)

## Capacity building, Awareness raising

(Journal, AI Olympics, Open education, Events, Formal Educational Programs on AI)

## Global Network of AI competences

(NAIXUS, IRCAI backbone network)



## FOCUSING ON SCIENTIFIC / TECHNOLOGICAL MAJOR CHALLENGES

Focus on scientific or technological major challenges, with the primary goal to reinforce capacity and progress in critical technologies

## SPREADING ADVANCED KNOWLEDGE

Building on existing efforts of already existing platforms, networks and projects, the network will develop mechanisms to spread the latest and most advanced knowledge to all the AI-labs in the United Nations five geographical regions: Africa, Americas, Asia and the Pacific, Europe and Central Asia, and the Middle East and prepare the next generation of talent in AI. Such mechanisms will be defined with all network partners

## DEVELOPING SYNERGIES AND CROSS-FERTILIZATION

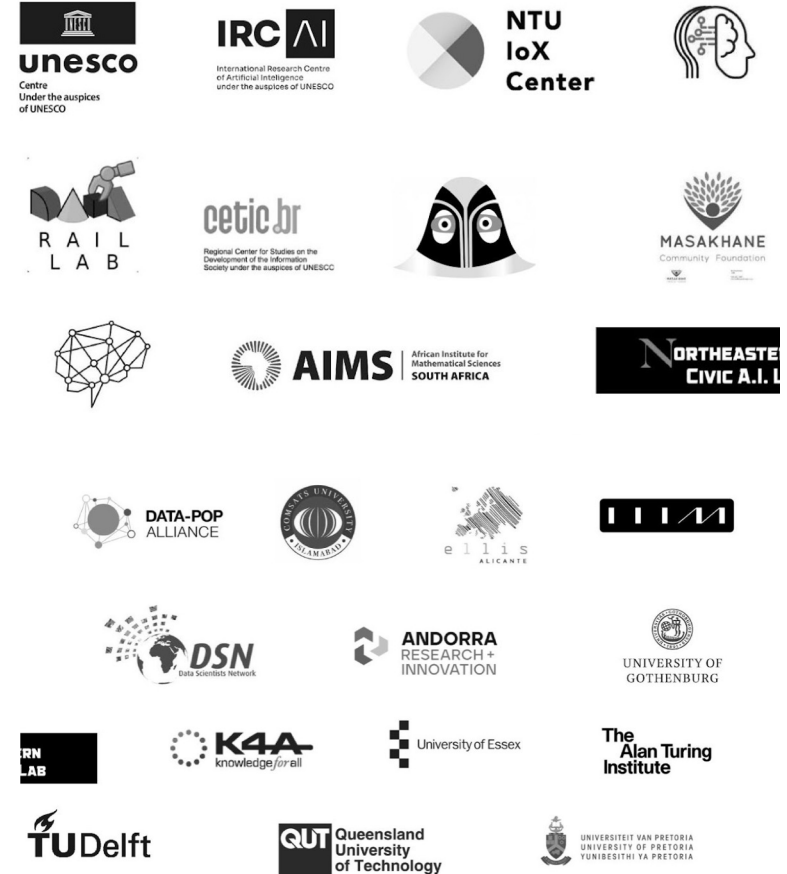
Another objective is also to develop synergies and cross-fertilization between industry and investors in the network of excellence centres, in particular through internships of academic staff (at all levels) in industry, or PhD programmes with industry;

## BECOMING A VIRTUAL LABORATORY

The network will form a common resource and will become a shared facility, as a virtual laboratory offering access to knowledge and expertise and attracting scientists, investors, policymakers and new talents. It should become a reference, creating an easy entry point to AI excellence across the Global South as well and should also be instrumental for its visibility

## CREATING A DATA DRIVEN BENCHMARK SYSTEM

Creating a data driven benchmark system for the measurement of SDGs to further create a social impact investment ecosystem for AI research and companies



# IRCAI AWARD



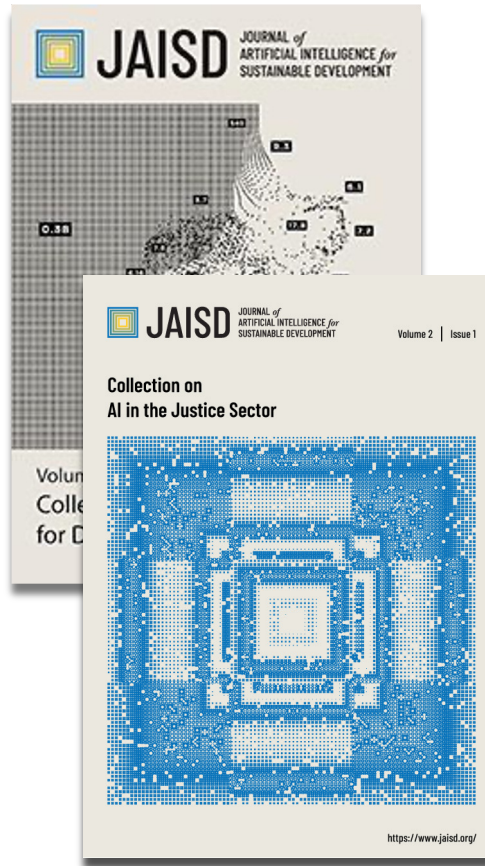
**Adriana-Eufrosina Bora**

Project AIMS

(Artificial Intelligence against Modern Slavery)



# IRCAI JOURNAL



# IRCAI LIBRARY



Find out more at [ircai.org](http://ircai.org)



Paid - AI Summit



## PANEL DISCUSSION | RESPONSIBLE AI IN ACTION: ACCELERATE AI FOR PUBLIC GOOD IN ASIA

As agentic AI expands across Asia, the key challenge is not scale alone, but relevance, inclusion, and trust. Frugal edge AI and participatory citizen science offer...

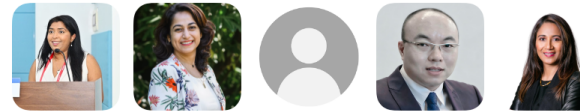


🕒 5:15PM - 5:55PM

📍 Main Stage

## Panel Discussion | Building the AI Vanguard: From Classroom Literacy to Olympiad Excellence

This panel explores the future of AI education across diverse global ecosystems. Moderated by the lead of South Africa's National AI Olympiad, the panel unites...



🕒 11:50AM - 12:15PM

📍 10X Stage

## Keynote | AI's Global Impact: Reshaping Society and Education



🕒 11:40AM - 11:50AM

📍 10X Stage

## Senior Showcase & Pitch Competition | AI4SDG3 Challenge With IRCAI Top100 Asia

Innovators representing IRCAI's Top 100 Asia take to the stage with their AI4SDG3 Challenge presentations. [3min]



🕒 12:25PM - 12:55PM

📍 10X Stage

## Junior Pitch Competition | AI4SDG3 Challenge High School Student Category

High school students take to the stage to present their AI4SDG3 Challenge submissions. [3min]



🕒 12:55PM - 1:10PM

📍 10X Stage

🕒 11:30AM -

📍 10X Stage





تحت الرعاية السامية لصاحب الجلالة الملك محمد السادس  
Under the High Patronage of His Majesty King Mohammed VI

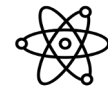
**GITEX**  
AFRICA  
Morocco

UNDER THE AUTHORITY OF  
REPARTNERSHIP WITH  
ORGANIZED BY  
#ADD  
KAOÛN  
INTERNATIONAL  
14 - 16 APRIL 2025 MARRAKECH

# THE AI FOR SUSTAINABLE AFRICAN CITIES AND COMMUNITIES CHALLENGE

is driving innovation to create  
**SMARTER** more resilient cities  
across Africa.

AI will be key in solving urban challenges  
such as climate change, resource  
management, and social inclusion



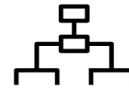
**1. Science Knowledge Extraction for Good Practices**



**2. Data Visualization Narratives on Smart Water**



**3. Appropriate application of LLMs**



**4. SDG6-focused classifiers**



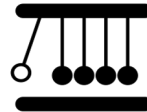
**5. Early detection of floods based on social media**



**6. Looking for trends with AI**



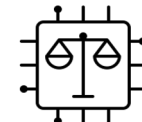
**7. Automated sentiment analysis on media**



**8. Causality analysis on water topics from social media**



**9. Analysis of waterborne disease on scientific articles**

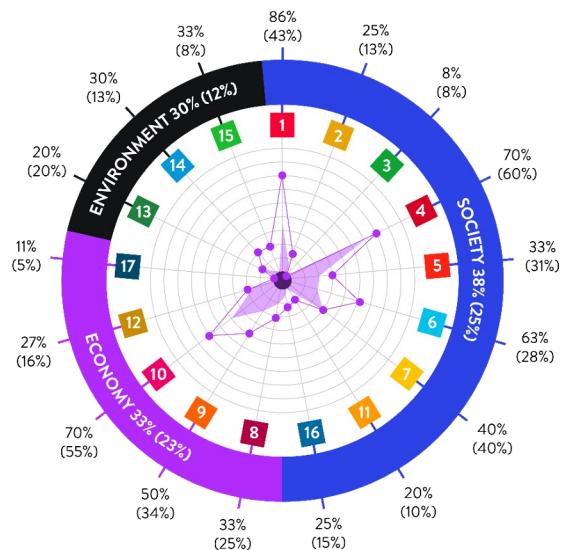
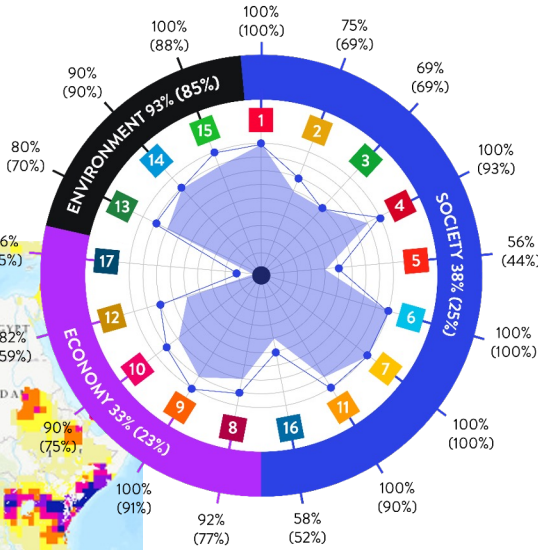


**10. Bias, Ethics & Responsible AI**

# SDG Agentic AI #AI4SDG

## A Positive Impacts of AI: 79% (71%)

## B Negative Impacts of AI: 35% (23%)



DeepLeaf



Welcome to Marrakech  
Your smart urban mobility companion

Ask GO

Traffic Conditions  
Last updated: 9:15 AM  
Avenue Mohammed V  
morning commute  
Delay: 5 min

Jemaa el-Fnaa area  
tourist activity  
Delay: 10 min

Route to Train Station  
clear roads

Quick Access

- Day Trip Planner
- View Map
- Train Schedule
- Top Attractions

Chat with SLTVerse Assistant

THE AI FOR SUSTAINABLE AFRICAN CITIES AND COMMUNITIES CHALLENGE is driving innovation to create **SMARTER** more resilient cities across Africa.

