

Introduction to Explainable AI

Helmholtz AI @ Jülich
14.05.2024

Who are we?

Helmholtz AI



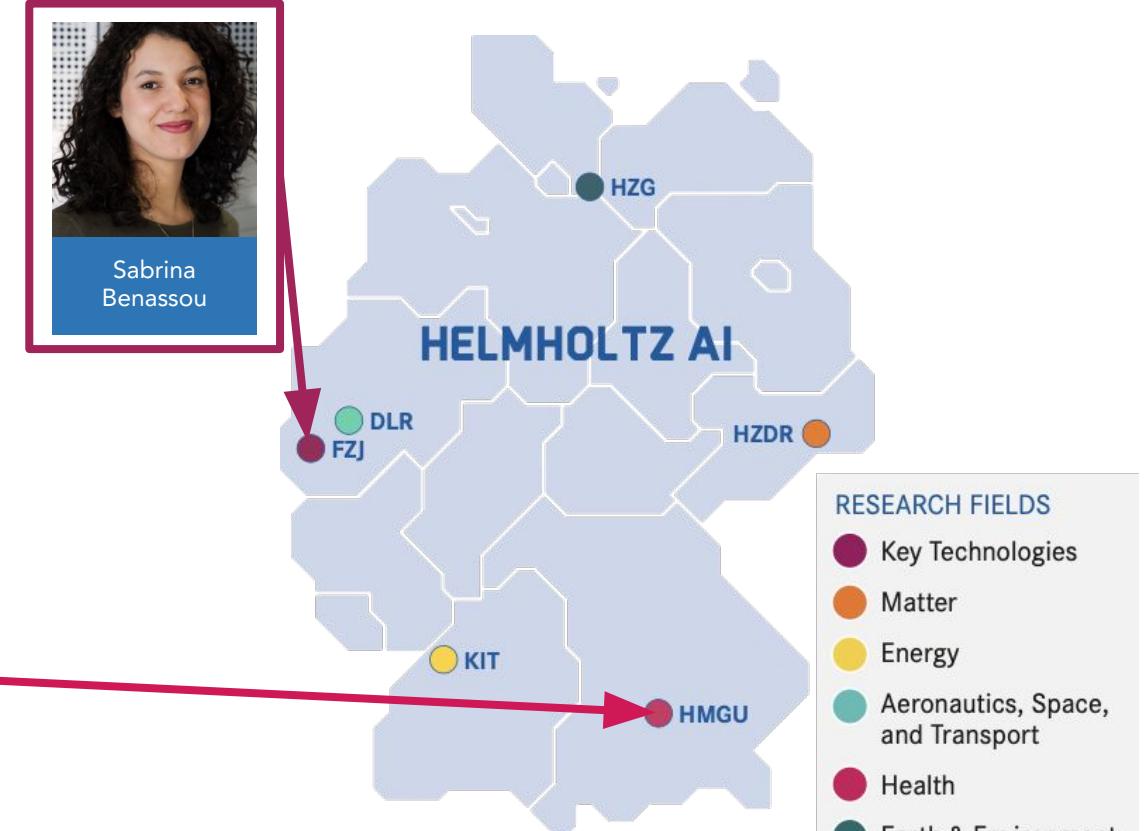
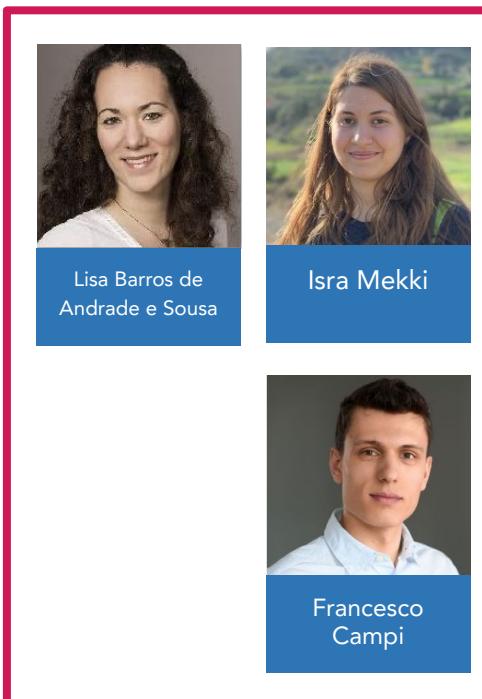
WHAT IS OUR MISSION?



Maximise research impact by
democratising access to AI

HOW DO WE DO THAT?

- Short- to mid-term scientific collaboration (2 weeks - 6 months)
- Free of charge
- Easy application



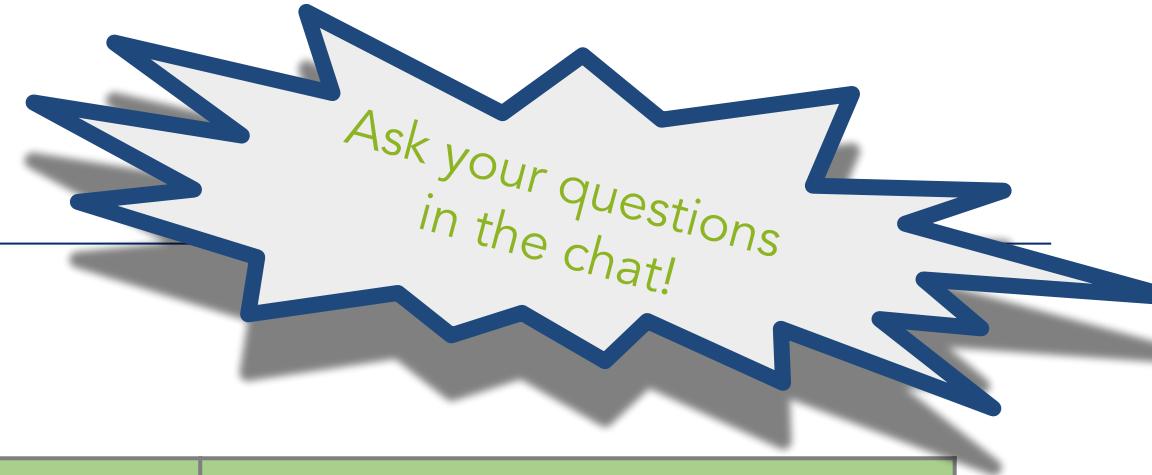


What is your field of study?

- ① Start presenting to display the poll results on this slide.

