# The CiTIUS-USC Z-Inspection® Lab

## "How to Evaluate, Audit, Certify... AI systems?"

**Jose María Alonso Moral**
**josemaria.alonso.moral@usc.es**

*Workshop on Social Impact of Artificial Intelligence*

*30 May 2024, Fundación Ramón Areces, Madrid*

# CITIUS

**Centro Singular de Investigación en Tecnoloxías Intelixentes**

USC — UNIVERSIDADE DE SANTIAGO DE COMPOSTELA

citius.gal

## SCIENTIFIC AREAS

**Green and smart computing and devices**
- Electronic design of intelligent devices
- High Performance Computing

**Support for machine intelligence**
- Data and process science and engineering
- Automatic learning and reasoning

**Intelligent technologies**
- Virtual and augmented reality
- Language technologies
- Computer vision
- Robotics

**Social, economic, ethical and political framework**
- Trustworthy AI

## R&D RESULTS 2021-23

**ARTICLES**
- **75**% in Q1
- **44**% in D1
- **27**% in top 10% most cited

**11,8** M€
- **21**% income from european projects

## TEAM

**+150 PEOPLE**
31 SENIOR RESEARCHERS

## SOME EUROPEAN PROJECTS

**NANOVR**
ERC Consolidator Grant
Nanoscale Design using Virtual Reality

**NL4XAI**
Interactive Natural Language for Explainable AI (Coordinator)

**MENELAOS NT**
Multimodal Fusion of Sensor Information (Coordinator)

**hybrids**
Hybrid intelligence to monitor, promote and analyse transformations in good democracy practices (Coordinator)

**MISEL**
Multispectral Intelligent Vision System with Embedded Low-Power Neural Computing

**iRead4Skills**
Intelligent Reading Improvement System for Fundamental and Transversal Skills Development

**SOS food**
Sustainability Optimization for Secure Food Systems

**EDiH DATAlife**
DIH for the deployment of AI and Data Analytics in SMEs in the primary, biotechnological and health sectors

