Julia Palma
julia.palma@ucd.ie

# Julia Palma

Multidisciplinary academic background with focus on **European innovation projects, public policy analysis, public affairs, and policy development** (Universidad de Sevilla, LaSalle, LSE & Cambridge).

**+300 collaborative projects** integrating research, innovation, and policy impact across digital and AI domains.

Leads **policy-focused initiatives** at Ireland's Centre for AI, contributing to **EU and national reports, consultations, and recommendations**.

Former participant in the **European Commission's Digitising European Industry** WG policy initiative.

Represents Ireland in **European Health Data Space (EHDS)** and broader **EU data space policy initiatives**.

Former chair of **BDVA's Data i-Spaces subgroup.**

Current lead of **BDVA's etami (Ethical and Trustworthy Artificial and Machine Intelligence)** subgroup.

Member of the Steering Committee of the **Global AI Policy Research Network (GlobAIpol)**

Member of **EU and international policy groups**, including the **AI Act Regulatory Sandboxes Expert Group, Apply AI Alliance, EU Health Policy Platform**, and **AI Transparency Code of Practice Working Group**.

Julia Palma
julia.palma@ucd.ie

# CeADAR
# Ireland's Centre for AI

Not-for-profit established in 2013

Funded by IDA Ireland and Enterprise Ireland

Based in University College Dublin

Julia Palma
julia.palma@ucd.ie

# CeADAR's Mission

CeADAR supports companies and organisations across Ireland to explore, experiment, and build innovative AI solutions into processes and products, transforming productivity, competitiveness, digitalisation and sustainability.
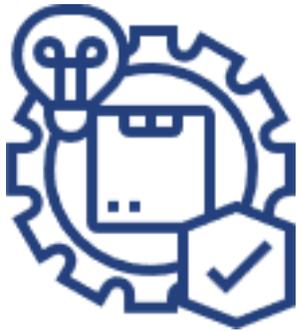
Julia Palma
julia.palma@ucd.ie

Julia Palma
julia.palma@ucd.ie

CeADAR
Ireland's Centre for AI

# The 7 Pillars of Trustworthy AI (HLEG)

Source: Ethics Guidelines for Trustworthy AI. Independent High-Level Expert Group on Artificial Intelligence. European commission, 8 April, 2019.

Julia Palma
julia.palma@ucd.ie

CeADAR
Ireland's Centre for AI

# The 7 Pillars of Trustworthy AI (HLEG) Requirements and sub-requirements

1 **Human agency and oversight:** fundamental rights

2 **Technical robustness and safety:** resilience to attack and security, fall back plan and general safety, accuracy, reliability and reproducibility

3 **Privacy and data governance :**respect for privacy, quality and integrity of data, and access to data

4 **Transparency:** traceability, explainability and communication

5 **Diversity, non-discrimination and fairness:** avoidance of unfair bias, accessibility and universal design, and stakeholder participation

6 **Societal and environmental wellbeing:** sustainability and environmental friendliness, social impact, society and democracy

7 **Accountability** auditability, minimization and reporting of negative impact, trade-offs and redress.

Julia Palma
julia.palma@ucd.ie

CeADAR
Ireland's Centre for AI

# THE LEGAL REALITY: THE EU AI ACT

- **Regulation:** Moving from "Ethics Guidelines" to "Enforceable Law."

- **Risk-Based Approach:** Systems need to meet proportionate and adequate requirements based on the applications (e.g. strict documentation and transparency standards.)

- **KBS Advantage:** Explicit knowledge representation simplifies compliance (Traceability by design) to meet the **"Data Governance"** and **"Technical Documentation"** requirements of the Act.

Julia Palma
julia.palma@ucd.ie

# THE EU AI ACT

The **EU AI Act** (passed in 2024) provides the legal enforcement using a **Risk-Based Approach:**

- **Unacceptable Risk:** Prohibited.
- **High Risk:** Strict compliance, logging, and human oversight required.
- **Limited/Minimal Risk:** Transparency obligations (e.g.,chatbot disclosing it's an AI).