Outline

Schedule and Tools



Day 1 9:00 - 13:00	Day 2 9:00 - 13:00	Day 3 9:00 - 13:00
XAI for Random Forests	XAI for CNNs	XAI for Transformers



Introduction

Terminology

Explainability or Interpretability?

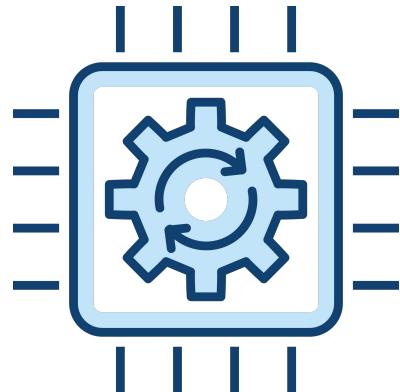


Introduction

Terminology

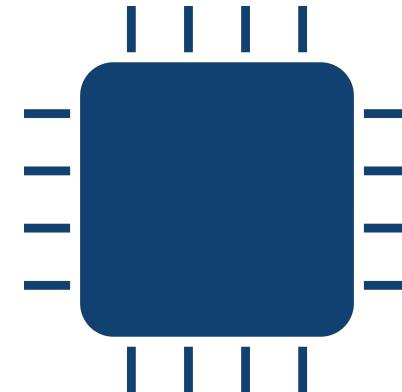
Interpretability

Understand exactly why and how the model is generating predictions by observing the inner mechanics of the AI/ML method.



Explainability

Focus on the decision-making process and try to explain the behaviour in human understandable terms.

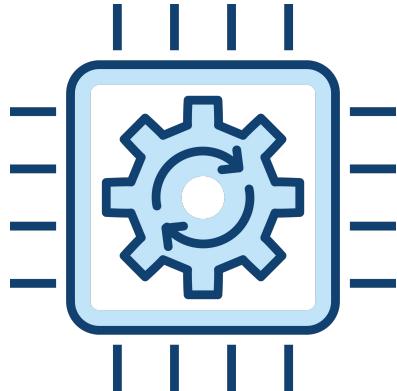


Introduction

Terminology

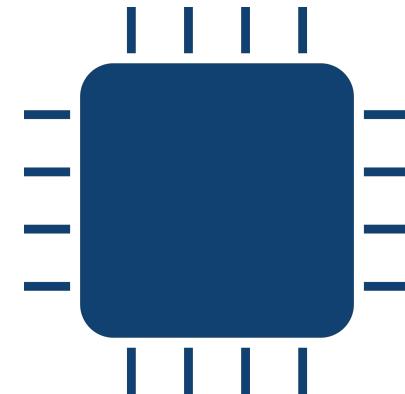
Interpretability

Understand exactly why and how the model is generating predictions by observing the inner mechanics of the AI/ML method.



Explainability

Focus on the decision-making process and try to explain the behaviour in human understandable terms.

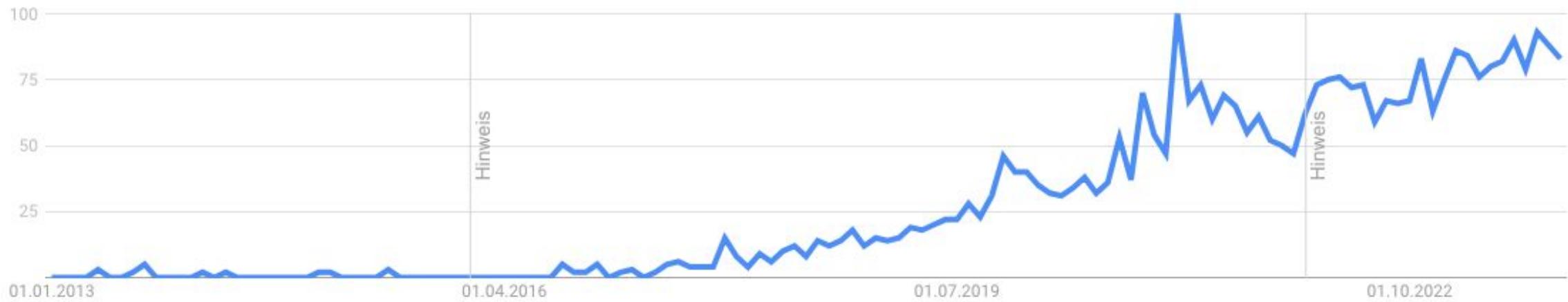


In this course, we will focus only on **eXplainable Artificial Intelligence (XAI)**.

Introduction

Why is explainability important?

Google Trends Popularity Index of the term *Explainable AI* over the last ten years (2013–2023)





Why is explainability important?

- ① Start presenting to display the poll results on this slide.

Introduction

Why is explainability important?

„The problem is that a single metric, such as classification accuracy, is an incomplete description of most real-world tasks.“ — (Doshi-Velez et al., 2017)

Introduction

Why is explainability important?

„The problem is that a single metric, such as classification accuracy, is an incomplete description of most real-world tasks.” — (Doshi-Velez et al., 2017)



Introduction

XAI is important for technology acceptance



Introduction

XAI is important to avoid ethical issues

NEWS | 24 October 2019 | Update [26 October 2019](#)

Millions of black people affected by racial bias in health-care algorithms

Study reveals rampant racism in decision-making software used by US hospitals – and highlights ways to correct it.

[Heidi Ledford](#)



Introduction

XAI is important for knowledge creation

What Does Deep Learning See? Insights From a Classifier Trained to Predict Contrast Enhancement Phase From CT Images

Kenneth A. Philbrick¹
Kotaro Yoshida
Dai Inoue
Zeynettin Akkus
Timothy L. Kline
Alexander D. Weston
Panagiotis Korfiatis
Naoki Takahashi
Bradley J. Erickson

OBJECTIVE. Deep learning has shown great promise for improving medical image classification tasks. However, knowing what aspects of an image the deep learning system uses or, in a manner of speaking, sees to make its prediction is difficult.