## SOME STRATEGIC PROJECTS

**nós**
Artificial Intelligence at the service of the Galician Language

**CAMELIA**
Plexus Tech - USC Chair in AI applied to personalised precision medicine

Televés - USC Chip Chair in microelectronic design

## TECHNOLOGY TRANSFER

**+30** contracts with industry

**85**% licences to international companies

**FINSA** soluciones en madera
Machine learning for digital twins and processes optimisation in industry

**Avincis** AVIATION CRITICAL SERVICES
Efficient processing of aerial images within the Civil UAVs Initiative

**Lacer**
Rapid analysis of the development of dental plaque

## SPIN-OFFS

**situm** INDOOR POSITIONING

**inverbis**

## CONGRESS

**ECAI** 50th ANNIVERSARY

# How to certify AI ?
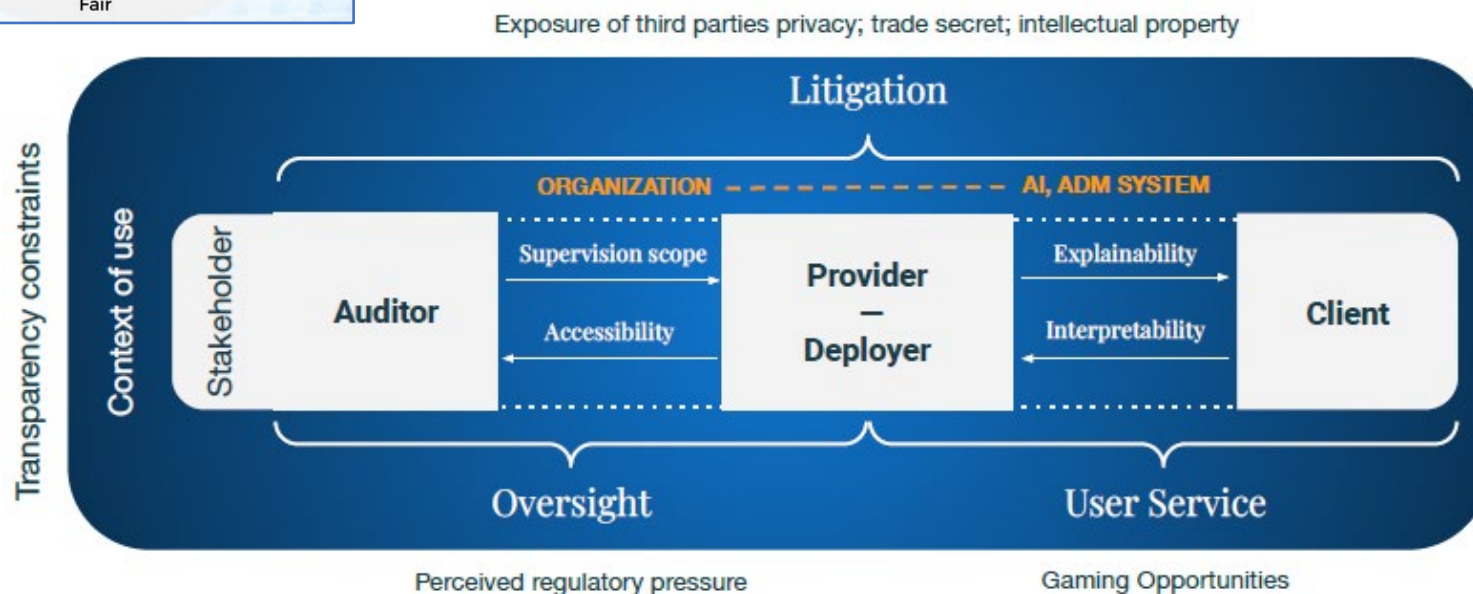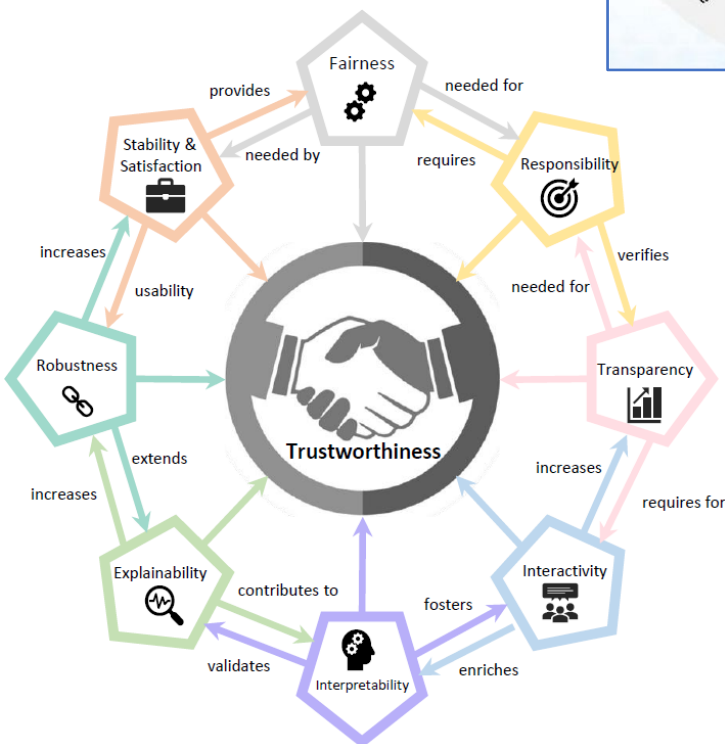
# Certification



L. Nannini, J. M. Alonso-Moral, A. Catala, M. Lama, S. Barro, "**Operationalizing Explainable AI in the EU Regulatory Ecosystem**", IEEE Intelligent Systems, 2024, https://doi.org/10.1109/MIS.2024.3383155
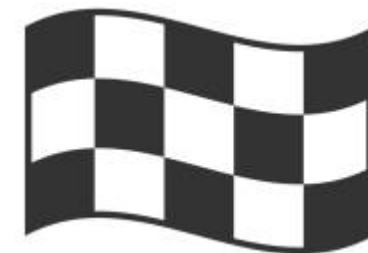
S. Alia, T. Abuhmed, S. El-Sappagh, K. Muhammad, Jose M. Alonso-Moral, R. Confalonieri R. Guidotti, J. Del Ser, N. Díaz-Rodríguez, F. Herrera, "**Explainable Artificial Intelligence (XAI): What we know and what is left to attain Trustworthy Artificial Intelligence**", Information Fusion, 2023, https://doi.org/10.1016/j.inffus.2023.101805
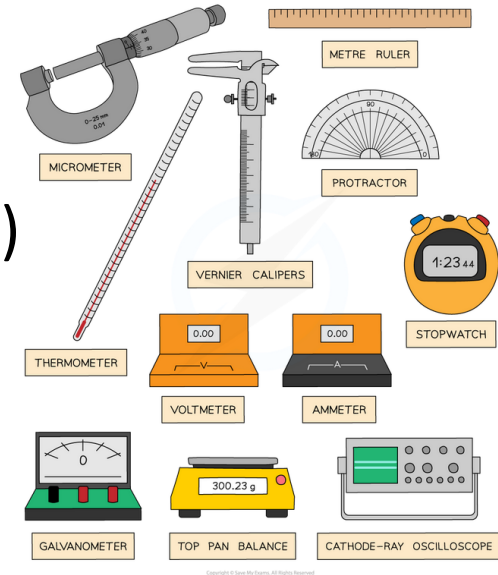
# How to audit and evaluate AI ?