AI will be key in solving urban challenges such as climate change, resource management, and social inclusion

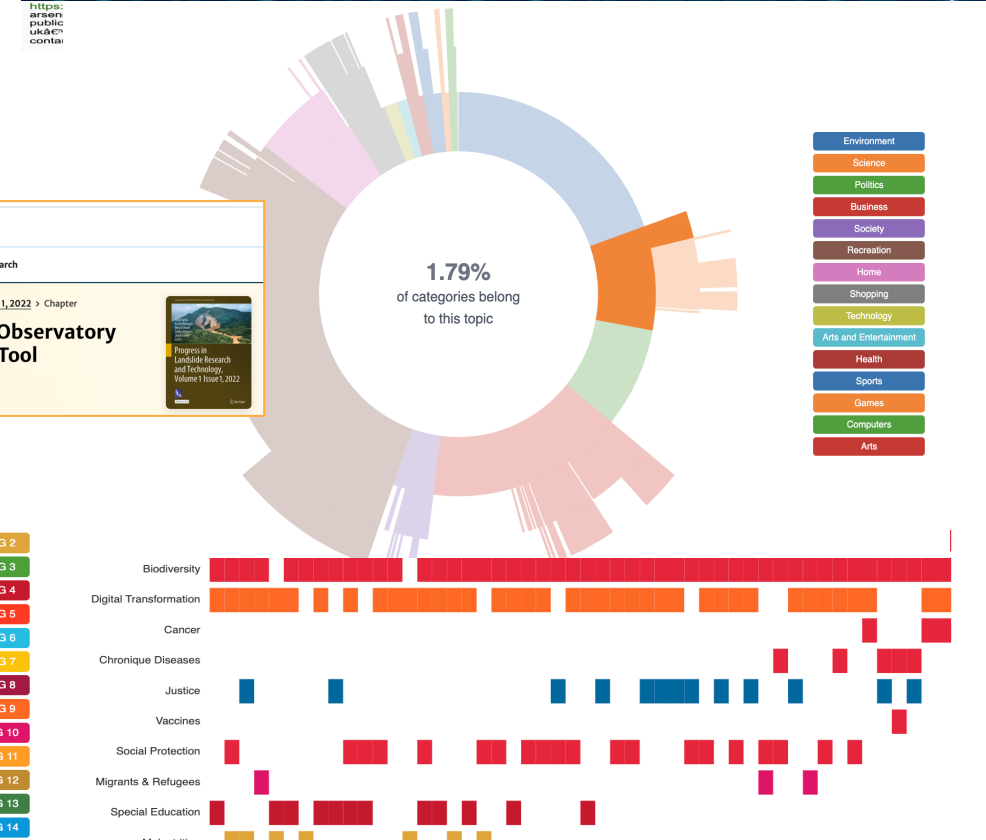


# Disaster Risk Reduction & Landslides

The Landslide Observatory is developed to provide a multifaceted view on landslide-related information worldwide. It includes data from daily multilingual news and published science, as well as yearly ingested statistical indicators, AI policies and open education resources.



<https://www.unesco-floods.eu/>



# Good Practices in Water Sustainability

The water sector is facing rapid development in the direction of the smart digitalisation of resources, much motivated and supported by the UN global initiative of the Sustainable Development Goal 6. In that context, the efforts to address the specific challenges related to water management data and priorities multiply globally.

NAIADES WATER OBSERVATORY

ABOUT

FAQ

## MONITORING WATER RELATED EVENTS TO EXPLORE RELEVANT WATER ISSUES

Enabling insights into water related issues through data analysis and AI. Built for water experts and the general public.

ABOUT

### EXPLORE WATER EVENTS THROUGH AVAILABLE TOOLS



INDICATORS



MEDIA



RESEARCH



RESOURCES



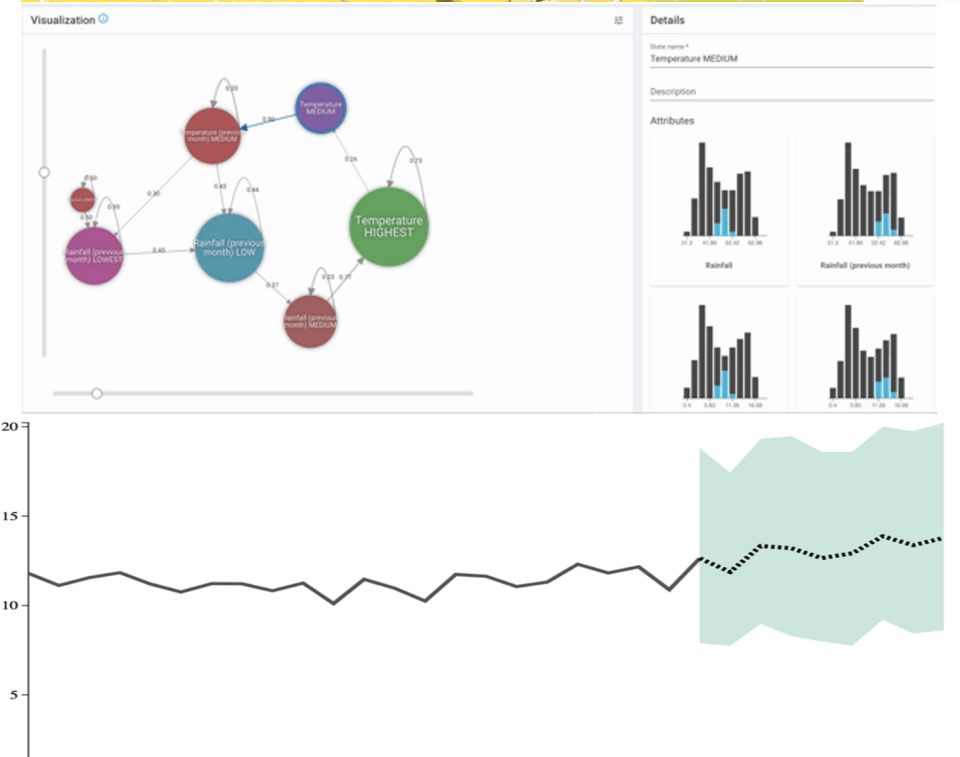
received funding from the European Union's research and innovation program under grant 820985

Supported by eventregistry



AI FOR WATER SUSTAINABILITY  
REVOLUTIONIZING RESOURCE MANAGEMENT  
CHALLENGE MOROCCO 2024  
FRIDAY 26TH OF APRIL, 3PM MOROCCO

BROUGHT TO YOU BY:  
IN AFRICA  
IRC AI  
unesco  
GITEX AFRICA Morocco  
AIMovement



# AI Impact on Sustainable development

## WORLD NEWS

Canarias24Horas.com - Tenerife será sed  
Wed Jul 10 2024, 11:33

OPINION - How can Türkiye make Developme  
Wed Jul 10 2024, 11:00

Features Invest in ferti  
Wed Jul 10 2024, 05:59

One Health Approach Can Help Ease Climat  
Wed Jul 10 2024, 07:28

Trilateral Powerhouse: India, Indonesia,  
Tue Jul 9 2024, 16:25

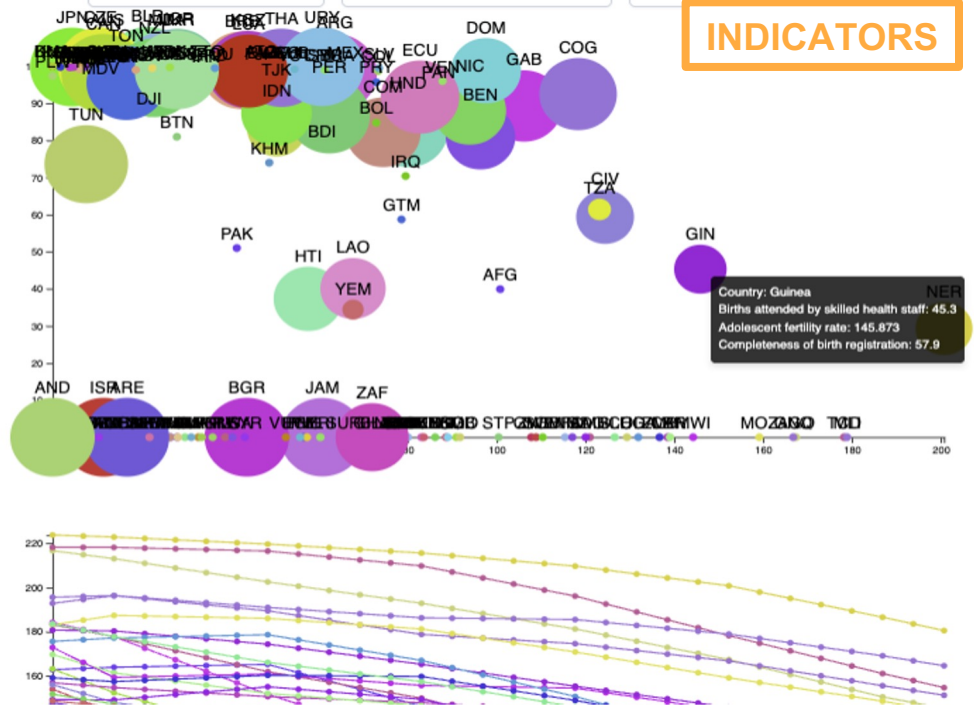
LG USA RENEWS PARTNERSHIP WITH TENNESSEE  
Tue Jul 9 2024, 15:57

Virtus and Quality Clouds Forge Strateg

group private sector digital \$ security  
tech software information world time  
energy use cooperation  
% company year trade  
growth cloud forward

Agricultural Services & Products

## INDICATORS



## SCIENCE



## POLICY

Battle-related deaths (number of people)

2022 18.9%  
2019 3.57%

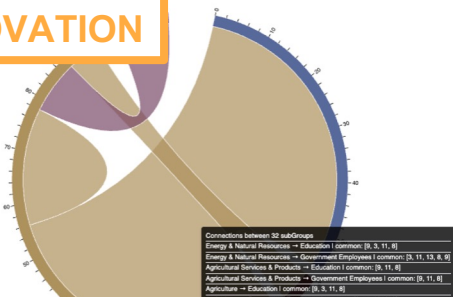
Risk Employment  
United Nations  
Human rights  
Aid Free will Law  
Sovereign state  
Other

SDGs and popularity in educational resources on dealing with Ethnic violence

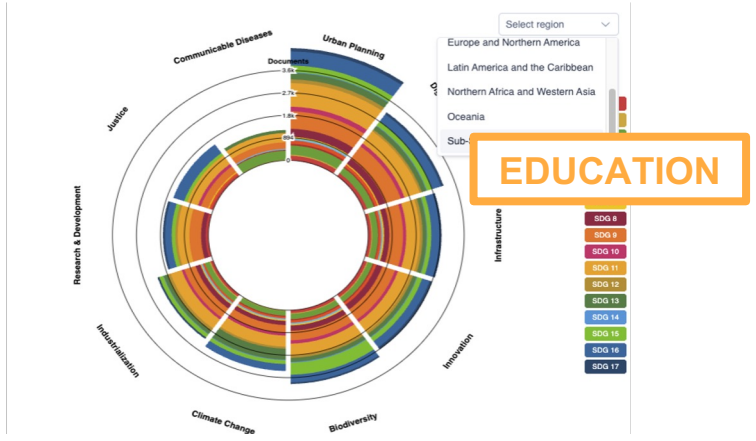
SDG 5	11.78%	SDG 8	10.9%
SDG 1	1.9%	SDG 2	10.35%
SDG 9	9.14%	SDG 15	10.05%
SDG 12	9.16%	SDG 13	9.88%
SDG 17	9.87%	SDG 7	9.88%