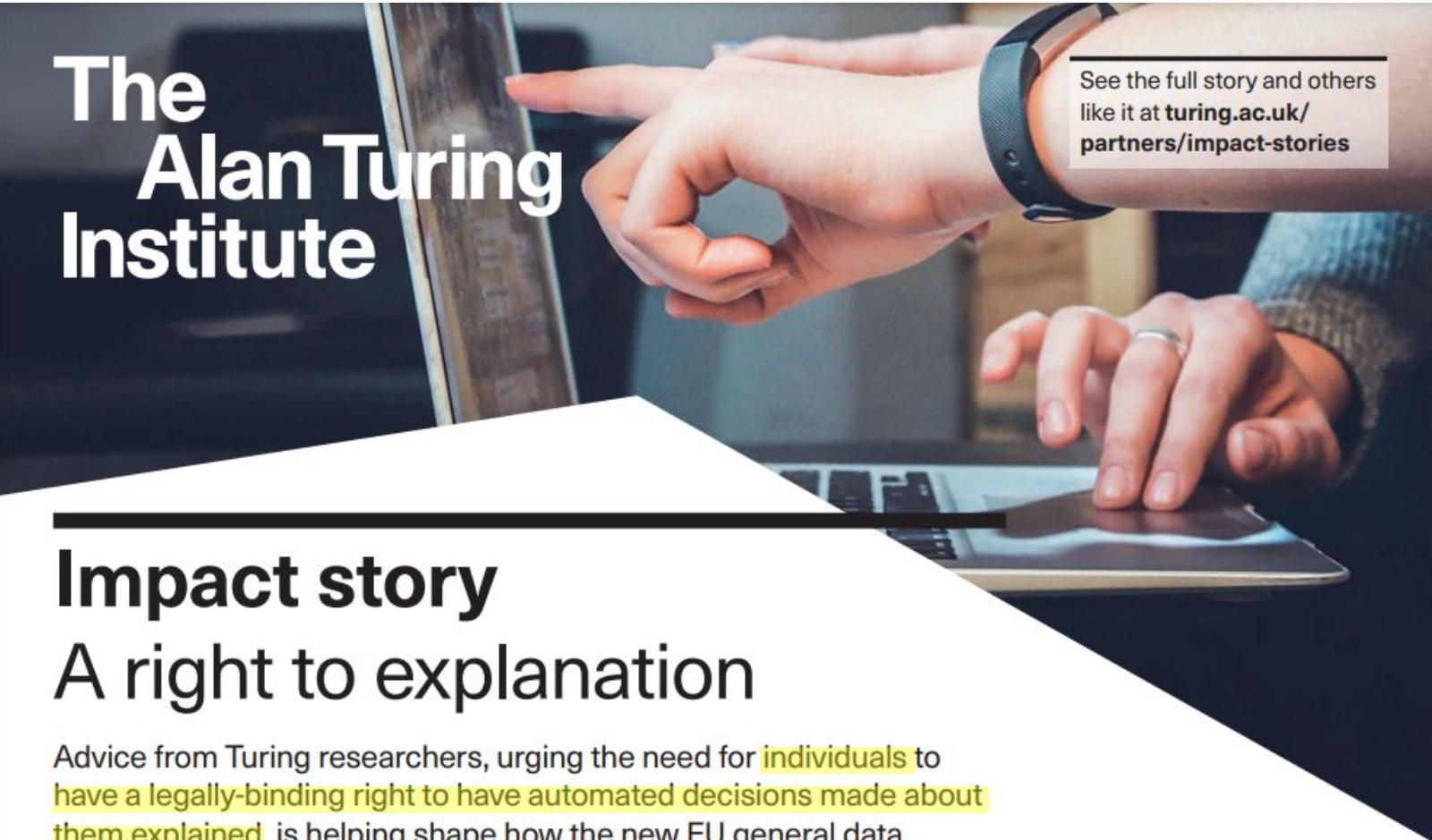
MATERIALS AND METHODS. Within a radiologic imaging context, we investigated the utility of methods designed to identify features within images on which deep learning activates. In this study, we developed a classifier to identify contrast enhancement phase from whole-slice CT data. We then used this classifier as an easily interpretable system to explore the utility of class activation map (CAMs), gradient-weighted class activation maps (Grad-CAMs), saliency maps, guided backpropagation maps, and the saliency activation map, a novel map reported here, to identify image features the model used when performing prediction.

RESULTS. All techniques identified voxels within imaging that the classifier used. SAMs had greater specificity than did guided backpropagation maps, CAMs, and Grad-CAMs at identifying voxels within imaging that the model used to perform prediction. At shallow network layers, SAMs had greater specificity than Grad-CAMs at identifying input voxels that the layers within the model used to perform prediction.

CONCLUSION. As a whole, voxel-level visualizations and visualizations of the imaging features that activate shallow network layers are powerful techniques to identify features that deep learning models use when performing prediction.

Introduction

XAI is important to meet regulatory requirements



The Alan Turing Institute

Impact story

A right to explanation

Advice from Turing researchers, urging the need for individuals to have a legally-binding right to have automated decisions made about them explained, is helping shape how the new EU general data protection regulations (GDPR) will be implemented.

See the full story and others like it at turing.ac.uk/partners/impact-stories

Introduction

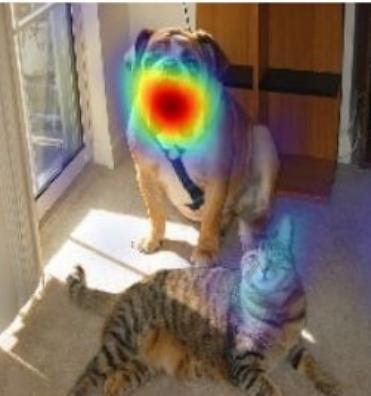
XAI is important as a defense strategy



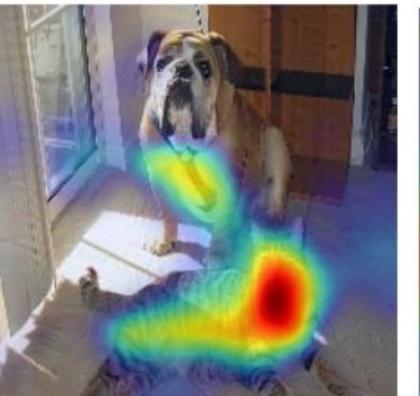
Boxer: 0.4 Cat: 0.2
(a) Original image



Airliner: 0.9999
(b) Adversarial image



Boxer: 1.1e-20
(c) Grad-CAM "Dog"



Tiger Cat: 6.5e-17
(d) Grad-CAM "Cat"



Airliner: 0.9999
(e) Grad-CAM "Airliner"



Space shuttle: 1e-5
(f) Grad-CAM "Space Shuttle"

[Home](#) > [Artificial Intelligence and Soft Computing](#) > Conference paper

Explainable AI for Inspecting Adversarial Attacks on Deep Neural Networks

Zuzanna Klawikowska, Agnieszka Mikołajczyk & Michał Grochowski

Conference paper | [First Online: 07 October 2020](#)