☐ **Metrology is the Science of Measurement and its application**

➢ To measure something, we need a Metric System (i.e., a standard reference and a related unit which is defined without any ambiguity)

☐ **SI: International System of Units** (accepted worldwide)

✓ Time is measured in seconds, minutes, hours…

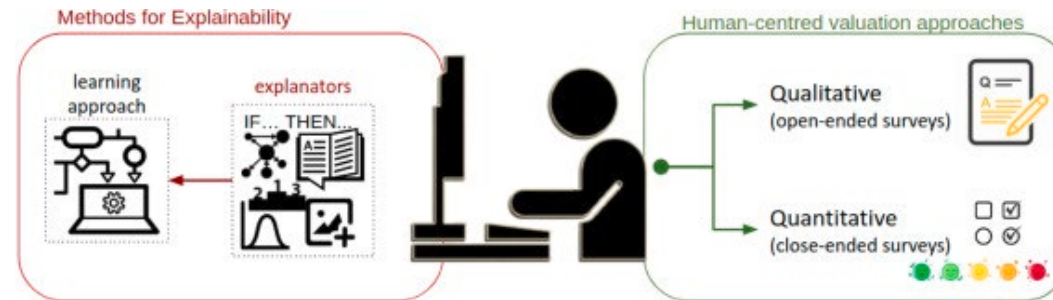✓ Length is measured in meters, inches…

✓ Mass is measured in Kg, Pounds…

**Human Evaluation**

- ✓ **Interviews**
- ✓ **Focus Groups**
- ✓ **Free-text Questionnaires**

- ✓ **Panel of Experts**
- ✓ **Consumer Testing**
- ✓ **Randomized Control Tests**

Methods for Explainability

learning approach    explanators
IF… THEN…

Human-centred valuation approaches

Qualitative (open-ended surveys)

Quantitative (close-ended surveys)

- ✓ **Complexity Metrics**
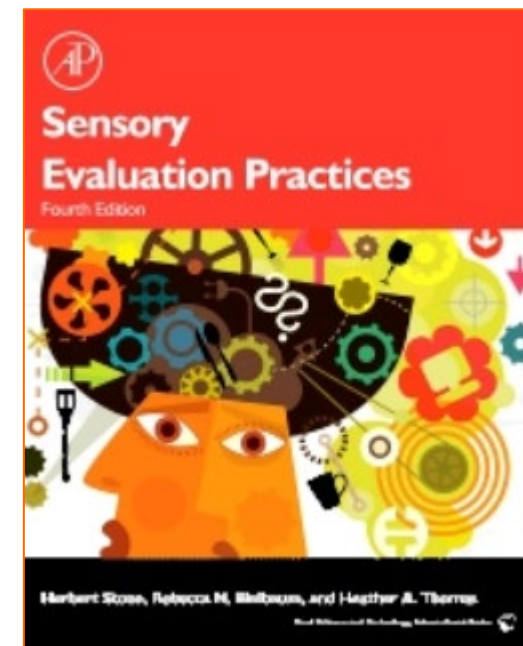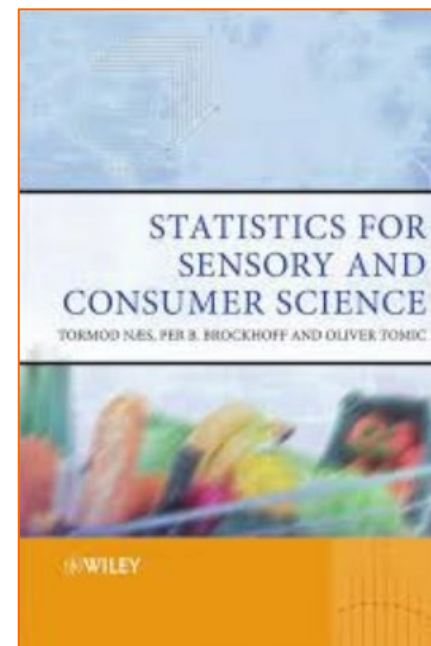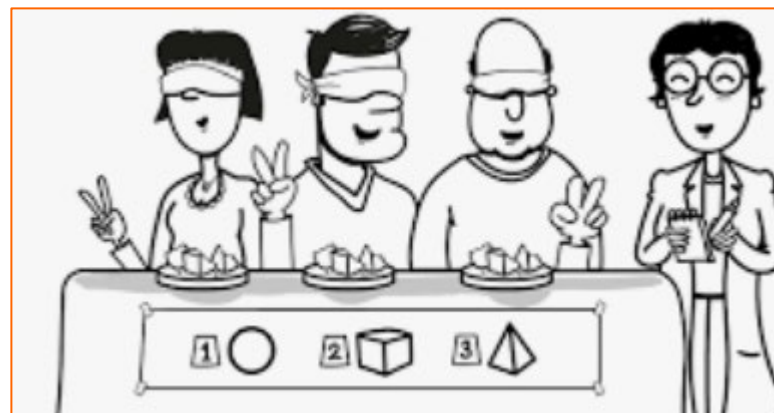- ✓ **Transparency Metrics**
- ✓ **Interpretability Metrics**