Top 50 values of annotations.title.keyword - Count of records

## INNOVATION

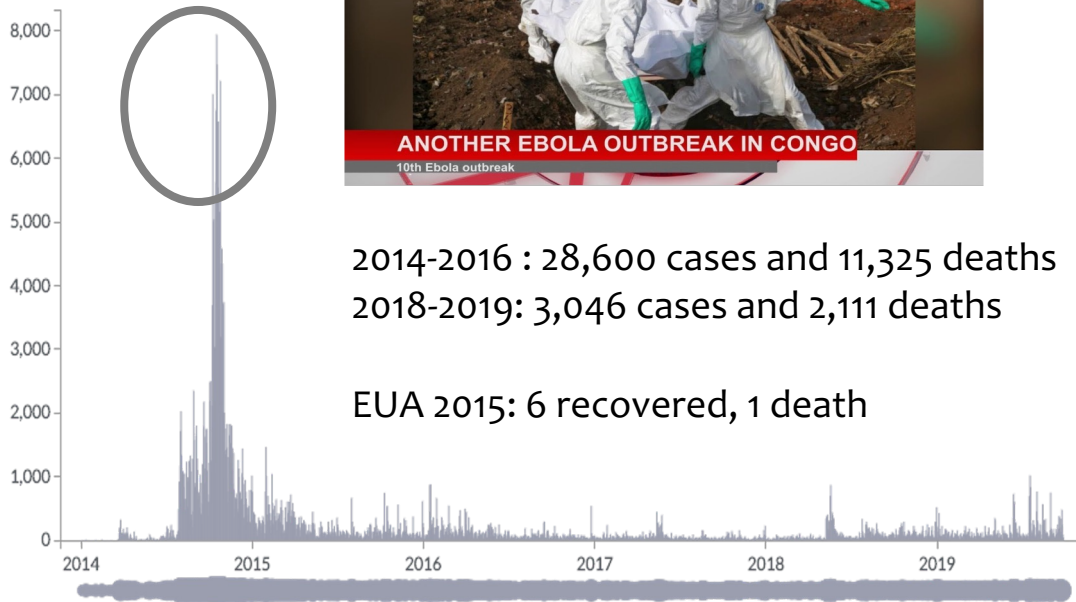


## EDUCATION

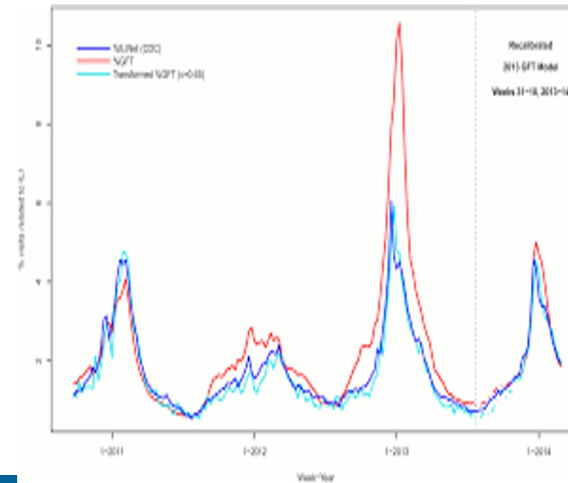


# Bias in News on Science

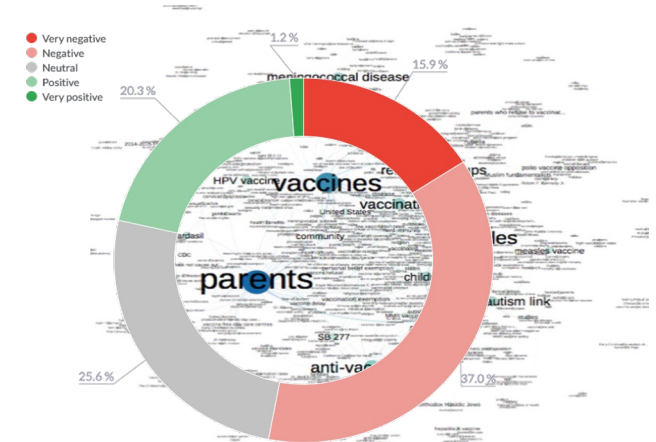
1



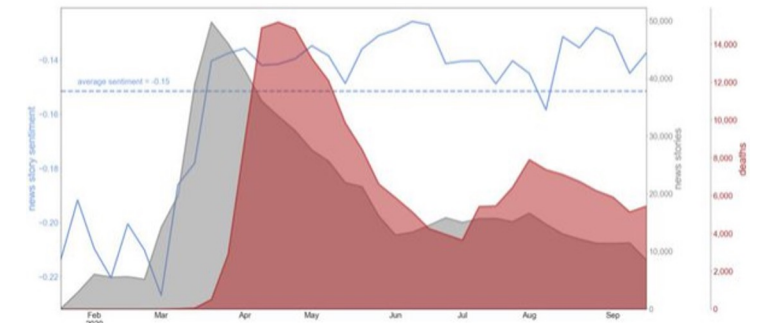
2



3



COVID-19 media coverage sentiment vs. deaths



INFORMATION  
COMPLEXITY

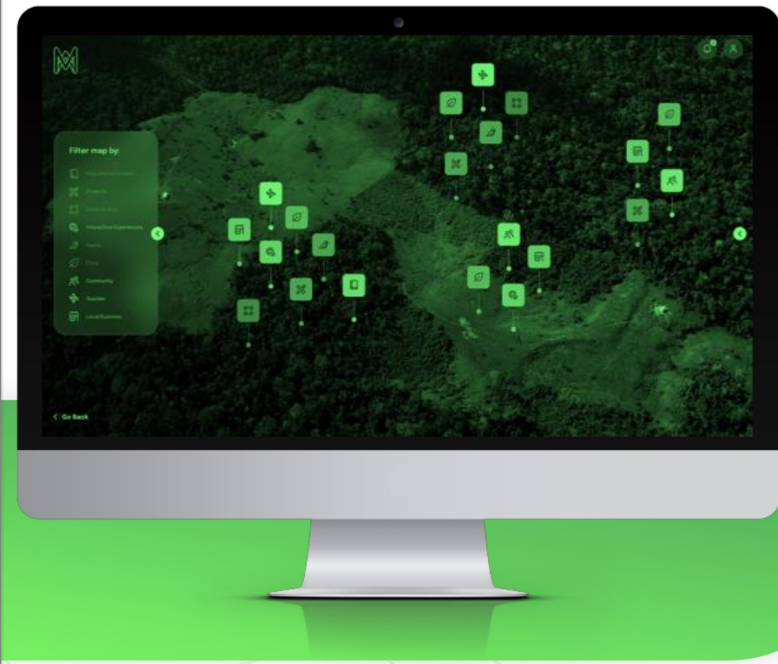
MANIPULATION /  
IDEOLOGY

INNACCURACY /  
MISUNDERSTANDING

ORIGIN LOCATION/  
LANGUAGE



# Road to COP30: AI MetAmazonia



## Real Florest

20,000 hectares that bring real information to the Platform



## Social Impact

Improving the lives of more than 450 families



## Reducing the Ecological Footprint

2.2 million tons of carbon to be mitigated



## Partnerships with Universities

UNIVERSITY OF CAMBRIDGE **Imperial College London**

UNITED NATIONS CLIMATE CHANGE CONFERENCE  
**COP30**  
**AMAZÔNIA**  
 CUIDAR DO PLANETA PARA O FUTURO DA HUMANIDADE

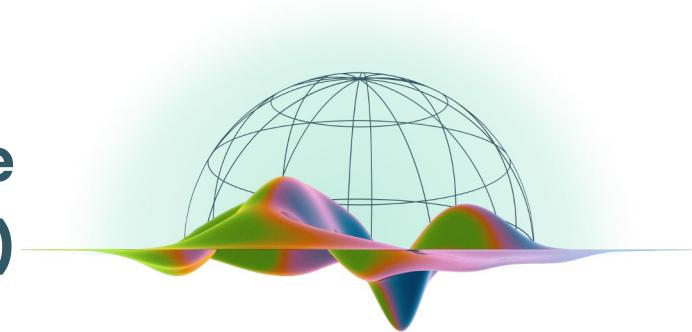




## Recommendation on the Ethics of AI

- 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)

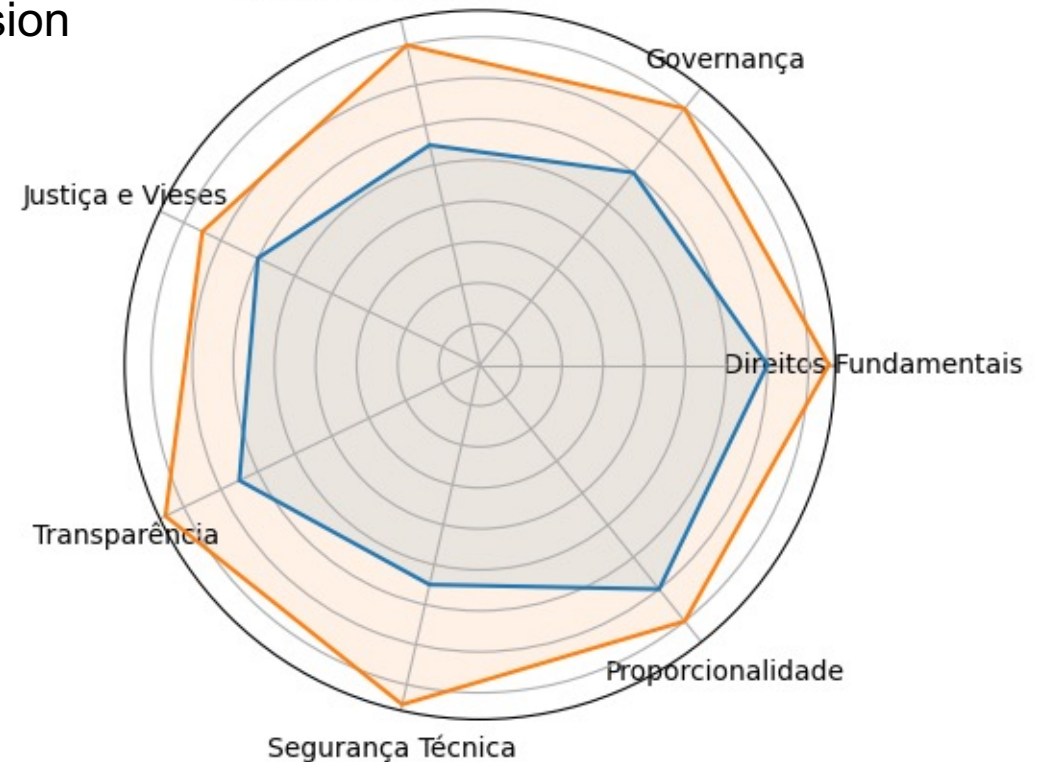


# AI REGULATORY SANDBOX & AI ETHICS

**AI regulatory sandbox** is a controlled environment established by competent authorities where AI systems can be developed, tested, and validated under regulatory supervision before being placed on the market or put into service.