Julia Palma
julia.palma@ucd.ie

# IMPACT ASSESSMENT FOR HIGH-RISK AI SYSTEMS

EU AI Act: Article 27: Fundamental Rights Impact Assessment for High-Risk AI Systems

- *Deployers prior to deploying AI of high-risk, shall perform an assessment of the impact on fundamental rights that the use of such system may produce.*

Julia Palma
julia.palma@ucd.ie

# Where's the line

LEGAL: A **fundamental rights-based approach** is more closely linked to **existing law** and focuses on aspects that are legally relevant and thus **enforceable**.

ETHICAL: Compared to this, an **ethics-based approach** is much broader and also more open to reflection on **potential** implications that may not be worth considering from a legal perspective.
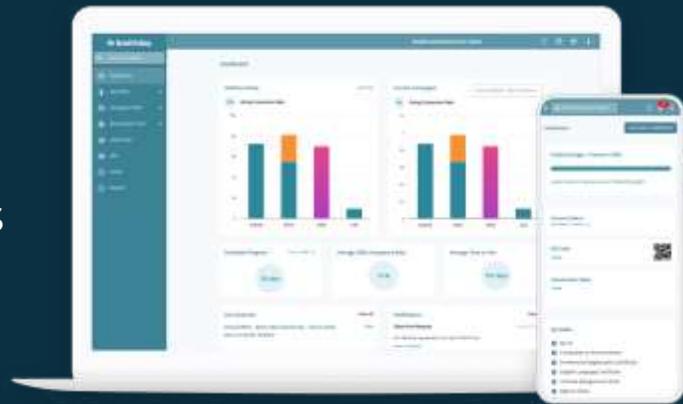
Julia Palma
julia.palma@ucd.ie

CeADAR
Ireland's Centre for AI

# EDIH

## European Digital Innovation Hubs Network

# First comes the Regulation; then, Trust Awareness

## 80% companies reaching out for 'compliance' and obligations of 'transparency'

**healthdaq**

LLMs to parse candidates

Human-oversight

An Roinn Fiontar, Turasóireachta agus Fostaíochta
Department of Enterprise, Tourism and Employment