**Automatic Evaluation**

# How to measure human-centric & subjective properties ?

❑ **What about subjective judgements (e.g., beauty, elegance, transparency,…) ?**

❑ **Which is the standard reference for sensory inputs (acid, sweet, salty, easy to read…) ?**

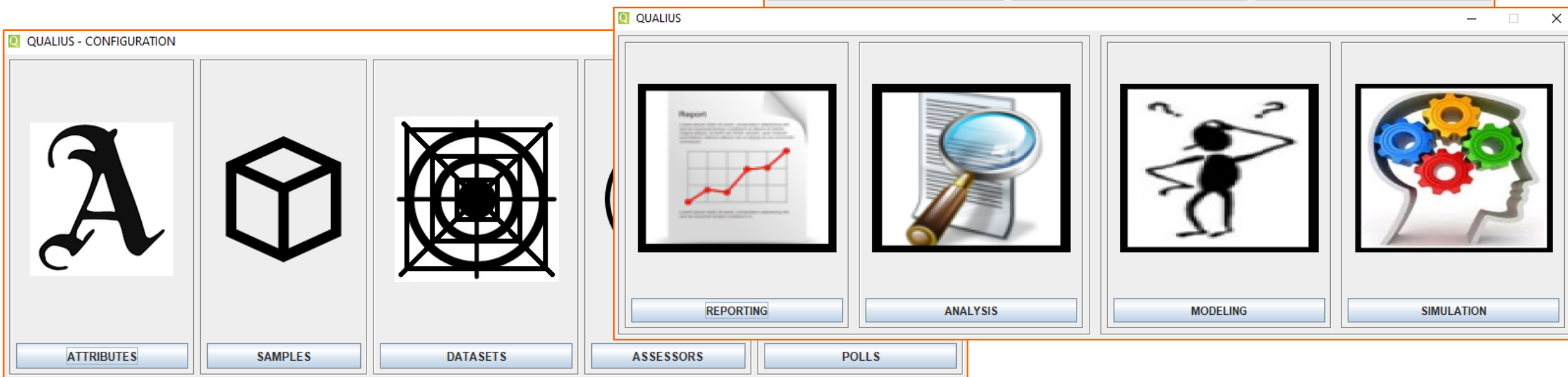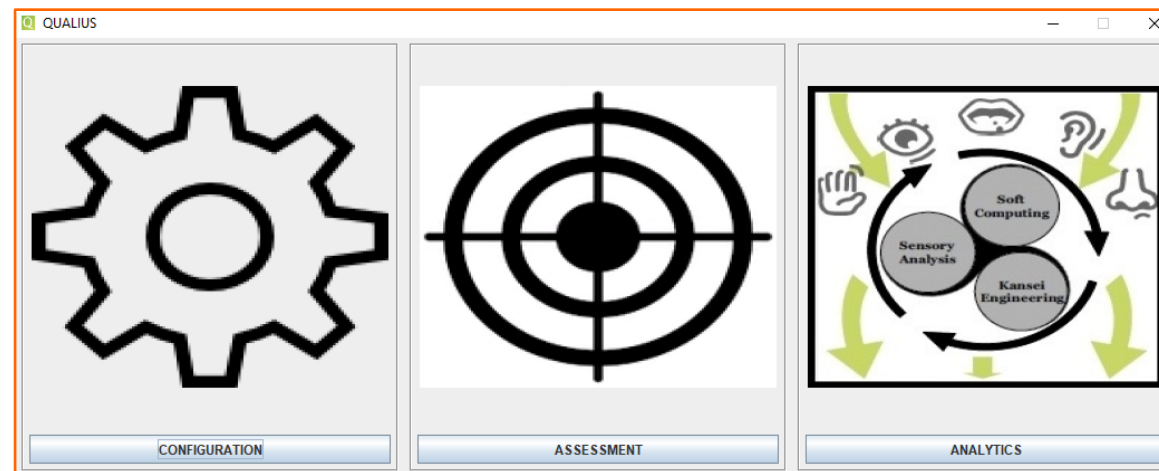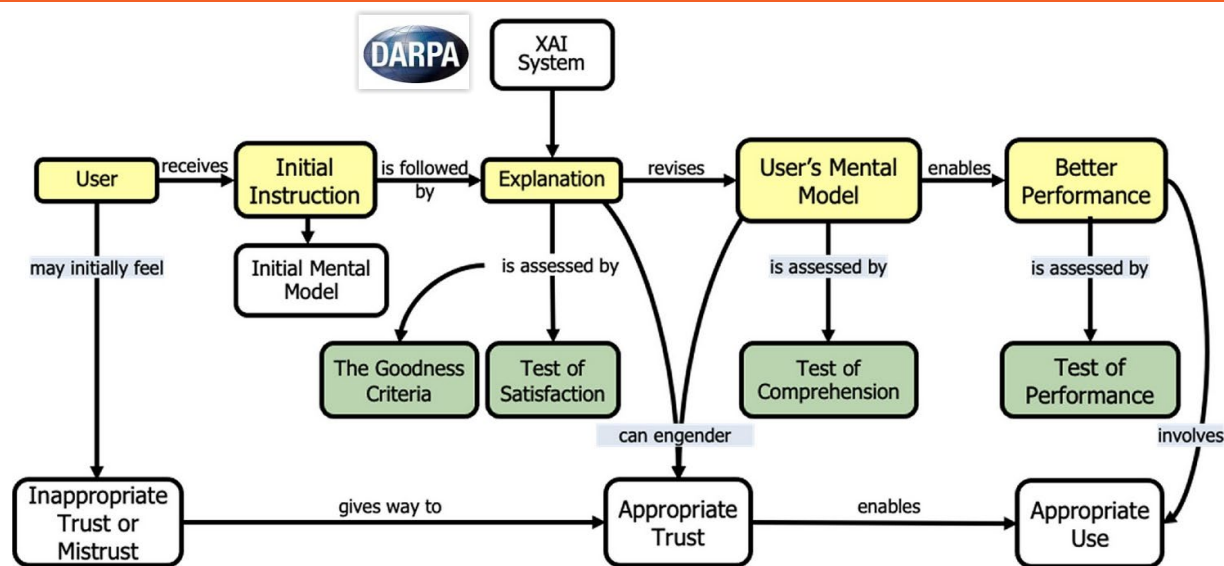❑ **How to avoid confounding effects ?**