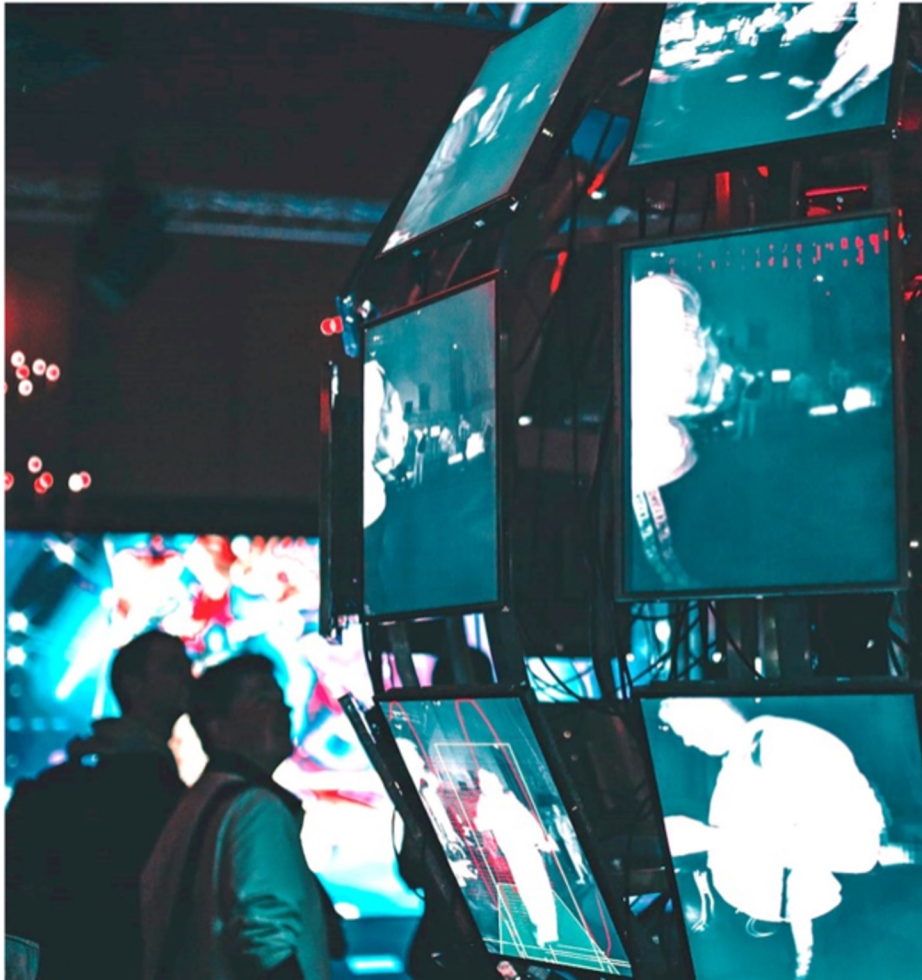
1. New to market and complying to current regulation
2. Centrality of human dignity and fundamental rights (Maximum priority – criterion for remaining in the sandbox)
3. Clear governance and accountability (Without governance → it is not an ethical sandbox)
4. Impact assessment (continuous) (Central element of the sandbox according to UNESCO)
5. Justice, non-discrimination, and algorithmic bias
6. Transparency, explainability, and communication
7. Technical safety and robustness
8. Proportionality and necessity

Radar Ético - SCORM AI vs IntelliDoctor.ai  
Avaliação de Impacto



Framework for AI Explainability, Transparency & Trustworthiness





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**





# Human Rights and AI Systems Interaction



**Human Rights Embedded  
in Functioning of the AI  
Systems**



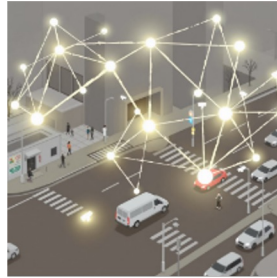
**AI Systems Can  
Violate Human Rights**



**AI Systems Can  
Contribute to Respect  
of Human Rights**



# Human Rights Embedded in Functioning of the AI Systems



## Right to Privacy and Data Protection

- Unauthorized data collection
- AI surveillance without adequate oversight
- Data breaches exposing sensitive data



## Right to Non-Discrimination

- Algorithmic biases that can lead to unfair treatment of vulnerable groups (excluding or targeting)
- Predictive policing or access to public services

## Right to a Fair Trial and Right to an Effective Remedy

- Lack of transparency in algorithmic decisions
- Decision of these systems are unclear and cannot be tested

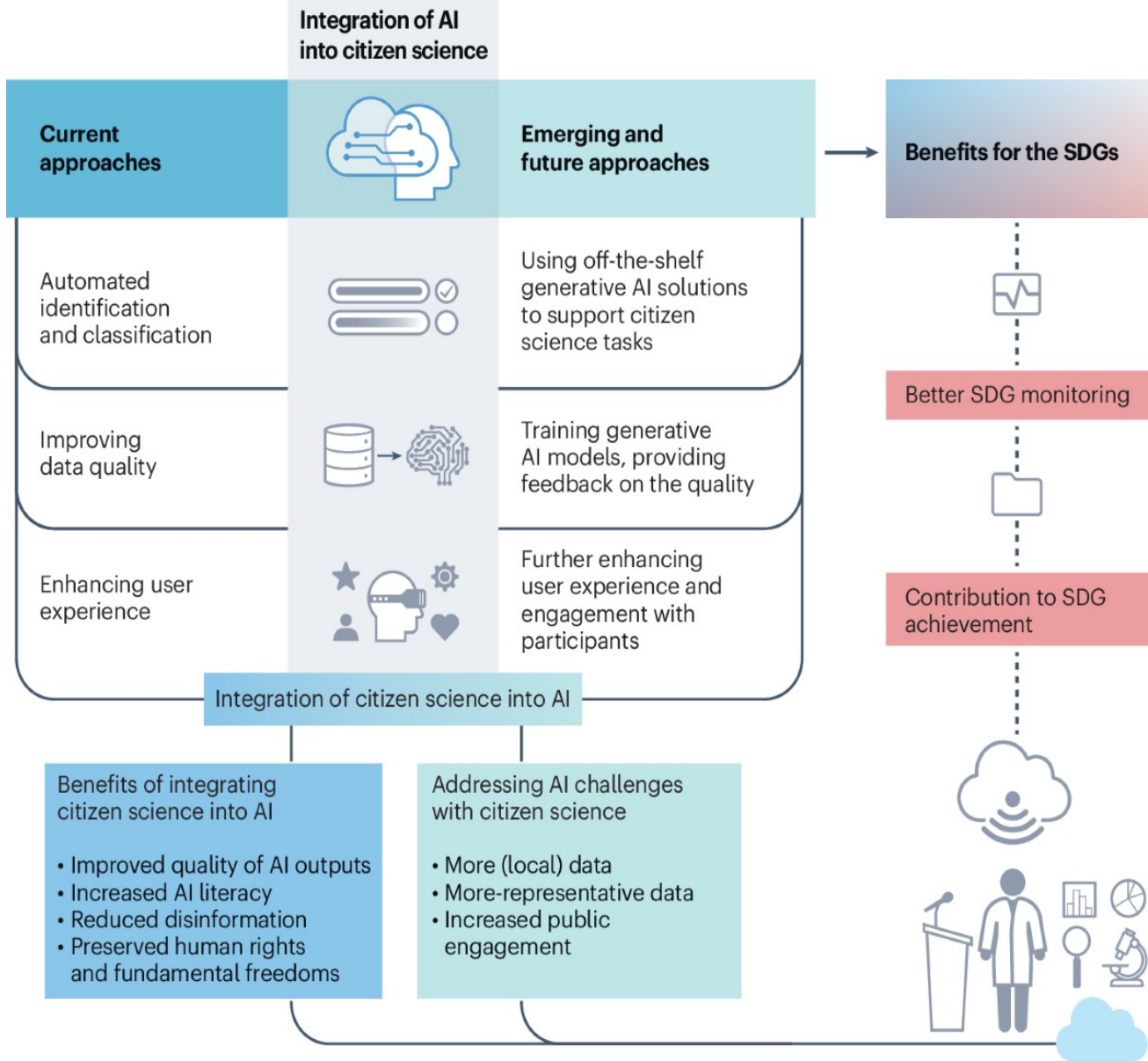


Key facts  
UNESCO's  
Recommendation on  
**the Ethics  
of Artificial  
Intelligence**  
Adopted on 23 November 2021

<https://unesdoc.unesco.org/ark:/48223/pf00000385082>



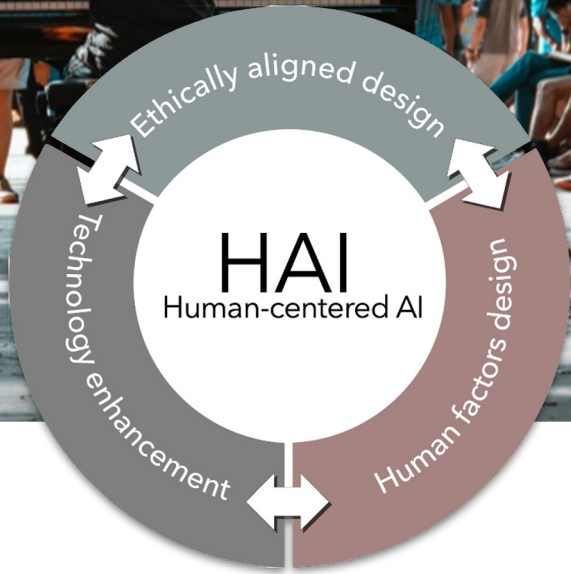
Leveraging the collaborative power of AI and citizen science for sustainable development, Nature Sustainability 2024





- 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.





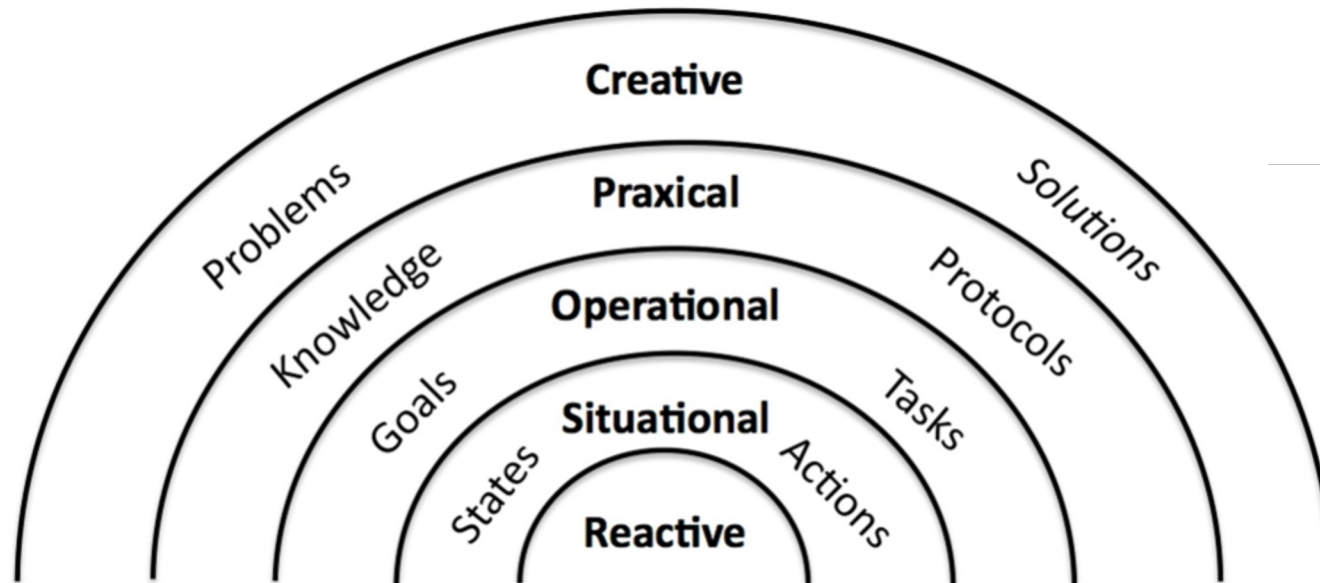
- Human-centric AI is at the centre of the European vision of a positive role for AI
- AI that empowers humans to be more effective, more creative, more understanding
- Human-centric AI (HCAI) has been the focus of the Humane AI Network of Excellence
- How does its research agenda envisage HCAI?