ONLINE TRAINING

AI for You: Introduction to AI and The EU AI Act

CeADAR
Ireland's Centre for AI

Julia Palma
julia.palma@ucd.ie

# Embedding Trustworthy AI across the innovation cycle

CeADAR
Ireland's Centre for AI



Figure 2: MANOLO Z-Inspection® Co-Design Framework

## Ethical Issues and Tensions

User privacy, safety and data protection vs Transparency

Accountability

User autonomy vs technical robustness/bias/justice
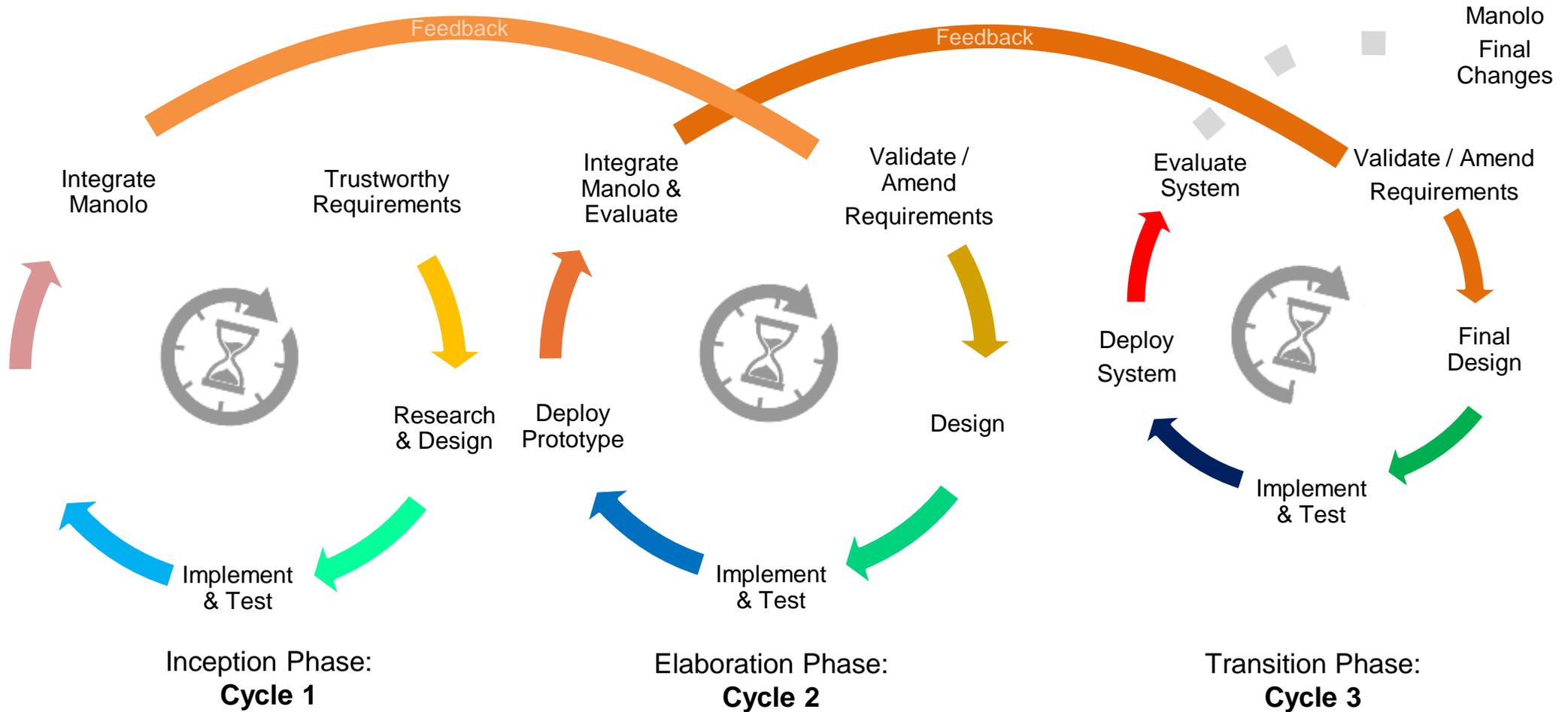
## Regulatory Issues and Tensions

Compliance with GDPR, AI Act, Cybersecurity, domestic/sectoral regulation

## KPI Assessment

## Mitigation Measures

Aligned technological innovation with ethical considerations, ensuring transparency, privacy, inclusivity, and accountability based on stakeholders' feedback
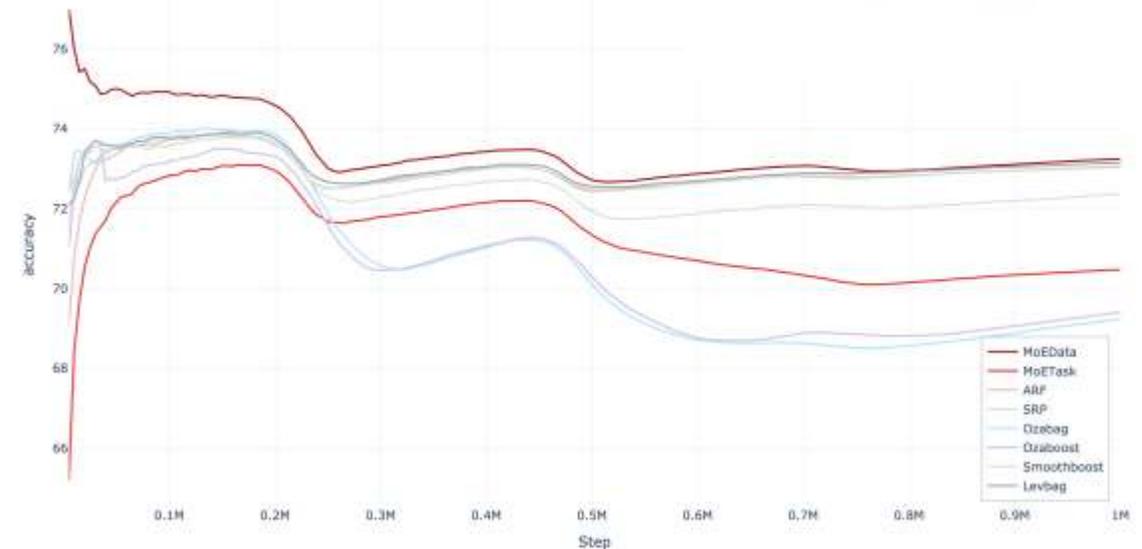
Integration Approach & Mapping