# QUALIUS

P. Quiros, J.M. Alonso, D.P. Pancho, **"Descriptive and comparative analysis of human perceptions expressed through fuzzy rating scale-based questionnaires"**, International Journal of Computational Intelligence Systems, 9(3):450-467, 2016, https://doi.org/10.1080/18756891.2016.1175811

J.M. Alonso, D.P. Pancho, L. Magdalena, D.A. Nuñez, D.S. Sánchez, P.F. Suárez, J. Mingot, V. Iglesias, **"QUALE: A new Toolbox for Quantitative and Qualitative Analysis of Human Perceptions"**, IFSA-EUSFLAT, 2015, https://doi.org/10.2991/ifsa-eusflat-15.2015.94



QUALIUS

CONFIGURATION | ASSESSMENT | ANALYTICS

QUALIUS

REPORTING | ANALYSIS | MODELING | SIMULATION

QUALIUS - CONFIGURATION

ATTRIBUTES | SAMPLES | DATASETS | ASSESSORS | POLLS

# Guiding Human Evaluation in the context of Explainable AI



D. Gunning, E. Vorm, J. Y. Wang, M. Turek, "**DARPA's explainable AI (XAI) program: A retrospective**," Applied AI Letters, 2(4):e61, 2021 https://doi.org/10.1002/ail2.61

R.R. Hoffman, S.T. Mueller, G. Klein, J. Litman, "**Metrics for Explainable AI: Challenges and Prospects**", https://doi.org/10.48550/arXiv.1812.04608