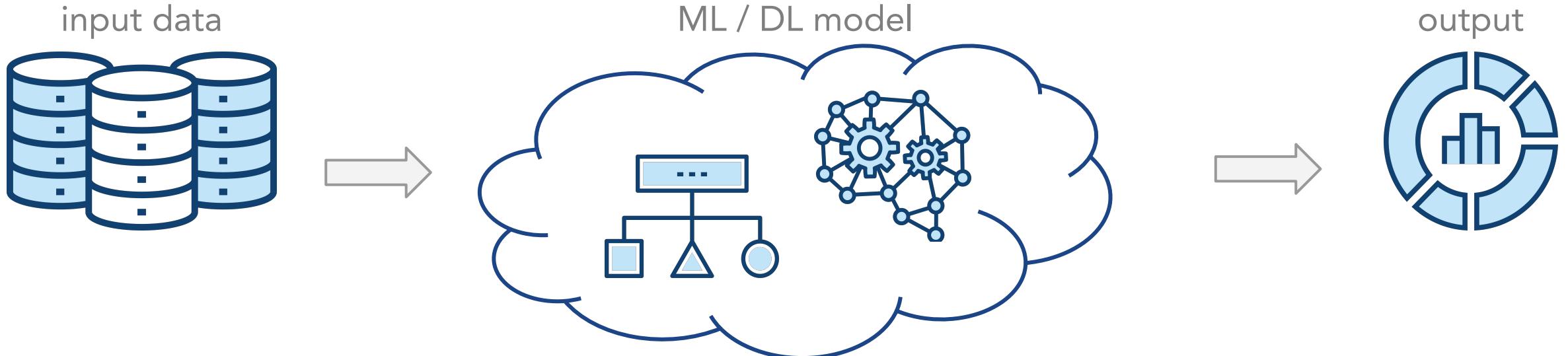
2252 Accesses | 1 Citations

Part of the [Lecture Notes in Computer Science](#) book series (LNCS, volume 12415)

https://campusai.github.io/_papers/Grad-CAM/adversarial.png

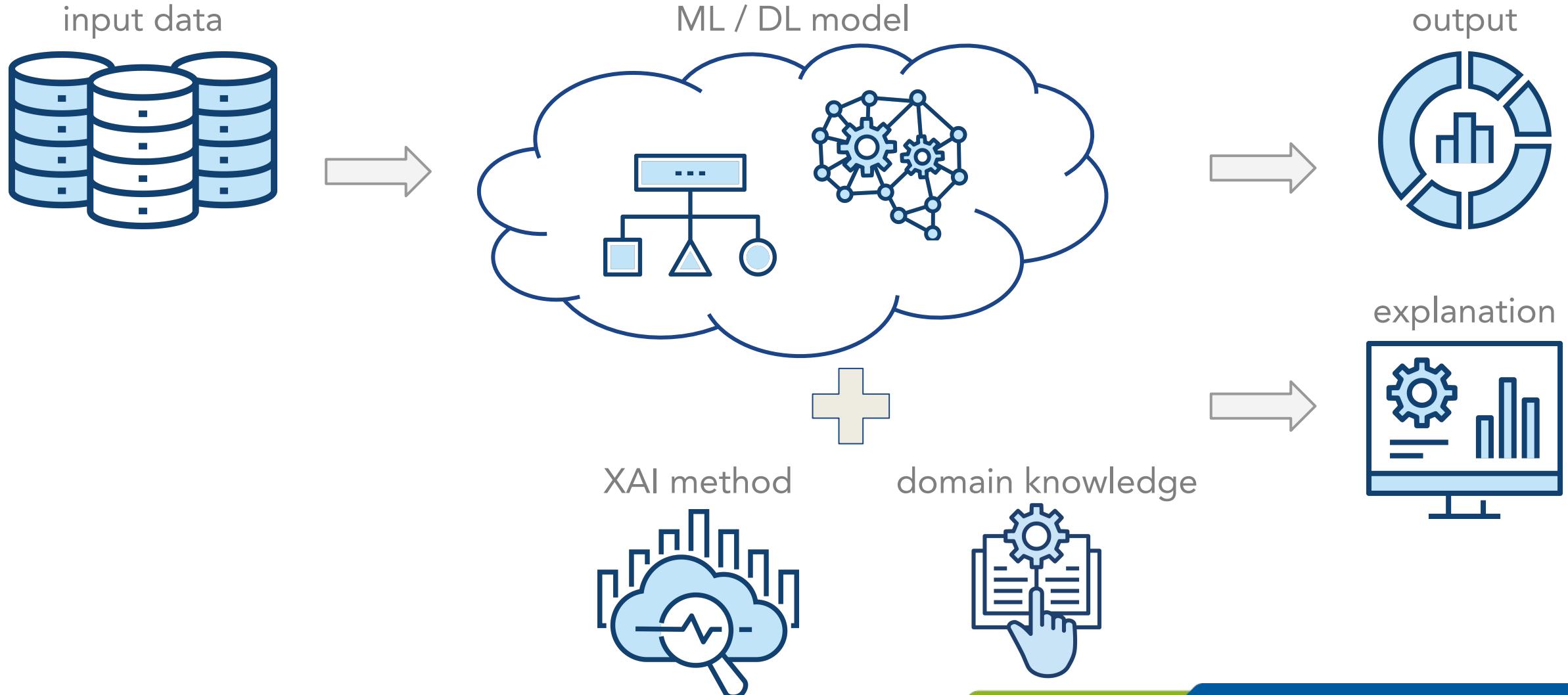
Introduction

XAI in your ML workflow



Introduction

XAI in your ML workflow



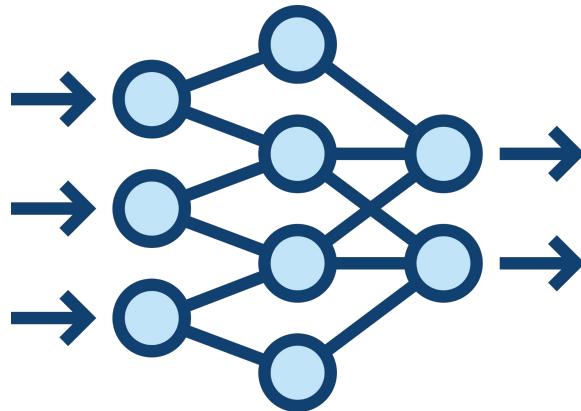
Introduction

XAI in your ML workflow

input data



ML / DL model



output



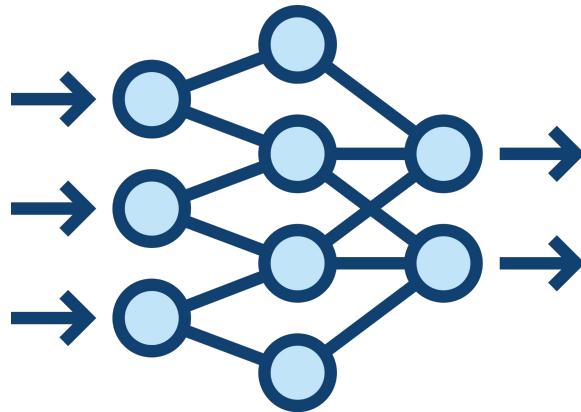
Introduction

XAI in your ML workflow

input data



ML / DL model



Current explanation:
This is a cat!