HUMANE  AI NET

**HumanE AI Net:**

**The HumanE AI Network**





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:**

**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  
Net. This document and its contents





- Fairness: several measures that are not equivalent
- Bias: AI has revealed bias in the data
- Privacy: techniques for ensuring personal data is not revealed through models developed by AI systems



An intelligent person



An unintelligent person



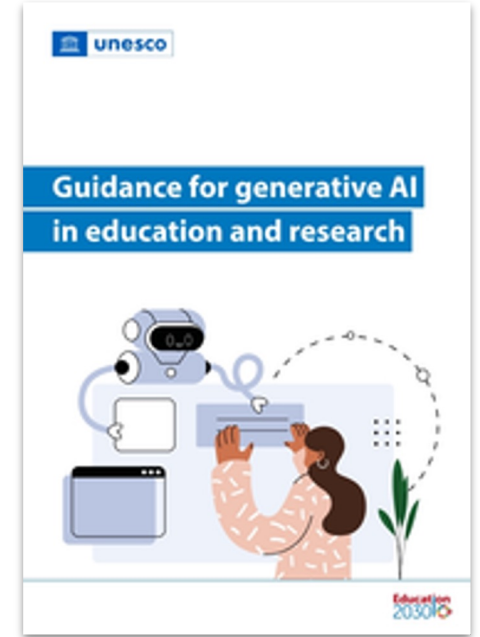
A very competent person



A very incompetent person

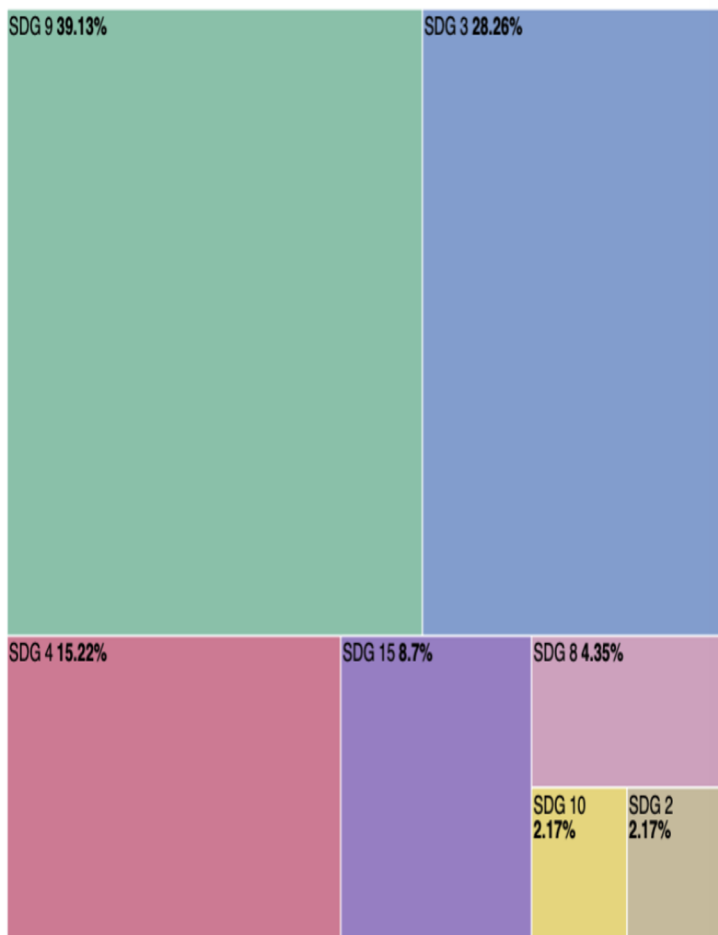
*Images created with DALL-E via Bing Image Creator*





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



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.



## Backend Structure Overview

```
backend/  
├── api/  
│   ├── routes.py           # API endpoints  
│   └── admin_routes.py    # Main API routes  
│                               # Admin API routes  
├── models/  
│   ├── research.py       # Data models  
│   └── admin.py          # Research-related mc  
│                               # Admin-related model  
├── services/  
│   └── llm_service.py    # Service layer  
│                               # LLM interaction ser  
├── utils/  
│   ├── file_manager.py  # Utility functions  
│   ├── admin_utils.py  # File operations  
│   ├── auth.py          # Admin utilities  
│   └── logger.py        # Authentication util  
│                               # Logging utilities  
├── config/  
│   └── config.py        # Configuration  
│                               # Application configu  
├── storage/  
│   ├── prompts/        # Storage for files a  
│   ├── outputs/        # Prompt templates  
│   └── logs/           # Generated research  
│                               # System logs  
├── main.py              # Application entry p  
├── requirements.txt     # Python dependencies  
├── README.md           # Documentation  
├── Dockerfile          # For containerizati  
└── .env.example        # Environment variabl
```

image: [https://www.youtube.com/watch?v=riyh\\_ClshTs](https://www.youtube.com/watch?v=riyh_ClshTs)





PHASE 1  
**DATA INGESTION**



PHASE 2  
**ALGORITHM**



PHASE 3  
**INTERPRETATION**



**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.





**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.

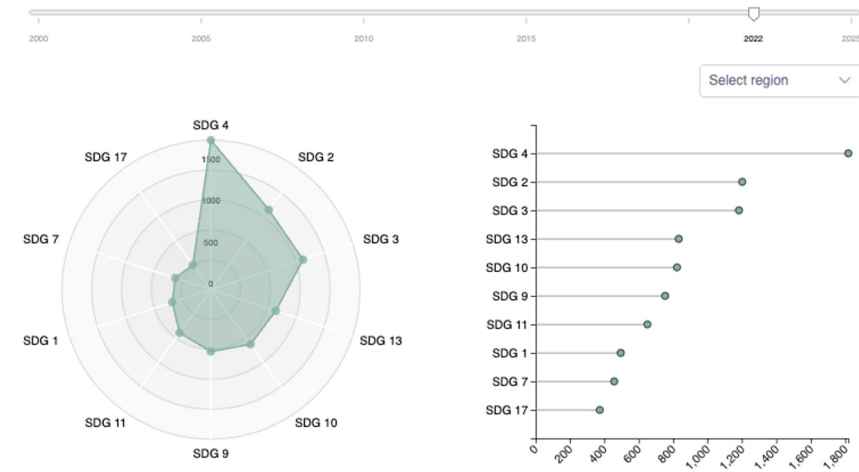
**Environmental Impact:** consider the environmental footprint and strive for sustainable AI research.

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

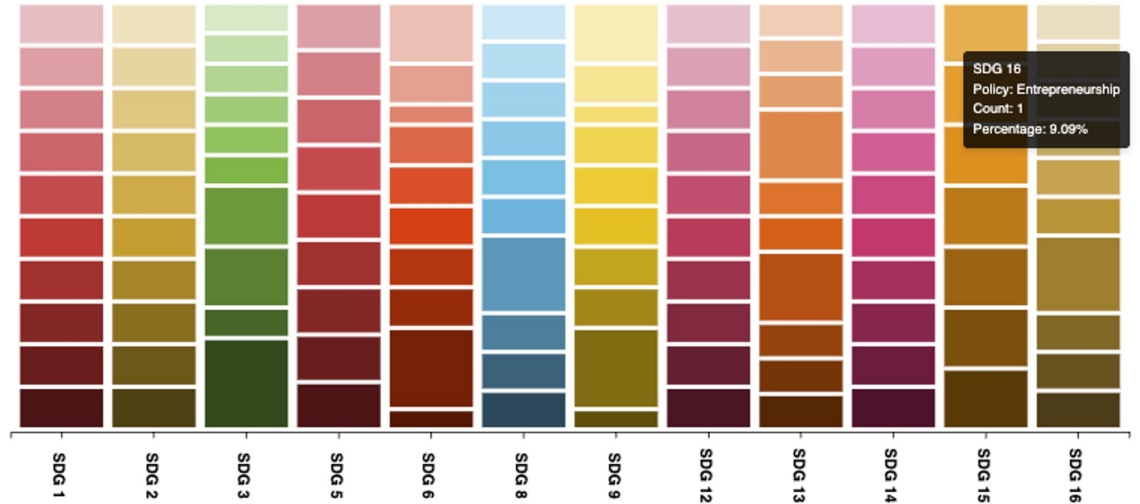
**Stakeholder Participation:** Involve relevant stakeholders, including those affected by AI systems, in the design, development, and deployment processes.



- 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.





**Avoiding Bias:** confirmation bias (favoring information that aligns with one's existing beliefs) or availability bias (relying on immediate examples).

**Transparency:** clearly communicate the methods and tools used in the analysis so that others can understand and potentially reproduce the results.

**Overreach:** Do not make claims or conclusions that go beyond what the data actually supports.

**Accuracy:** Ensure that interpretations are based on accurate and robust findings. Avoid cherry-picking data or results to support a particular narrative.

**Full Disclosure:** Present any limitations, uncertainties, or assumptions associated with the analysis. This allows for a more nuanced understanding of the results.

**Conflict of Interest:** Be transparent about any potential conflicts of interest that might influence the interpretation or presentation of the results.

**Sensitivity:** Be aware of the broader social and cultural context when interpreting data.

**Feedback Loop:** Consider the potential feedback loop effects where the interpretation of data analysis can influence future data collection and results.

**Stakeholder Consideration:** Think about who is affected by the interpretation of the data and ensure that their interests and perspectives are considered.

**Cultural Context:** Especially in global or multicultural studies, understand the cultural nuances and be wary of interpreting data through only one cultural lens.

**Ethical Implications:** Consider the broader ethical implications of any conclusions drawn, especially if they may lead to actions or policies that impact people's lives.

**Peer Review:** Encourage peer review of findings and interpretations to ensure validity and mitigate potential biases.

**Translating to Policy:** If data interpretation is used to inform policy or business decisions, ensure that the recommendations are ethically sound, justifiable, and in the best interest of those affected.

**Honesty:** Always approach data interpretation with honesty, and avoid the temptation to manipulate or selectively present results for personal or organizational gain.