R. Confalonieri, J.M. Alonso-Moral, "**An operational framework for guiding human evaluation in Explainable and Trustworthy AI**", IEEE Intelligent Systems, 2023, https://doi.org/10.1109/MIS.2023.3334639
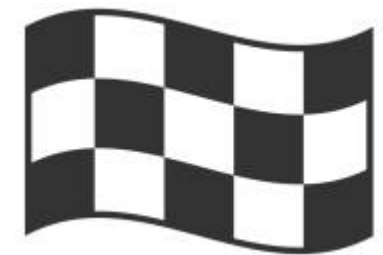


https://reprohum.github.io/

A. Belz et al., "**Missing information, unresponsive authors, experimental flaws: The impossibility of assessing the reproducibility of previous human evaluations in NLP**", Proceedings of the 4th ACL Workshop on Insights from Negative Results in NLP, 2023, https://aclanthology.org/2023.insights-1
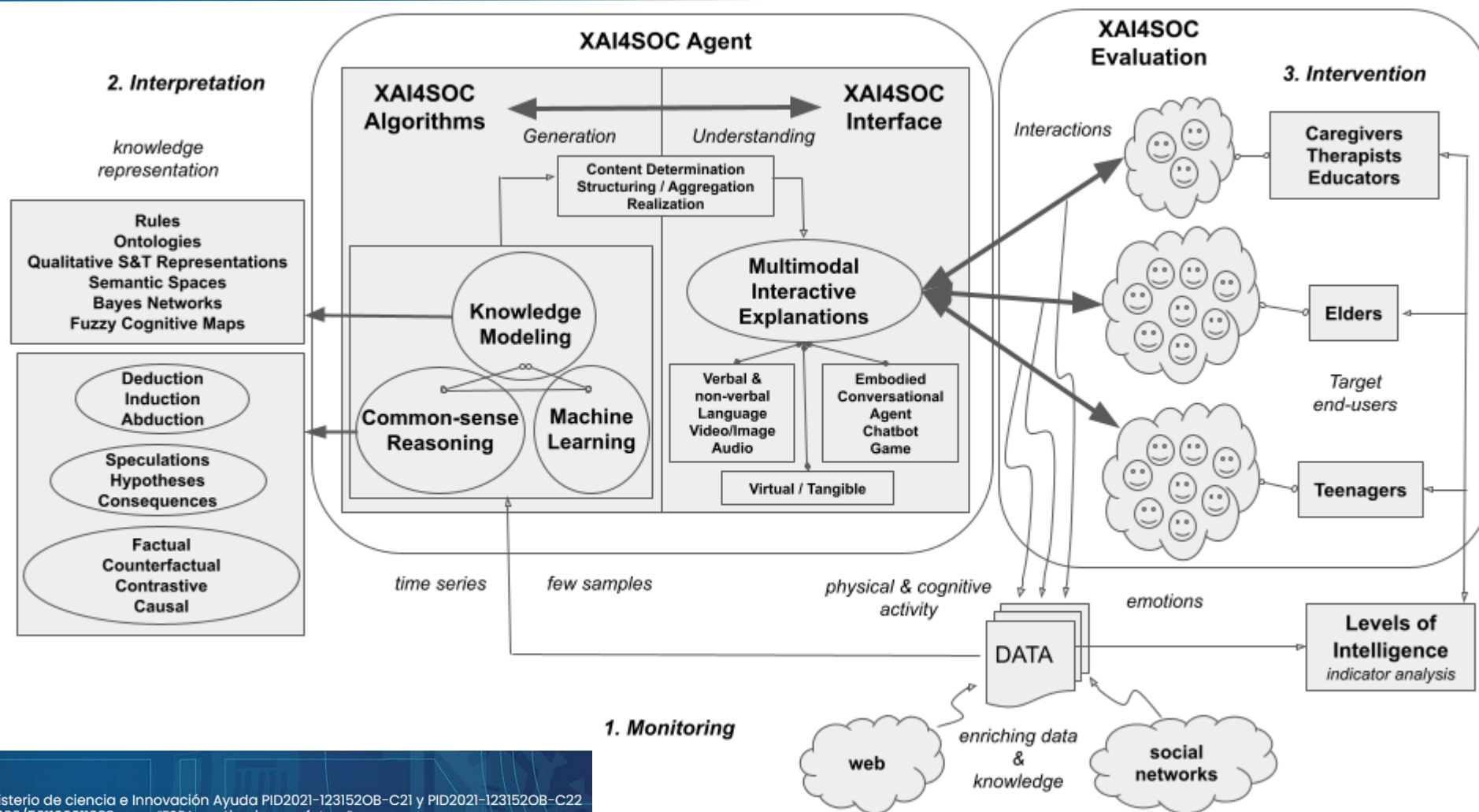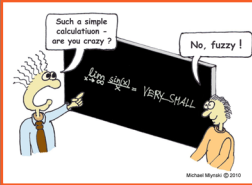
# Use Case