output

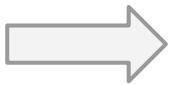


Cat

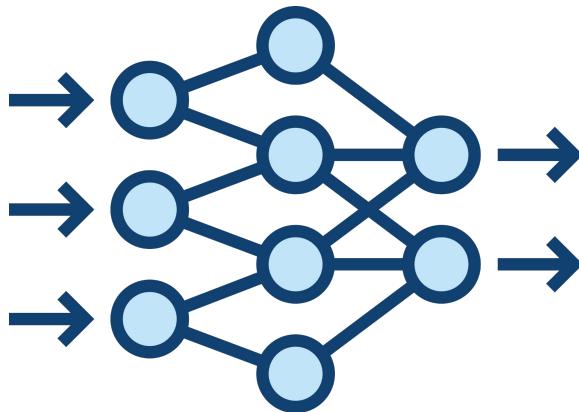
Introduction

XAI in your ML workflow

input data



ML / DL model



Current explanation:
This is a cat!

output

Cat



XAI method +
domain knowledge

XAI explanation:

- it has fur, whiskers, and claws
- it has this feature



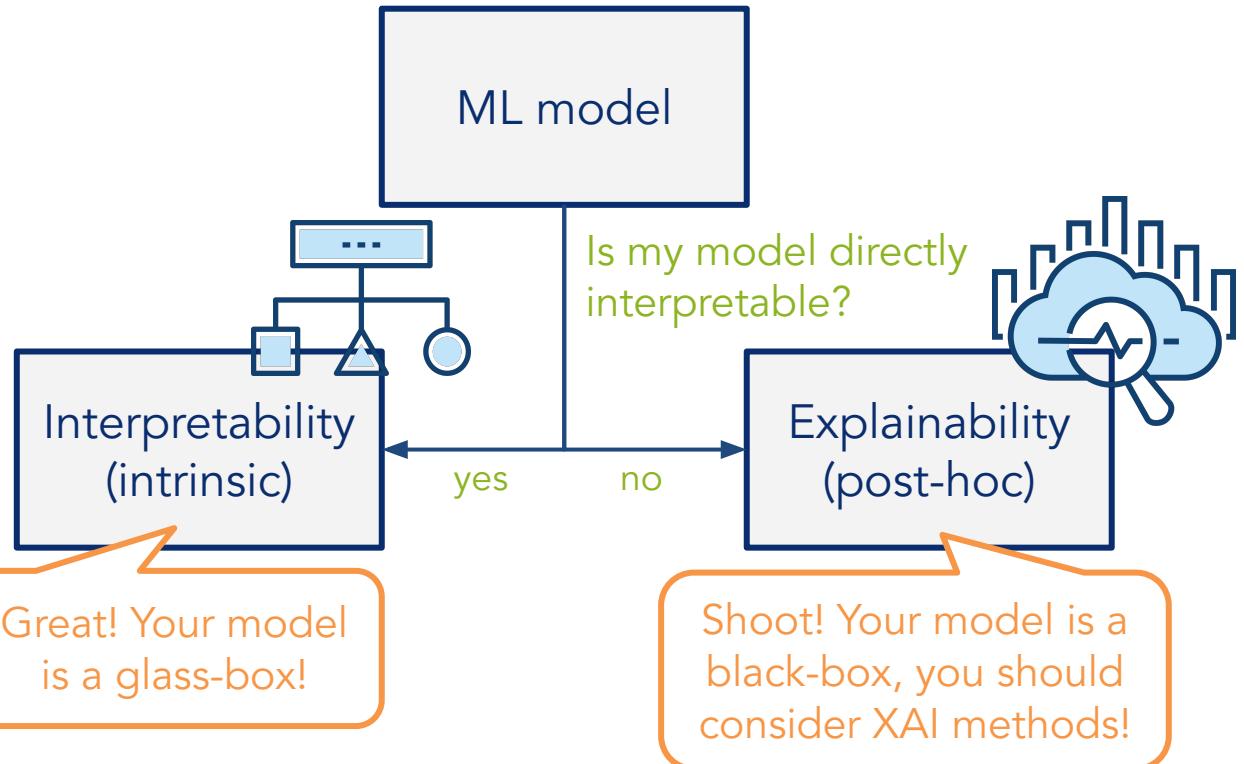
Introduction

Taxonomy of XAI methods