Julia Palma
julia.palma@ucd.ie

## MoE-drift (Mixture of Experts Approach to Handle Concept Drift).

**The Problem:** Concept drift (when real-world data changes and models fail).

**Function:** A framework using a neural router to detect when specific "expert" nodes in a system are no longer accurate due to data shifts.

**Impact:** Ensures **Technical Robustness** in non-stationary environments.



Accuracy over time plot for LED$_g$ dataset

Julia Palma
julia.palma@ucd.ie

CeADAR
Ireland's Centre for AI

# **UPCAST** is a Horizon Europe project where CeADAR developed the **AI Trustworthiness Assessment Tools and Dashboard**
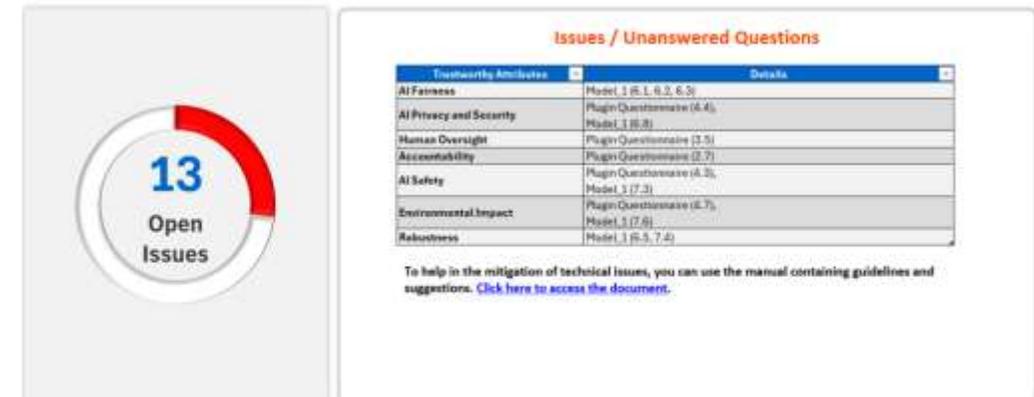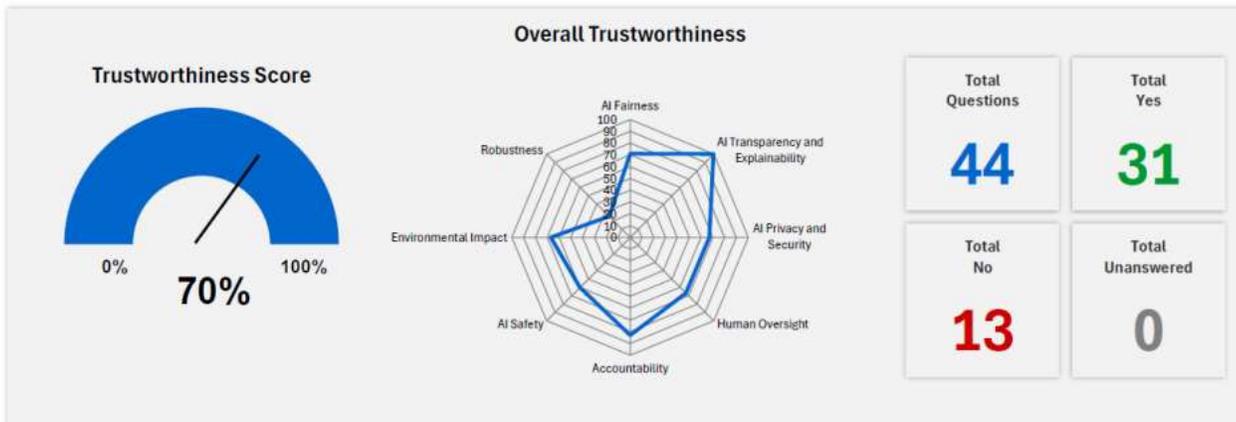