**Continual Learning:** Recognize that interpretations might change as more data becomes available or as methods evolve. Be open to updating or revising interpretations in the light of new evidence.

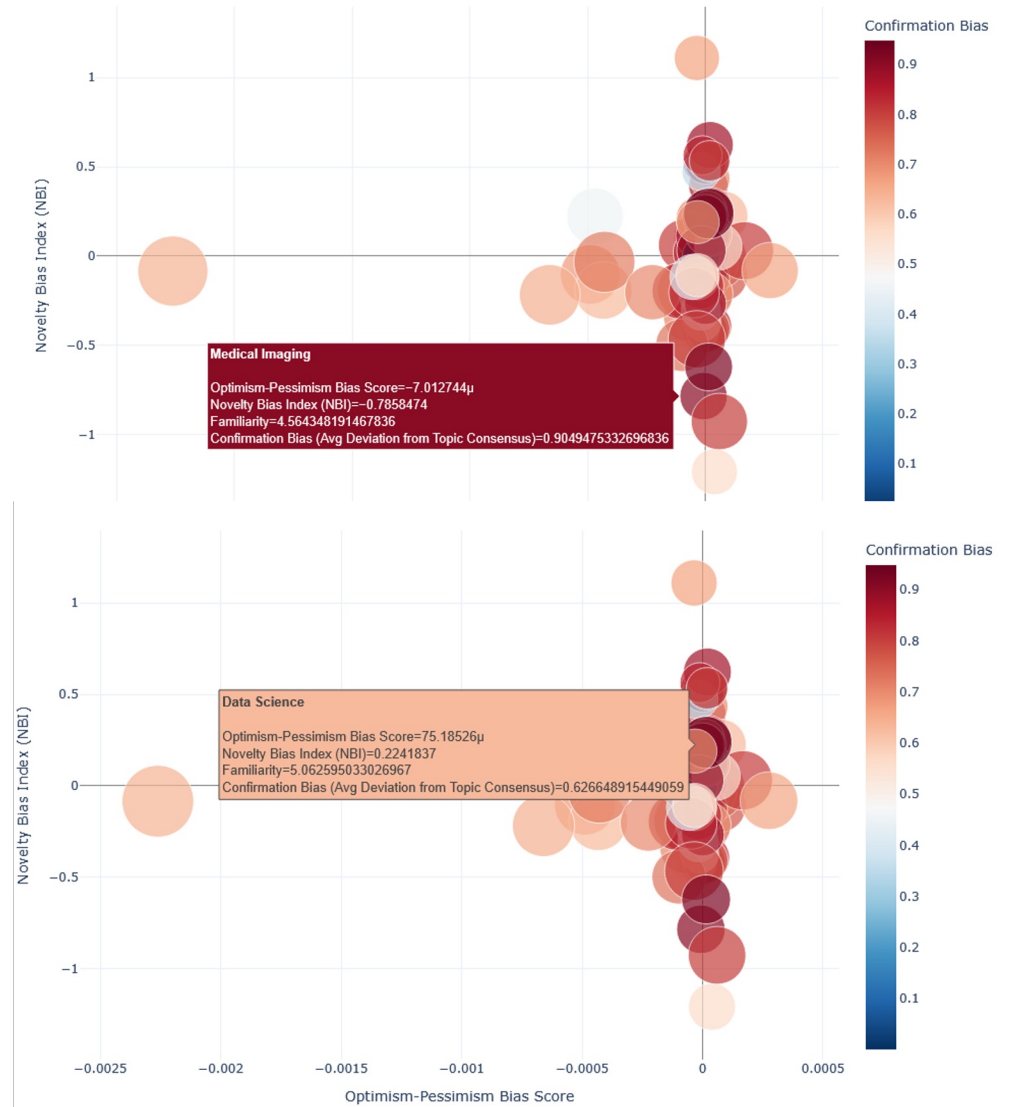


- Integration of cognitive science with computational social science.
- Empirical framework for bias detection in AI discourse.
- Use of semantic linking to align media content with AI taxonomy.
- Identification of cognitive mechanisms shaping AI adoption narratives.
- First large-scale framework connecting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

## Cognitive Bias Aware.

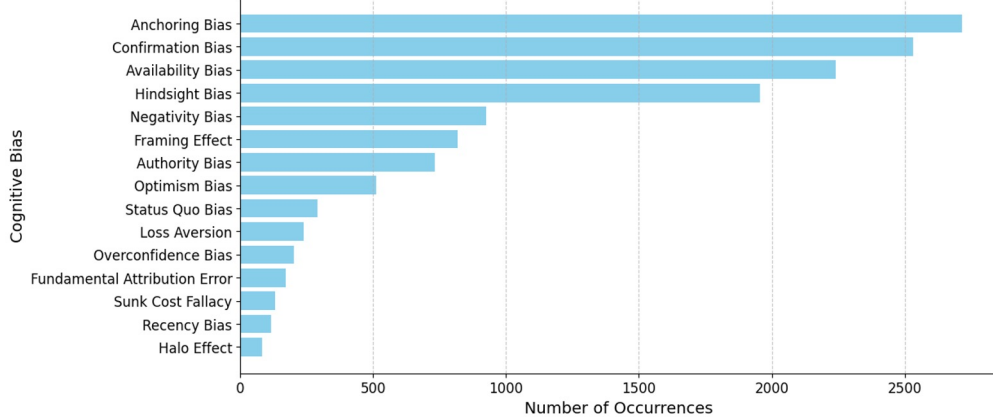


elias-ai.eu 101120237

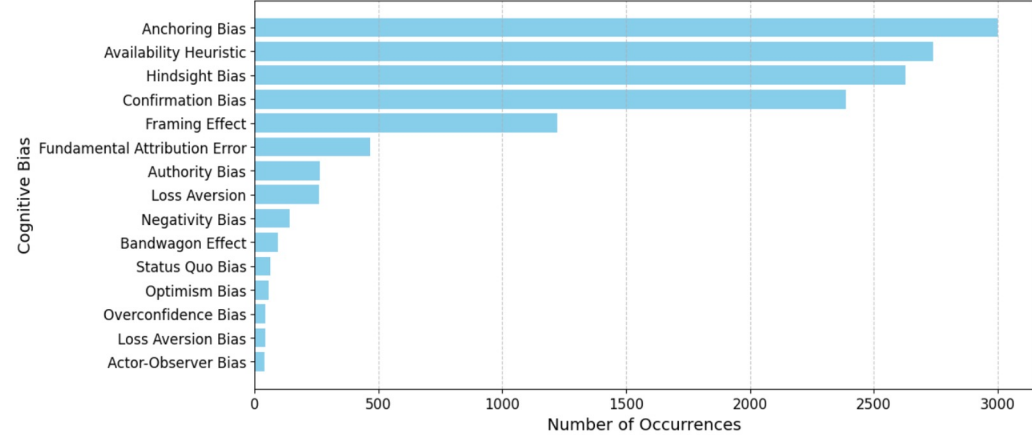


# ELIAS: Cognitive Bias on News

Top 15 Most Common Cognitive Biases



Top 15 Most Common Cognitive Biases

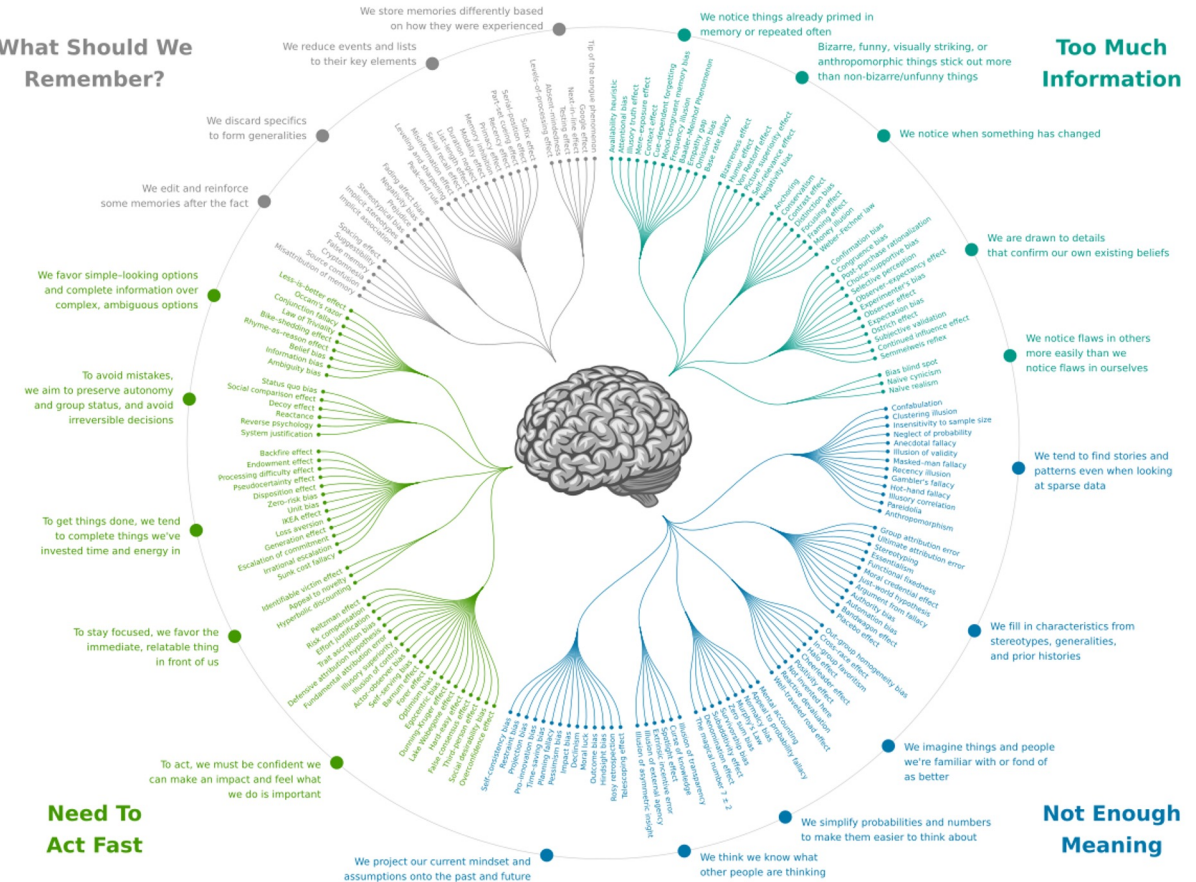


DEEPSEEK  
Topic: self-driving cars

LLAMA  
Topic: self-driving cars

## THE COGNITIVE BIAS CODEX

What Should We Remember?





