



**chist-era**



# CHIST-ERA Projects Seminar 2022

## *Explainable Machine Learning-based Artificial Intelligence (XAI)*

***Speaker: Slawomir Nowaczyk***  
**(project XPM)**

**March 30th, 2022**



Programme co-funded by the  
EUROPEAN UNION

# Introduction: Definition of the Topic

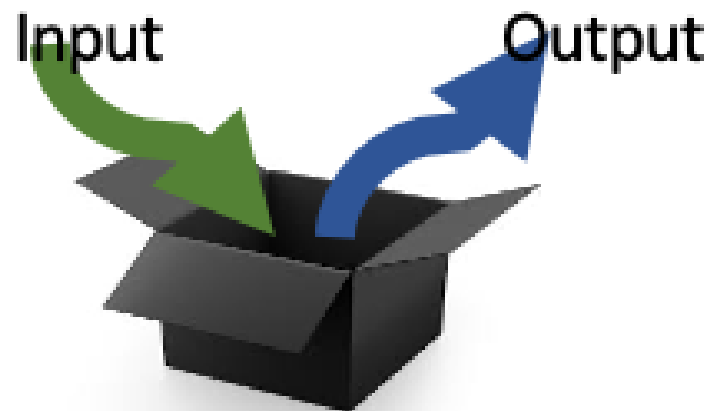
Most AI methods are “black boxes”, i.e., it’s difficult to know how and why they work.

XAI are tools to explain the models and their decision, so that humans can trust them.

Modern AI models, such as deep learning, are trained in an end-to-end manner from “big data” → it’s unclear what exactly the model is learning.

Important keywords:

- ❖ Trust
- ❖ Human-centric
- ❖ Transparency
- ❖ Personalization
- ❖ Interpretability of the models
- ❖ Explanation negotiation/argumentation
- ❖ Prediction / Classification
- ❖ Data-driven
- ❖ Neuro-Symbolic integration





**chist-era**

# Introduction: Projects of the Topic

**ANTIDOTE:** Provide a unified computational framework for jointly learning clinical predictions and the associated argumentative justifications.

**CausalXRL:** Learn two-level Causal model of environment to enable Reinforcement Learning agent to suggest eXplainable actions in critical environments.

**CIMPLE:** Automatically detect misinformation & generate creative explanation response

**COHERENT:** Generate and combine explanations and decide when to deliver them during robotic collaboration tasks.

**EXPECTATION:** Generate personalized explanations based on heterogeneous knowledge during iterative interactions with users and agents.

**GraphNEx:** Graph neural networks for XAI & user-guided explainability

**INFORM:** Interpretability of DNNs in oncology applications on data extracted from medical images (radiomics), and evaluation of interpretability methods

**iSEE:** Retrieve, personalise and reuse explanation strategies on different use case based on past explanation experiences evaluation against persona intents

**MUCCA:** Quantifying strengths and solving weaknesses of new and state of the art XAI methods using heterogeneous use cases

**SAI:** Decentralised ML-based AI in human social complex systems

**XAIface:** Explain black box facial recognition methods and fight against biased results.

**XPM:** How explanations support decision making in predictive maintenance systems



# Major Achievements and Outputs

## Evaluating XAI methods:

- User-guided explainability & co-creation of explanations
- Evaluating the reliability of explanations
- Explanation experience ontology

## XAI algorithms & methods:

- Resource-efficient training
- Incorporate ontologies to standardize and shape explanations
- Approaches for Generating Personalized Explanations

## Platforms & Frameworks:

- EREBOTS - Chabot personalizing explanations via argumentation
- Symbolic knowledge extraction & injection (PSyKE & PSyKI)
- Simulators for large-scale decentralised AI evaluation
- Explanation strategies recommendation and evaluation platform



# Upcoming Challenges and Needs

Generalization of XAI methods to different domains

Acceptance and implementation of AI systems

Challenging for end-users to trust developed AI models

XAI certification guiding organisations using AI systems

Explanation(s) alignment (in case of disagreement) in distributed envs.

Ontologies alignment across different parties/actors

Integration of different nature/data type explanations

Explanation personalization according to users' profiles & needs

Subjectivity in the explanation evaluation

Involve user in the loop



**chist-era**

# Possible Roadmap

Multidisciplinary approach, incorporating different perspectives and techniques inspired by existing state-of-the-art work, and sharing the knowledge and experience gained through the project with other researchers

Experiments with users to iteratively get their feedback, for example, design an interactive decision support system for Predictive Maintenance that can be used to evaluate different types of explanations for different types of users

Exploiting connection between the predictions and the underlying physics governing the system, which is generally assumed to be well-understood by the expert and can therefore provide common ground between the AI and the human.

Investigating different levels of explanations based on user expertise

Understand explainability at the collective level in decentralised, collaborative AI

## Expected impact:

overall higher acceptance of AI systems in all areas of human life,  
and more knowledge created through use of AI



Very interesting and relevant topic

(plus community-driven decision process for defining them)

Bringing the researchers together → let's meet for real next year!

Facilitate the contact between researchers of different projects to identified complementariness (ex: new use case for XAI algorithm)

One of the few programmes looking for long-term research ideas

Idea: take statistics about average timings (and delays) for fundings on a per-funding-agency basis, and make that public

Idea: propose rules concerning different funding categories (e.g., min travel budget), and properly communicate them to the funding agencies



## Integrated good practices for RRI

- ❖ Publish all papers in institutional repositories (open access)
- ❖ Share the data and software in zenodo/other repositories
- ❖ Conducting human experiments under the supervision of ethical committees

## Major hurdles to implement RRI

- ❖ Gender balance in technical domains is difficult due to the low percentage of women in computer science, robotics and AI
- ❖ Predatory journals are not well identified and create an unfair environment to develop science
- ❖ Cost of open-access journals → we can publish the author version in your institutional repository and publish without paying





Public dissemination through scientific publications and conferences

Free online access or open gold model

New data collected is released as public accessible databases

Whenever possible, at least

Open source and reproducible research platforms are used

Data: FAIR principles are followed

(Findable, Accessible, Interoperable, Reusable)

Explore EU platforms: OpenAIRE, ..., AI4EU, and more



Establishing **community** and social network of XAI researchers

Open Source **platform** & European XAI **compliance certification** framework

**Cross-application** assessment of available approaches to XAI

Find a **common representation** for heterogeneous explanation methods

Create an **European XAI compliance certification** for AI software

**Integration** of novel solutions into existing commercial solutions

Actions should involve all CHIST-ERA member countries:

- Identify transferable technologies
- Survey of potentially interested stakeholders
- Participate/organize events (workshops / training schools)



# Questions ?