Figure 6: Open Issues

Julia Palma
julia.palma@ucd.ie

# UPCAST Valuation Plug-in

- **The Principle:** Fairness and Accountability. Transparency.
- **Function:** for data providers - setting fair prices for their data products; for data consumers - estimate costs for data acquisition. Using explainable AI techniques (XAI), also provides transparency by identifying the key factors influencing the price.
- **Impact:** Prevents power asymmetries in the Common European Data Spaces.



Figure: Overall plugin architecture

Julia Palma
julia.palma@ucd.ie
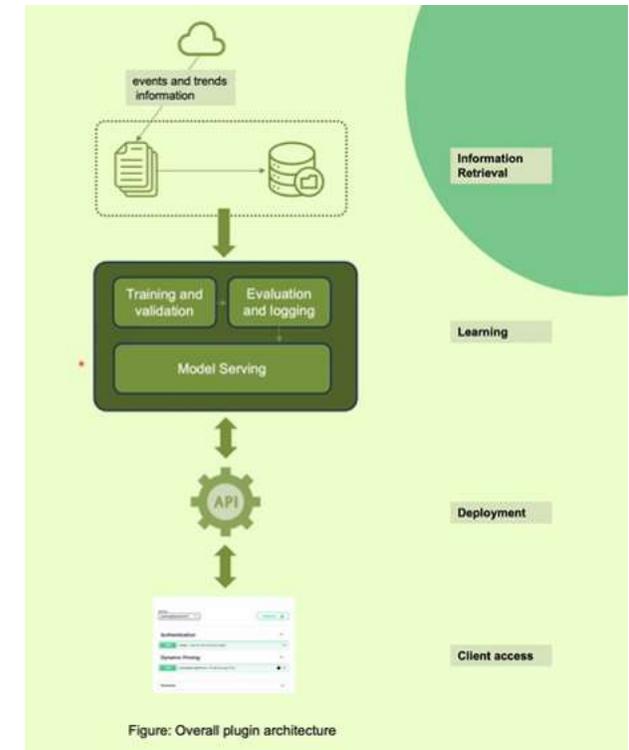
# UPCAST Environmental Plug-in

- **The Principle:** Societal and Environmental Well-being. Transparency.
- **Function:** estimates energy costs for data storage and processing, identifies and explains the factors affecting energy consumption using explainable AI techniques (XAI). It also provides real-time monitoring of energy usage during data processing.
- **Impact:** Enables "Green Data Spaces" where users choose low-impact reasoning paths.

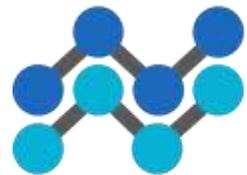Enables estimation of key environmental metrics such as energy consumption, carbon emissions, and costs

Offers two primary endpoints to estimate energy profiles:
- ➢ One for dataset storage energy profile at idle
- ➢ One for data processing workflow (DPW) energy profile

Julia Palma
julia.palma@ucd.ie

**CeADAR**
Ireland's Centre for AI

# CyclOps

## We are the partner providing….

- AI-based orchestration (AIMP)
- AI-Trustworthiness tools (XAI)
- Requirement collection coordination (T2.3 lead)

## …in WPs

- WP3: Runtime Layer
- WP5: IKB Exploitation

## We are in the project since…

CyclOps facilitates **collaboration** with diverse partners in **data space ecosystems**. Providing technological innovation with industrial relevance through the development of **automated AI-based trustworthy** components.

## As an outcome we expect…

- Facilitating the connection to various data sources and repositories through the **CyclOps AI tools & models marketplace**.

- Early access to **knowledge and CyclOps tools** for the management of the end-to-end data lifecycle for data space interoperability.

- Building **ecosystem connections** for potential collaborations, knowledge exchange, and opportunities in future projects.