# ELIAS: AI Cognitive Bias

## Example: self driving cars

source	negative	positive	neutral	other	count	% negativ	% positiv	% neutral	% other
forbes.com	8	9	15	17	49	0.16	0.18	0.31	0.35
dailymail.co.uk	11	3	11	7	31	0.35	0.10	0.35	0.23
finance.yahoo.com	8	5	2	4	19	0.42	0.26	0.11	0.21
theregister.com	9	5	2	2	18	0.50	0.28	0.11	0.11
techtimes.com	7	4	2	3	16	0.44	0.25	0.13	0.19
teslarati.com	7	4	1	3	15	0.47	0.27	0.07	0.20
jalopnik.com	5	0	2	7	14	0.36	0.00	0.14	0.50
electrek.co	5	3	3	3	14	0.36	0.21	0.21	0.21
marketscreener.co	4	4	2	2	12	0.33	0.33	0.17	0.17
news.yahoo.com	6	3	2	1	12	0.50	0.25	0.17	0.08
futurism.com	5	0	4	2	11	0.45	0.00	0.36	0.18
theguardian.com	3	2	0	6	11	0.27	0.18	0.00	0.55
sfist.com	4	3	1	3	11	0.36	0.27	0.09	0.27
benzinga.com	2	2	4	2	10	0.20	0.20	0.40	0.20
nbcbayarea.com	8	0	0	2	10	0.80	0.00	0.00	0.20
businessinsider.co	3	4	2	1	10	0.30	0.40	0.20	0.10
washingtonpost.co	5	2	2	1	10	0.50	0.20	0.20	0.10
techcrunch.com	6	2	0	2	10	0.60	0.20	0.00	0.20
wonderfulenginee	3	1	3	2	9	0.33	0.11	0.33	0.22
ca.sports.yahoo.co	0	3	4	2	9	0.00	0.33	0.44	0.22
carscoops.com	4	5	0	0	9	0.44	0.56	0.00	0.00
autoevolution.com	5	1	0	3	9	0.56	0.11	0.00	0.33
abc7news.com	2	2	2	3	9	0.22	0.22	0.22	0.33

on average, more **negative** sentiment towards hypothesis related to self-driving cars

possible confirmation bias

source	date	sentiment llm	sentiment classic
futurism.com	5/30/2023	neutral	-0.098039216
futurism.com	6/9/2023	The sentiment towards the hypothesis is negative.	0.066666667
futurism.com	7/10/2023	negative	0.098039216
futurism.com	9/16/2023	The hypothesis carries a negative sentiment highlighting potential safety risks associated with Musk's decisions on sensor technology for autonomous driving.	-0.011764706
futurism.com	10/4/2023	cautious	-0.341176471
futurism.com	10/22/2023	Mixed	-0.215686275
futurism.com	10/27/2023	negative	0
futurism.com	11/7/2023	Cruise's safety performance around children has proved to be inadequate contradicting the main hypothesis.	0
futurism.com	11/9/2023	neutral	-0.105882353
futurism.com	11/16/2023	The sentiment towards the main hypothesis is negative.	-0.129411765
futurism.com	1/29/2024	The sentiment towards the main hypothesis is mixed with some perceiving it as a safety improvement and others as inadequate.	-0.090196078

no confirmation bias

source	date	sentiment llm	sentiment classic
theguardian.com	4/20/2022	positive	0.066666667
theguardian.com	8/7/2022	Skeptical	0.043137255
theguardian.com	5/16/2023	The sentiment towards the hypothesis is negative as it involves the theft of sensitive technology.	0.152941176
theguardian.com	7/7/2023	The general sentiment towards the hypothesis seems to be in favor of the safety and environmental concerns raised by driverless cars.	0.356862745
theguardian.com	7/26/2023	Concerned about risks and dominance of self-driving cars supportive of public transportation	0.019607843
theguardian.com	8/10/2023	positive	-0.160784314
theguardian.com	8/14/2023	negative	-0.184313725
theguardian.com	10/31/2023	The general sentiment towards the hypothesis is controversial due to conflicting evidence.	0
theguardian.com	11/8/2023	cautiously optimistic	-0.121568627
		The sentiment towards the hypothesis is cautious optimism as the focus seems to be on	

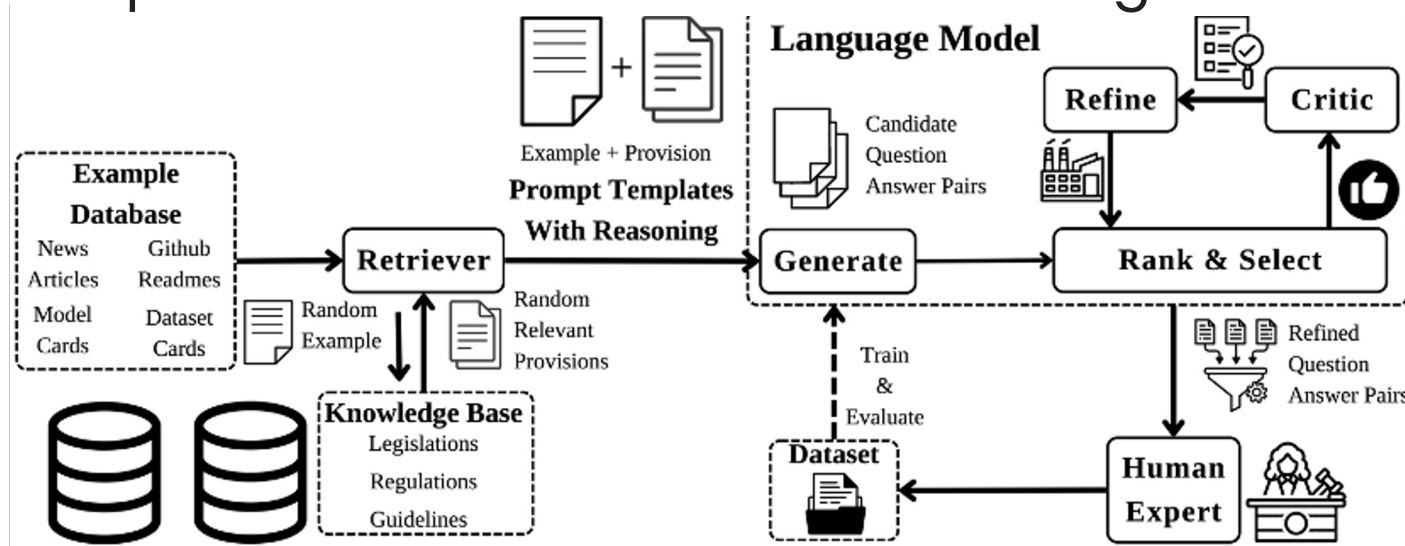


# ELIAS: AI Algorithmic Fairness



**Create an LLM assistant that can check an AI system for compliance with provided legal & ethical guidelines (incl. fairness)**

- ☒ Starting with the dataset
  - ☒ Generating examples for annotation
- ☒ An algorithm for compliance text indexing and querying that preserves structure while retrieving relevant provisions



*Preparing AI for Compliance: Initial Steps of a Framework for Teaching LLMs to Reason About Compliance*

**Makovec, Rei, and Novalija**

*RuleML+RR'24: Companion Proceedings of the 8th International Joint Conference on Rules and Reasoning (To appear)*

# ELIAS/DATAPACT: AI Algorithmic Fairness

## Indexing

1. Split document into pre-defined **articles** and **paragraphs**
2. Further segment paragraphs into overlapping **passages**
3. Encode (**Embed**) each passage using a dense retrieval embedding model

## Querying

1. Encode the query using the same model
2. Compute (dot product) similarity between query and passage **embeddings**
3. Retrieve top-*k* **passages** with scores above threshold *t*
4. Map retrieved passages back to their respective **articles**



## Prompt Generation

1. Use a predefined template incorporating up to *n* **articles**
2. Prompt the **LLM** + *“Let’s think step by step”*

Consider the following articles of legislation, provided between triple backticks, and nothing else:

```
```EU AI Act: Articles: 6, 73, 82, 111, 80, 20, 60, 99``` [text omitted]
```

Under these articles and only these articles and ignoring those that are not applicable, as a legal compliance expert, answer: what are the implications of the provided articles to the following incident, in triple backticks:

```
```a 50-year-old man driving on Autopilot was killed when his Tesla plowed under a semi truck.```
```

Answer article by article and don't provide any additional introduction or conclusions. Let's think step by step.

## ELIAS: AI Algorithmic Fairness

Example

### Article 73: Reporting of Serious Incidents

Should the Autopilot system be classified as a high-risk AI system, the provider (Tesla) would be required to report the incident to the market surveillance authorities within the stipulated time frame if a causal link between the AI system and the accident is established or reasonably suspected.

### Article 82: Compliant AI Systems Which Present a Risk

If after an evaluation, it is found that the Autopilot system, despite compliance with regulations, presents a risk to safety, health, or fundamental rights, the market surveillance authority could mandate Tesla to take appropriate measures to mitigate the risk.

### Article 80: Procedure for Dealing with AI Systems Classified by the Provider as a Not High-Risk in Application of Annex III

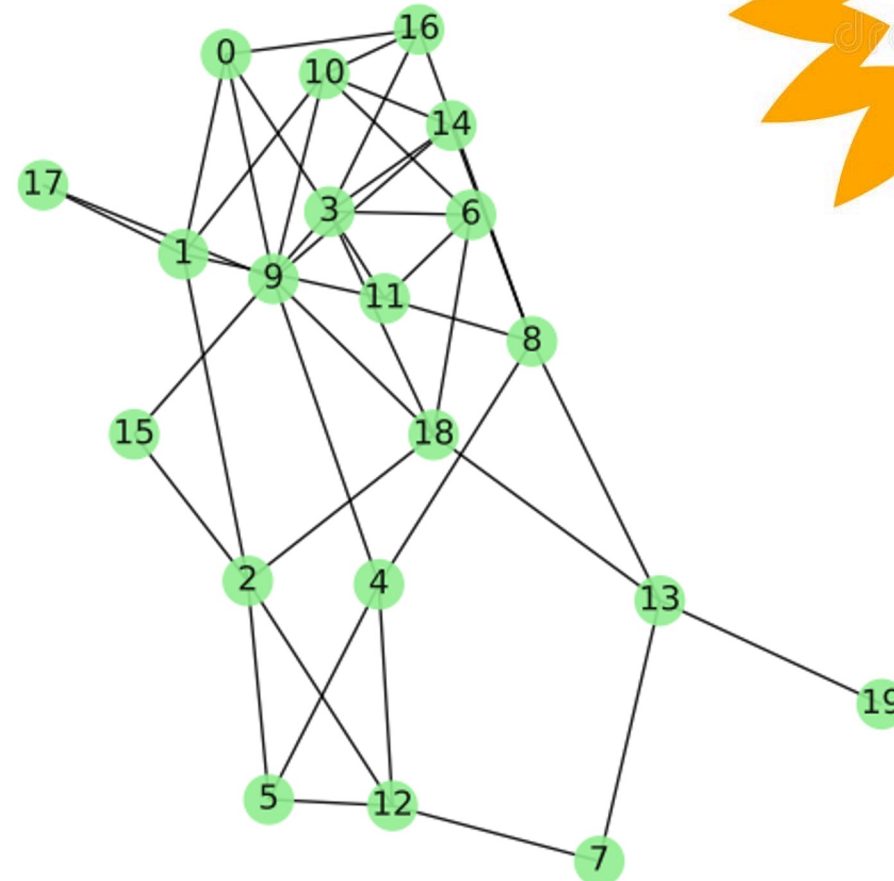
If Tesla had classified the Autopilot system as not high-risk and an evaluation by the market surveillance authority finds otherwise, Tesla would be required to comply with high-risk AI system regulations, including taking necessary corrective actions.