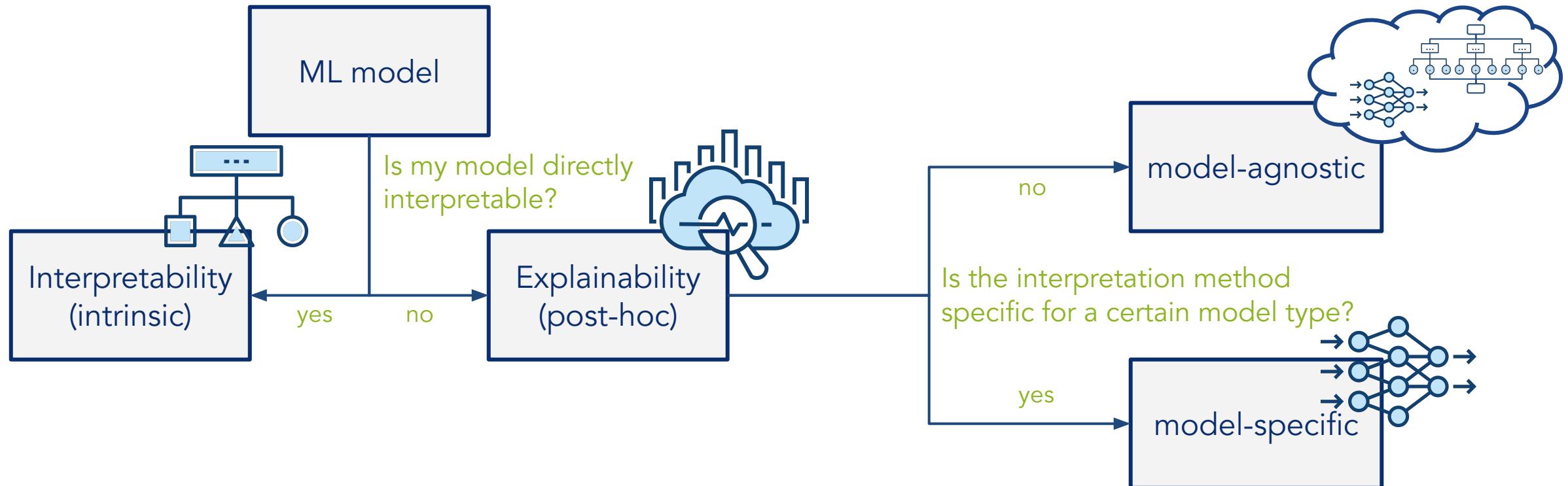
Introduction

Taxonomy of XAI methods



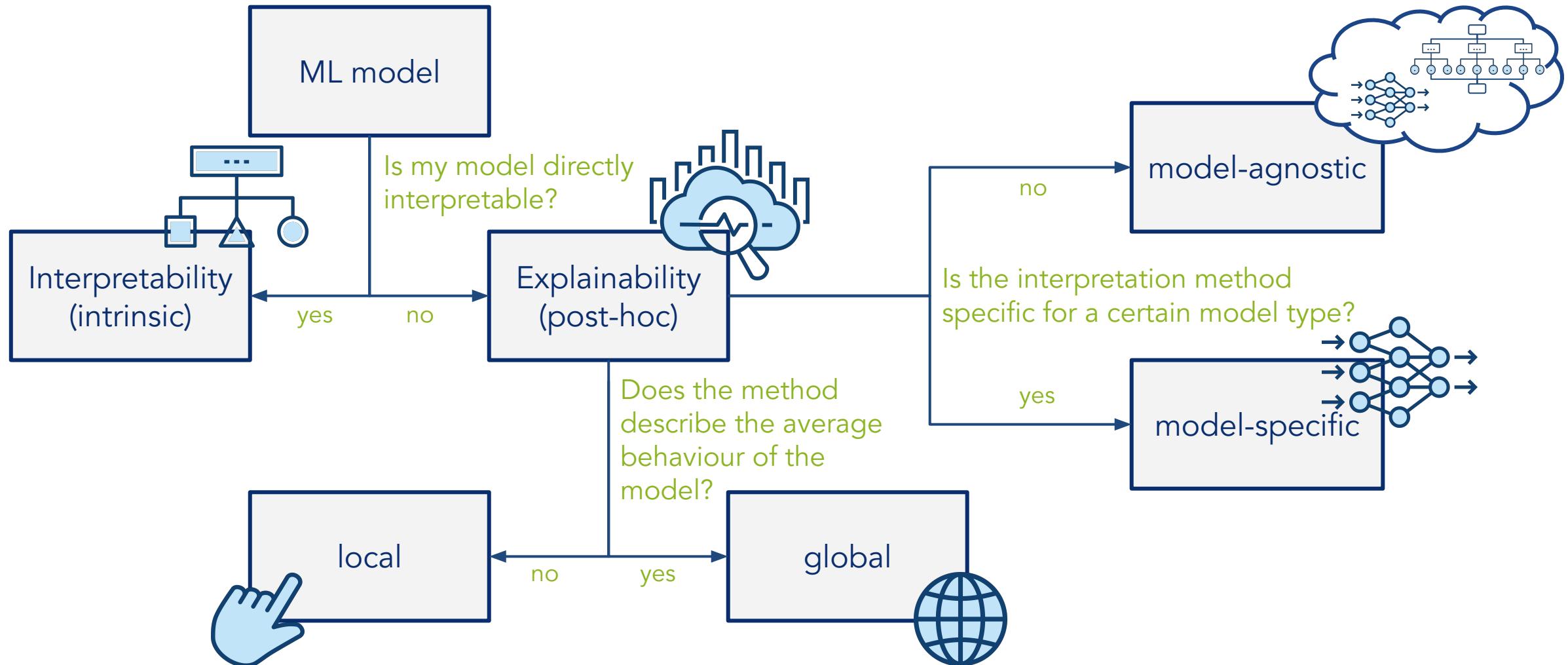
Introduction

Taxonomy of XAI methods



Introduction

Taxonomy of XAI methods



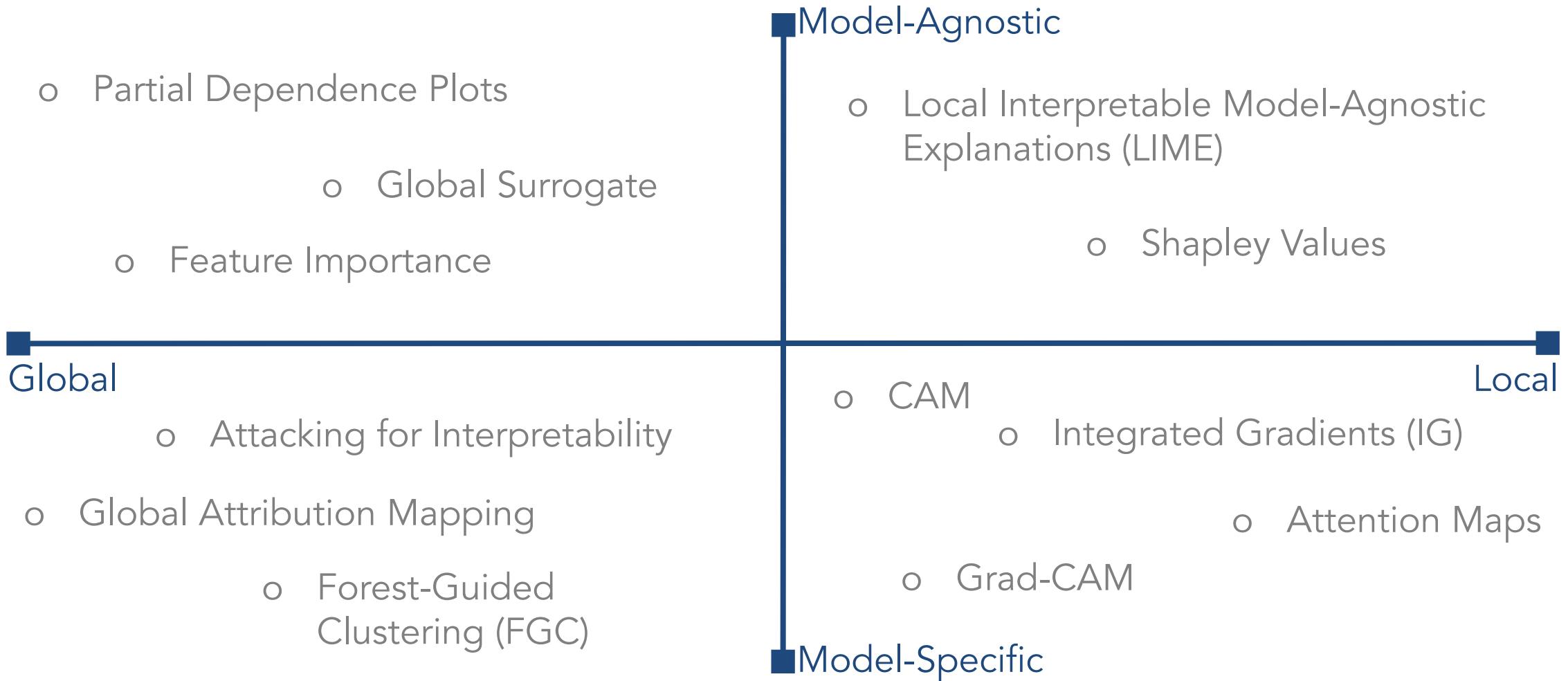
- —
- —
- —

To understand what impact blood pressure has on the survival rate of patient John Doe in a Random Forest model, we need:

- ① Start presenting to display the poll results on this slide.