## Our key people are…

*Ricardo Simón-Carbajo*
*Director of Innovation and Development*

*Julia Palma*
*Innovation Programmes Manager*

*Hanene Jemoui*
*Senior Software Engineer*

*Vicky Kumar*
*Research Software Engineer*

*Saul Burgess*
*Research Software Engineer*

Julia Palma
julia.palma@ucd.ie

# AI Marketplace (AIMP)

**CyclOps AI Marketplace Module**

•**The Principle:** Robustness. Transparency.

• **Function:** AI-based dashboard supporting automated algorithm matchmaking for users to efficiently select algorithms from a large catalogue of algorithms in use case scenarios such as Tourism, Green Deal, Public Procurement, and Manufacturing

•**Impact:** Efficiency for data processing and AI-based solutions development

Julia Palma
julia.palma@ucd.ie

CeADAR
Ireland's Centre for AI

CyclOps
**AI Marketplace (AIMP)**

- Algorithm Metadata Schema {
  "id": "Identifier for the algorithm",
  "repo_id": "Identifier for the repository hosting the algorithm",
  "dataset_id": "Identifier for the selected dataset",
  "name": "Name of the algorithm",
  "description": "Textual description of the algorithm",
  "category": ["Machine Learning", "Deep Learning", "Other"],
  "type": ["Classification", "Regression", "Clustering", "Anomaly Detection", "Forecasting", "Recommendation"],
  "applications": ["Remote sensing", "Land cover classification", "Change detection", "Tourism data segmentation"],
  "data_format": ["Numerical", "Categorical", "Textual", "Imagery", "Time series", "Sequential", "Acoustic", "Geospatial"],
  "optimizer": ["sgd", "mini-batch gradient descent", "adam", "adamw"],
  "url": "Web address of the algorithm's source code or library"
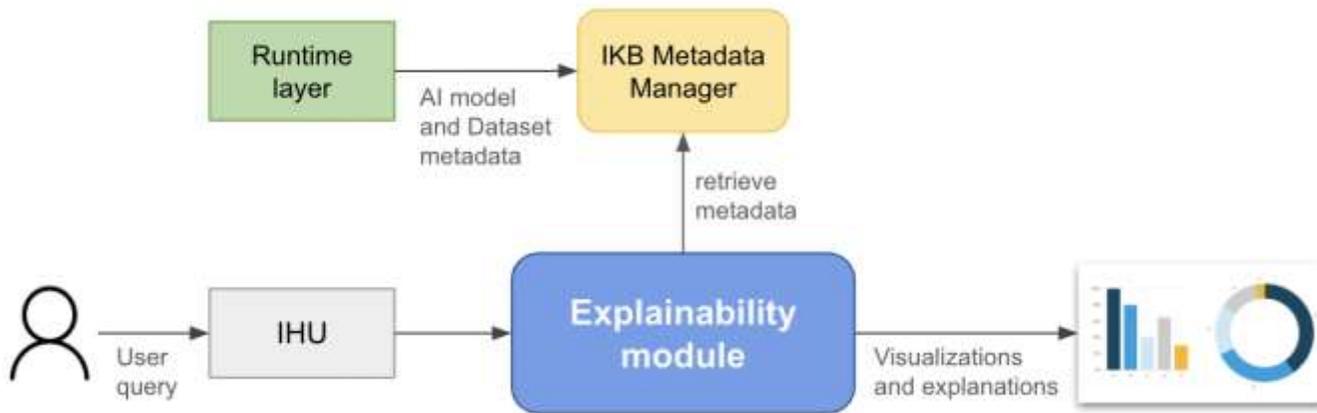}