### Article 60: Testing of High-Risk AI Systems in Real World Conditions Outside AI Regulatory Sandboxes

If the Autopilot system is under testing as a high-risk AI system, specific conditions and approvals are required for

# ELIAS: Multi Agent Learning

- ⊠ Regret minimization with  $N$  agents over a communication network  $G$
- ⊠ Agents can only talk to their neighbors (no message passing)
- ⊠ A random number of agents is active at each time step
- ⊠ We study the extent to which communication compensates for the partial feedback
- ⊠ Special cases: full information (experts) and bandit feedback



N. Cesa-Bianchi, T. Cesari, R. Della Vecchia.  
Cooperative Online Learning with Feedback  
Graphs.  
Submitted, 2024.

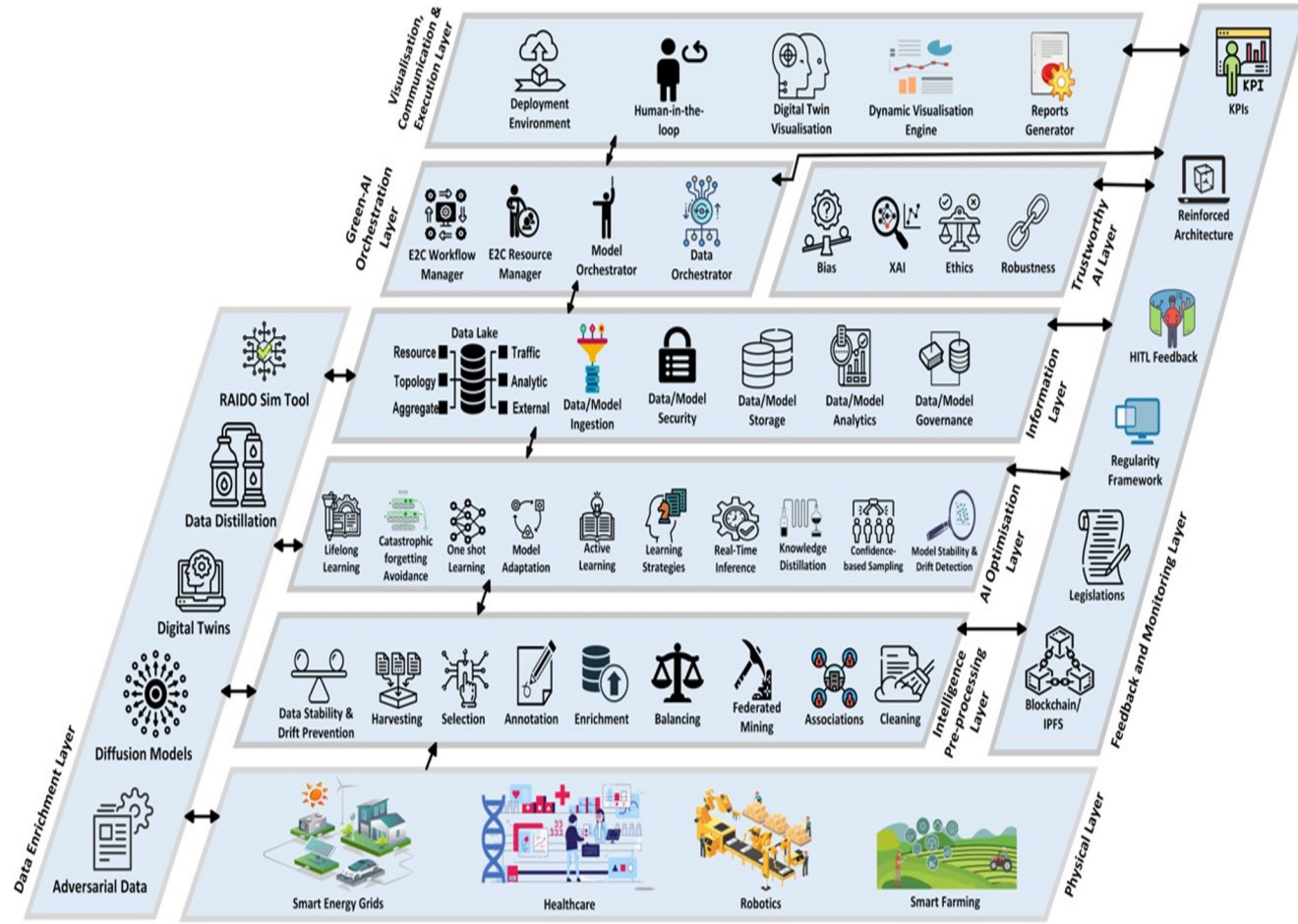


# RAIDO

raido-project.eu

## Reliable AI and Data Optimization

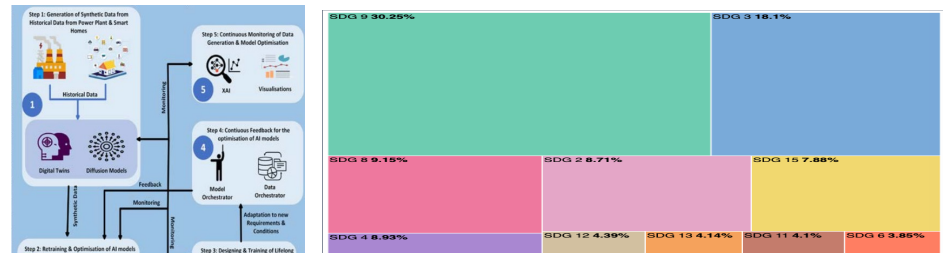
An integrated platform for Green and Energy-efficient data and model related operations



### Green AI Orchestration Bridging Trustworthy AI and Edge AI through tinyML for Frugal Intelligence

Joao Pita Costa<sup>a,\*</sup>, Marco Zennaro<sup>b</sup>, Ioana Ntinou<sup>c,1</sup>, John Shawe-Taylor<sup>a</sup>

<sup>a</sup>International Research Centre on AI under the auspices of UNESCO, Jamova cesta 39, Ljubljana, 1000, Slovenia



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

## AI Ethics & Principles by Design Framework

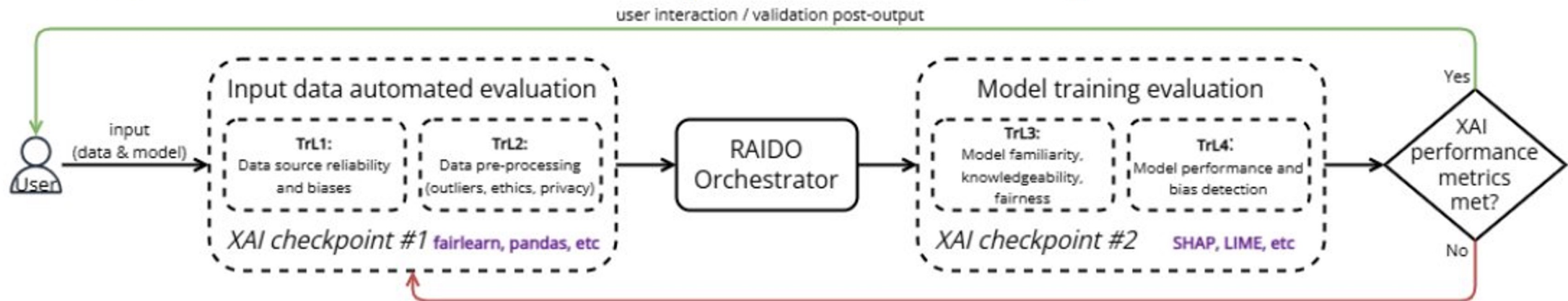
## Reinforced Benchmarking & Feedback-based Progress Monitoring

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

# OECD Bias Tools

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

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 STAKEHOLDER(S)

PURPOSE(S)

**List of tools (921)**

[Eticas Bias](#)

Technical United Kingdom 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(s)**

[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(s)**

[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(s)**

[Transparency](#)

**Related lifecycle stage(s)**

[Operate & monitor](#), [Verify & validate](#)

[COMPL-AI](#)

Technical Switzerland, European Union Uploaded on Jan 24, 2025

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

Show metrics Show use cases

SEARCH

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!

**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:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

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](#)

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. This is particularly important in high-risk environments, e.g. healthcare, finance, insurance and defence. Having community-based prompts to validate models as fit for purpose is safer in a world where LLMs are not static.

The PLIM method consists of question-and-answer prompts that can be applied to specific purposes validated by the community the Prompt-LLM output seeks to support. These prompts are shared widely across sector leads for validation purposes (for example in a healthcare context, this would be senior clinicians, NICE and

[Website](#) [Github](#) [Hugging Face](#)

**About the tool**

You can click on the links to see the associated tools

**Developing organisation(s):** [CarefulAI](#)

**Objectives(s):** [Performance](#)

**Impacted stakeholders:** [Consumers](#) [Regulators](#)

**Country of origin:** [United Kingdom](#)

**Type of approach:** [Educational](#)

**Maturity:** [Implemented in multiple projects](#)

**Usage rights:** [Fee-based](#)

**License:** [Apache 2.0](#)

**Target users:** [Data scientist](#) [Developer](#) [Other](#)

**Group:** [Academia](#) [Business](#) [Government](#)

[Always up to date](#)

**Scope:** [International](#)

ai responsible  
performance  
benchmarking  
llm



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



# AI4GOV Bias Tools

AI4Gov Platform

- Home
- Use Cases
- Bias Detector Toolkit
  - Introduction
  - Bias Detector Catalogue
  - Rare Diseases Bias Results

## 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 o 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
Source: <a href="https://github.com/KenSciResearch/fairMLHealth">https://github.com/KenSciResearch/fairMLHealth</a> Type: MITIGATION Programming Language: PYTHON Description: FairMLHealth is a healthcare-specific tool for bias analysis. It provides machine-learning fairness, healthcare applications, and variation analysis. Applicability: HEALTHCARE Limitations: The 'fair' range to be used for these metrics requires judgement on the part of the analyst. References: Ahmad et al., (2020). Fairness in Machine Learning for Healthcare. Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. <a href="https://doi.org/10.1145/3394486.3406461">https://doi.org/10.1145/3394486.3406461</a> .		
<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.eu/projects.net/>



# AI4GOV MOOC



## Course learning outcomes

- Understand the concepts of artificial intelligence, machine learning, and deep learning.
- Recognize the transformative potential and ethical implications of AI in various domains.
- Identify the types and sources of bias that can manifest in AI systems.
- Understand the real-world consequences of biased AI.
- Learn how data collection and pre-processing can introduce bias into AI models.
- Implement best practices for collecting, cleaning, and preparing data to reduce bias.
- Study and understand regulatory frameworks governing AI, such as GDPR and AI Act.

Course description	<b>Course content</b>	Course reviews
--------------------	-----------------------	----------------

## Course content

Below is the course content. You can click on any section here and it will take you through to this section of the course. If you are signed in and enrolled on this course we can track your progress.

Course progress  50% completed

<b>Module 1: Introduction to Trustworthy and Democratic AI (Fundamentals)</b>	>
<b>Module 2: Bias in AI - Understanding Bias</b>	>
<b>Module 3: Data and Bias</b>	>
<b>Module 4: Ethical AI Governance</b>	>
<b>Resources &amp; Authorship</b>	>



<https://ai4gov-project.eu/home/resources/>





**unesco**

Centre  
Under the auspices  
of UNESCO



International Research Centre  
of Artificial Intelligence  
under the auspices of UNESCO

**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/>