Introduction

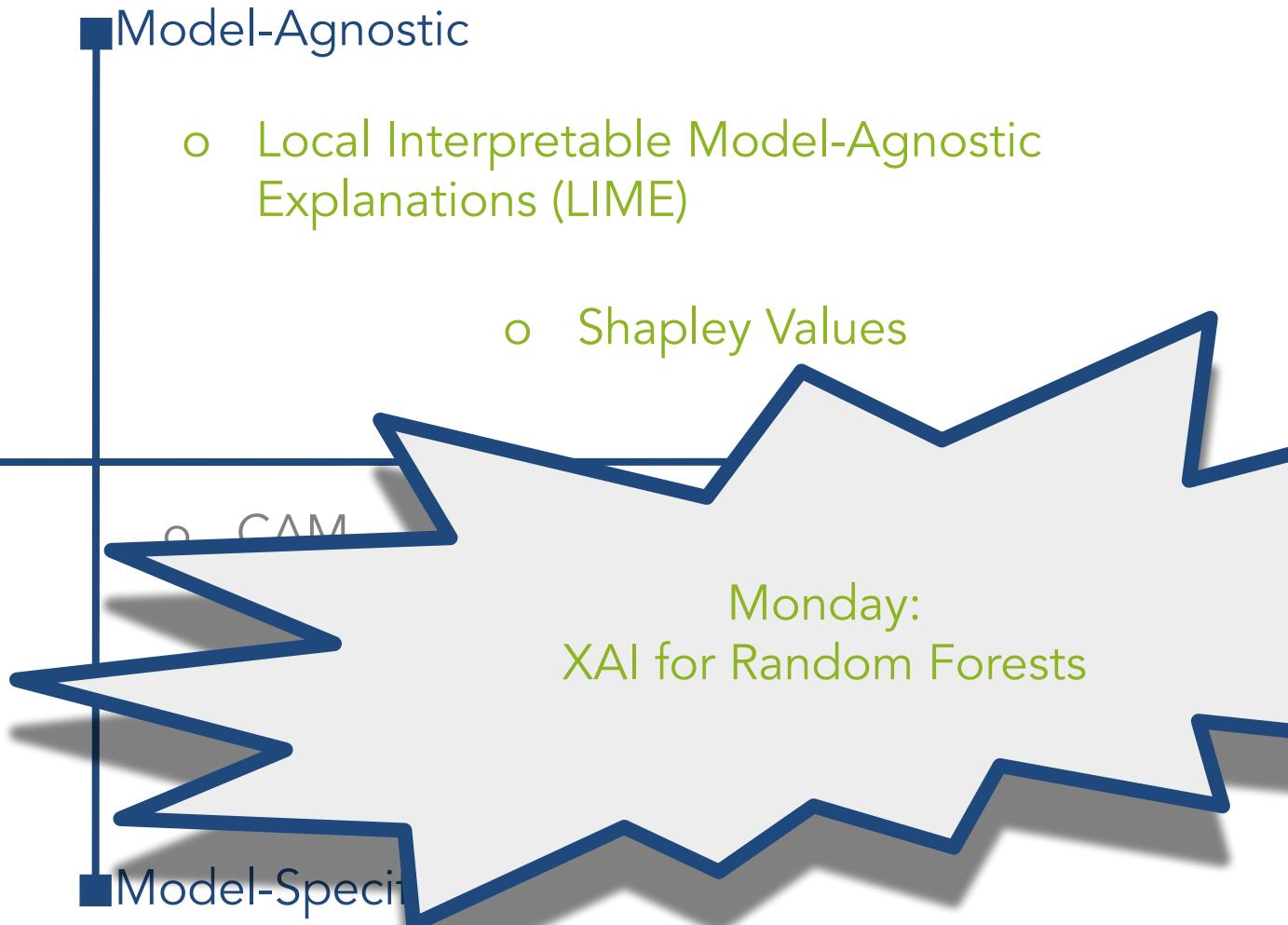
Overview on post-hoc methods



Introduction

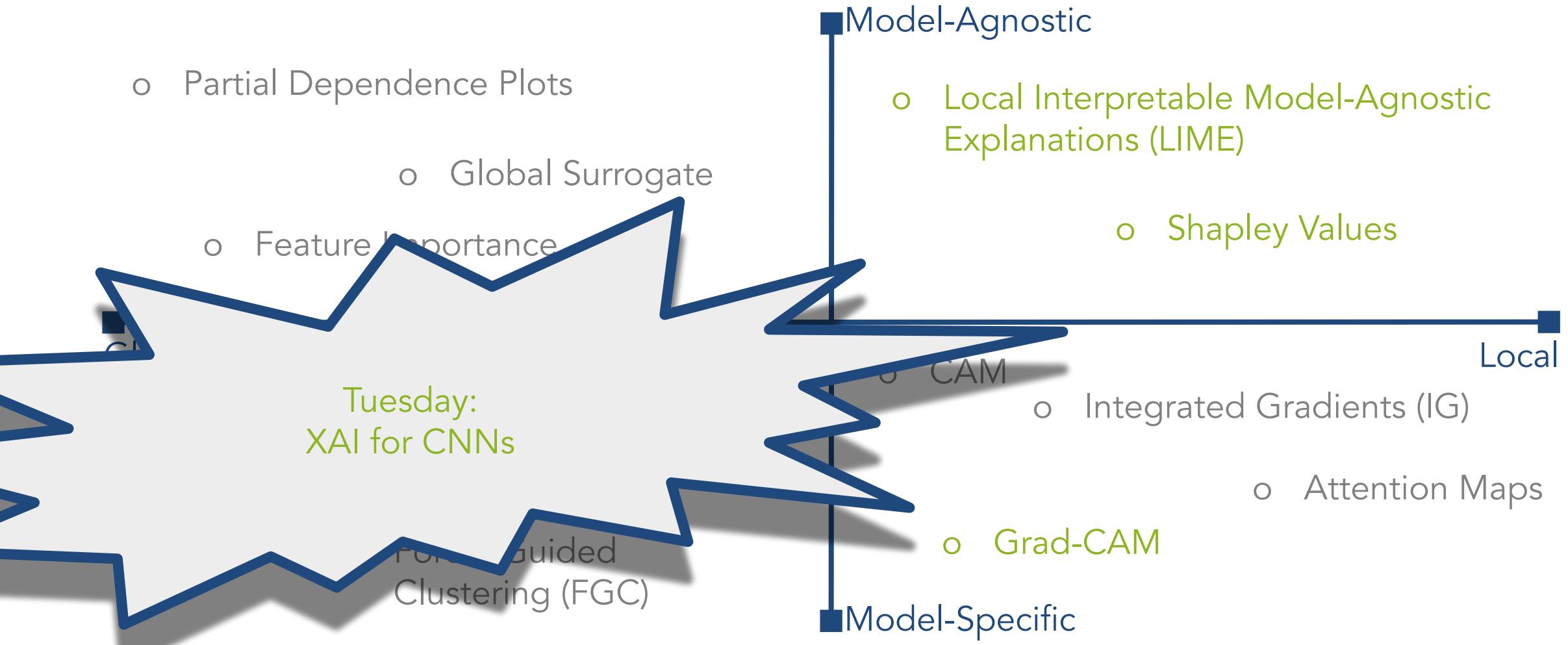
Overview on post-hoc methods

- Partial Dependence Plots
- Global Surrogate