- User Intent Metadata Schema

{
  "user_intent": {
    "id": "Identifier for the user intent",
    "session_id": "Identifier for the user session",
    "dataset_id": "Identifier for the selected dataset",
    "application_requirements": {
      "application_type": "Selected application type such as Classification, Regression, Clustering, etc",
      "application_domain": "Selected application domain such as Tourism, Green Deal, Public Procurement, etc",
      "description": "Textual description of the user query"
    },
    "functional_requirements": {
      "accuracy": ["High", "Medium", "Low"],
      "training_time": ["High", "Medium", "Low"],
      "inference_time": ["High", "Medium", "Low"],
      "scalability": ["High", "Medium", "Low"]
    },
    "infrastructure_constraints": {
      "cpu_cores": "Number of CPU cores available",
      "ram_gb": "Amount of RAM available in GB",
      "gpu_gb": "Amount of GPU memory available in GB",
      "platform": ["on-premise", "cloud", "edge"]
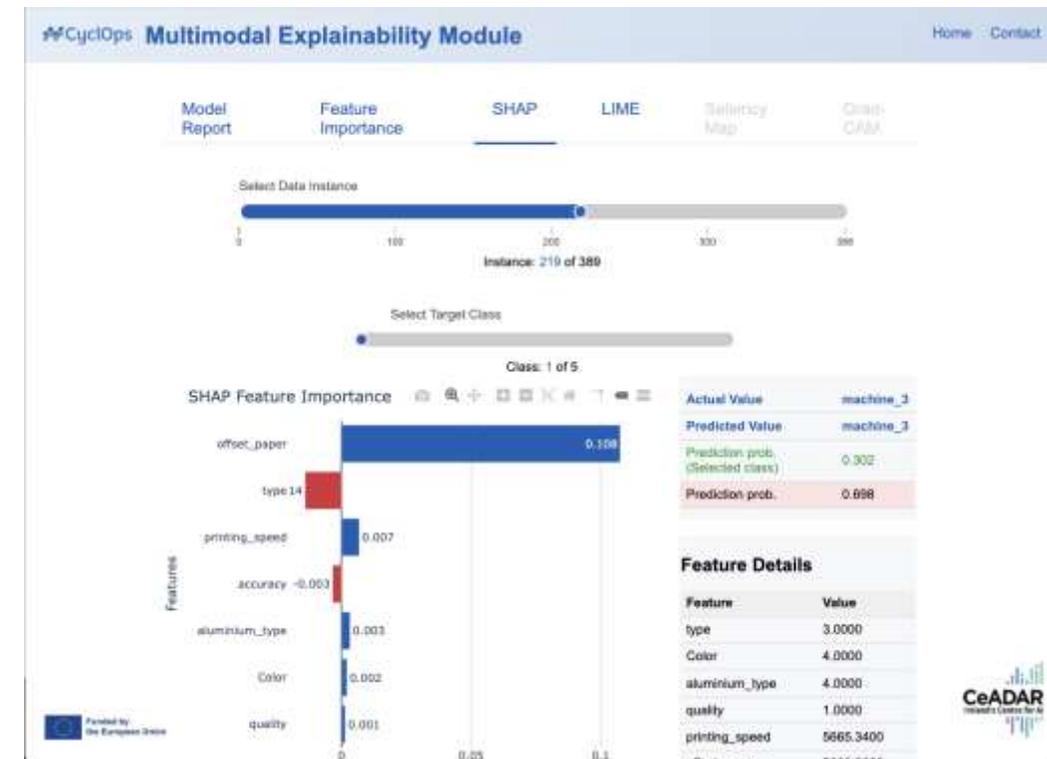    }
  }
}

Julia Palma
julia.palma@ucd.ie

# Multimodal Explainability Module (MXAI)



## CyclOps Multimodal Explainability Module
•**The Principle:** Transparency. Human oversight.
• **Function:** AI-based dashboard supporting automated algorithm matchmaking for users to efficiently select algorithms from a large catalogue of algorithms in use case scenarios such as Tourism, Green Deal, Public Procurement, and Manufacturing
•**Impact:** Explainability, enhancing trust and interpretability.

Julia Palma
julia.palma@ucd.ie

CeADAR
Ireland's Centre for AI

Takeaways

Shift focus from optimizing performance to **optimizing for accountability**

Integrate **sustainability** and **fairness** into the very first lines of your architecture

AI and Knowledge-Based Systems are the **means**; human-centric wisdom and societal well-being are the **ends**