**XAI4SOC will create an open-source framework for the design, development, validation, evaluation and deployment of a new generation of ubiquitous XAI agents (algorithms and interfaces) for enhancing the prominence of their target users, with the support of caregivers, therapists and educators, in pursuing a healthy ageing and wellbeing development**
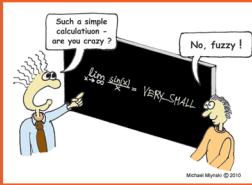
https://xai4soc.es/

Fragment of a real conversation between TimeVersa and a user. Underlined, temporal expressions relevant for the convsration.

| Move | Participant | Utterance |
|------|-------------|-----------|
| $m_1$ | User | Hello! I would like to be in Cíes Islands around 12 PM. What time should I take the boat from Vigo? |
| $m_2$ | TimeVersa | There is a boat at 11.15 AM that is suitable for you. |
| $m_3$ | User | Why do you think this boat is suitable for me? |
| $m_4$ | TimeVersa | The boat trip from Vigo to Cíes Islands lasts about 35 minutes, so if you take the boat at 11.15 AM you can be there around 12 PM. |
| $m_5$ | User | Could you specify what "around 12 PM" means for you? |
| $m_6$ | TimeVersa | It means between 11.50 AM and 12.10 PM. |
| $m_7$ | User | OK. Book the boat for me, please. |
| $m_8$ | TimeVersa | Could I help you with anything else? |
| $m_9$ | User | Yes. Once we have reached Cíes Islands around 12 PM, we had thought of doing the hiking route "Faro de Cíes". Do you think we can be back in Vigo before 2 PM? |
| $m_{10}$ | TimeVersa | No, this is not possible. |
| $m_{11}$ | User | Why not? |
| $m_{12}$ | TimeVersa | Because the route will take you about 3 hours. If you arrive the Islands around 12 PM, you can't be back in Vigo before 2 PM. |
| $m_{13}$ | User | OK, I understand. So, what time can we be back? |
| $m_{14}$ | TimeVersa | You can be back in Vigo around 4.35 PM if you take the boat at 4 PM. |

https://tec.citius.usc.es/survey/

https://www.prolific.com/

# Evaluating the consistency of explanations by LLMs



https://tec.citius.usc.es/survey/

M. Canabal-Juanatey, J.M. Alonso-Moral, A. Catala, A. Bugarín-Diz, **"Enriching Interactive Explanations with Fuzzy Temporal Constraint Networks"**, International Journal of Approximate Reasoning, 109128, 2024, https://doi.org/10.1016/j.ijar.2024.109128

# Evaluating ...LMs

Below, you can find a conversation between a fictitious user, Mary, and a conversational agent oriented to book hostels and provide information about the stages of the French St James' way. Mary is interested in the last two stages of the way and booking a hostel in the village "Pedrouzo".

The agent knows the following information about the last two stages:

⏱ 07H 00'          ⏱ 06H 30'

Arzúa          Pedrouzo          Santiago de Compostela

And the agent knows the following information about the available hostels in Pedrouzo:

**Peregrina**
33 euros per night
Pilgrims should arrive between 11:00 AM and 6.30 PM and leave it before 8.30 AM.

**Triskel**
55 euros per night
No constraints about arrival and leaving times.

**O Burgo**
20 euros per night
Pilgrims should arrive between 11:00 AM and 6.30 PM and leave it before 8.30 AM.

**Xacobeo**
25 euros per night
No constraints about arrival and leaving times.

Now, please read carefully the following conversation between Mary and the agent and answer the question below.

Mary:
Hello! We are planning to do the two latest stages of French St. James' way next weekend, so we would like to book a hostel in Pedrouzo for Saturday.

Agent:
The hostel "Peregrina" in Pedrouzo seems suitable for you.

Mary:
Could you provide us more information about it?

Agent:
Sure! The price is 33 euros per night, and you should enter between 11:00 AM and 6:30 PM and leave it before 8:30 AM on Sunday. Would you like to book it?

Mary:
Yes. Book it for me, please.

Agent:
The hostel has been successfully booked. Can I help you with anything else?

Mary:
Yes. We would like to visit the village for a couple of hours before arriving at the hostel. What time should we start the stage to ensure we have enough time?

Agent:
You could start the stage from Arzúa to Pedrouzo at 8.30 AM.

Mary:
Perfect. So, we plan to start the last but one stage from Arzúa to Pedrouzo at 8:30 AM. We will also need to rest for at least 16 hours between this stage ending and the next stage starting, from Pedrouzo to Santiago de Compostela. Do you think we can arrive in Santiago de Compostela before 12:00 PM?

Agent:
No, this is not possible.

Based only on the temporal information involved in the conversation, please assess the consistency of the agent's responses.

○ Very inconsistent    ○ Inconsistent    ○ Somewhat consistent    ○ Consistent    ○ Very Consistent

* Required

Next

https://tec.citius.usc.es/survey/

...sational agents

..., "Reducing, Reusing and Recycling large pre-...dge Representation and Reasoning". This PhD ... de Tecnologías Inteligentes da Universidad de ... under the supervision of the researchers Alberto

...pleted. In this case, they will only be used for the

...rina.canabal.juanatey@usc.es).

...ejudice at any time.
...ocessing and publication for the purpose of this
...'9 of the European Parliament and of the Council

...answered before starting the experimental tasks.

**...zzy Temporal**
...r.2024.109128

...orthy AI

### Instructions

In the survey you will fi...
like the one in the ex...
screen, you can se...
conversation between ...
conversational agent o...
hostels and provide inf...
the stages of the French...

Before to show the con...
relevant information av...
agent to the elaboration...
will be displayed: a...
hostels and duration o...
interest. This informati...
same in the three sc...
conversation below will c...

Finally, you will be asked...
question about the con...
the agent's responses...
correctly, you should...
account the tempor...
included in the conversat...

**VERY IMPORTANT**: be...
the questions, it is es...
read the conversation be...
the agent carefully, even...
if necessary, to make sur...
understood each senten...

M. Canaba...
**Constraint...**

Such a simple calculatiuon – are you crazy ?
No, fuzzy !
from logical to VERY-SMALL
Michael Mysski © 2010

## 130 subjects



(a) Gender.

(b) Age.

(c) Country of residence.

(d) The highest education level completed.

**Experimental design:** two factors as independent variables at two levels each (2x2)



**H1:** "In case of involving vague temporal information, the interaction with TimeVersa is perceived as more consistent than the interaction with the OpenAI API."

**H2:** "In case of involving vague temporal information, the interaction with TimeVersa incorporating the temporal reasoning model is perceived as more consistent than the interaction with TimeVersa incorporating the crisp version of the temporal reasoning model."

M. Canabal-Juanatey, J.M. Alonso-Moral, A. Catala, A. Bugarín-Diz, **"Enriching Interactive Explanations with Fuzzy Temporal Constraint Networks"**, International Journal of Approximate Reasoning, 109128, 2024, https://doi.org/10.1016/j.ijar.2024.109128

# NL4XAI

Interactive *Natural Language* Technology for eXplainable Artificial Intelligence

Supporting *the right to explanation by* AI POWERED SYSTEMS

NL4XAI.eu

# Jose Maria Alonso-Moral

iD https://orcid.org/
0000-0003-3673-421X

## Profesor Titular de Universidad

GSi · USC UNIVERSIDADE DE SANTIAGO DE COMPOSTELA · CAMPUS VIDA CAMPUS DE EXCELENCIA INTERNACIONAL · CITIUS

josemaria.alonso.moral@usc.es
https://citius.gal/team/jose-maria-alonso-moral
https://gitlab.citius.usc.es/jose.alonso/xai

### International Ambassador

FM square

https://fmsquare.org/

IEEE Computational Intelligence Society

Task Force on Explainable Fuzzy Systems

https://sites.google.com/view/tf-explainable-fuzzy-systems/

Z-Inspection®

https://z-inspection.org/

HYBRIDS MARIE CURIE ACTIONS

Combating Disinformation and Abusive Language with Human and Artificial Intelligence

MINISTERIO DE CIENCIA E INNOVACIÓN · AGENCIA ESTATAL DE INVESTIGACIÓN

**XAI for Healthy Aging and Social Wellbeing**

XAI 4 SOC
EXPLAINABLE ARTIFICIAL INTELLIGENCE FOR HEALTHY AGING AND SOCIAL WELLBEING

https://xai4soc.es/

MINISTERIO DE CIENCIA E INNOVACIÓN · AGENCIA ESTATAL DE INVESTIGACIÓN

**TELSEC4TAI**   https://iafiable.es/

**R&D Network on Trustworthy Artificial Intelligence: Technical, Ethical, Legal, Cultural and Socio-economic Challenges**

**DeepR3: Reducing, Reusing, and Recycling large models for developing Responsible and Green Language Technologies**

MINISTERIO DE CIENCIA E INNOVACIÓN · Financiado por la Unión Europea NextGenerationEU · Plan de Recuperación, Transformación y Resiliencia · AGENCIA ESTATAL DE INVESTIGACIÓN

http://ixa2.si.ehu.eus/deepr3/

https://hybridsproject.eu/