- Feature Importance
- Attacking for Interpretability
- Global Attribution Mapping
- Forest-Guided Clustering (FGC)



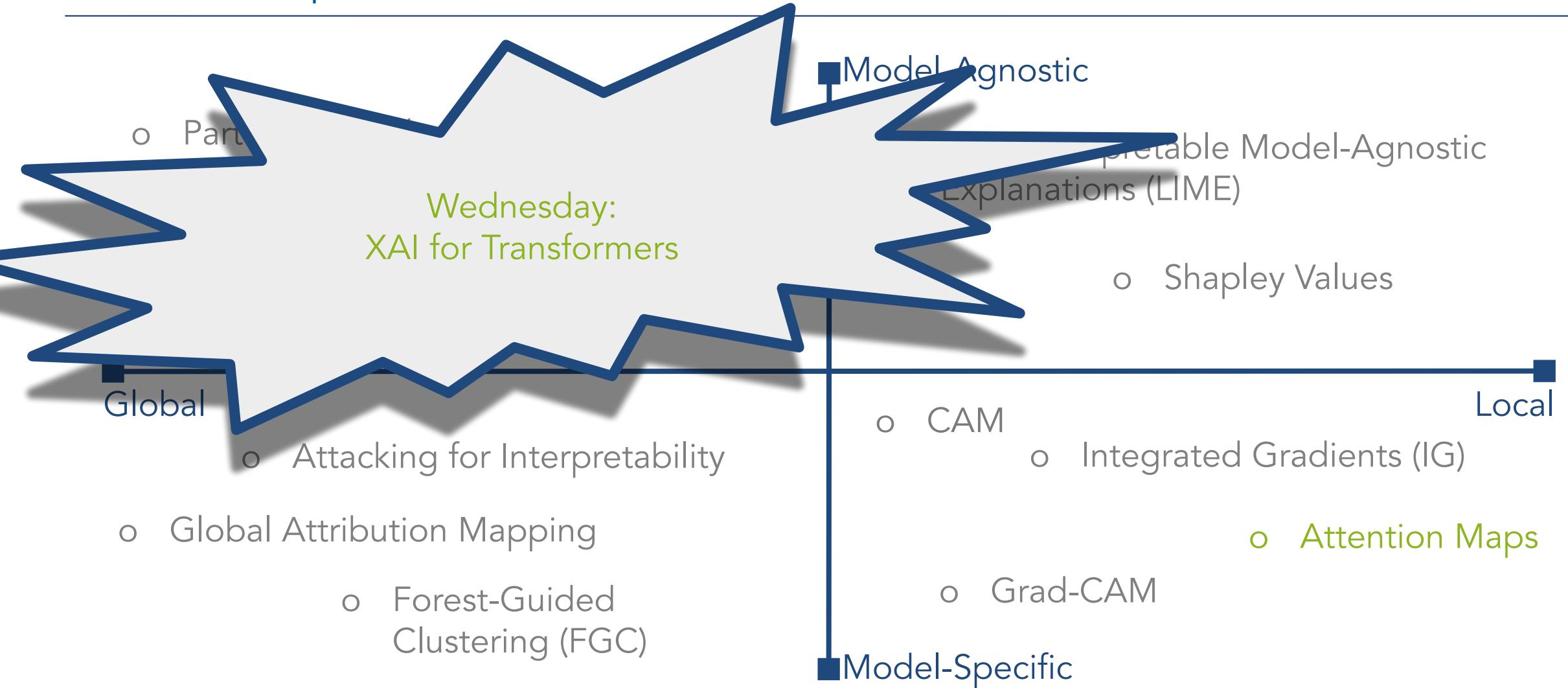
Introduction

Overview on post-hoc methods



Introduction

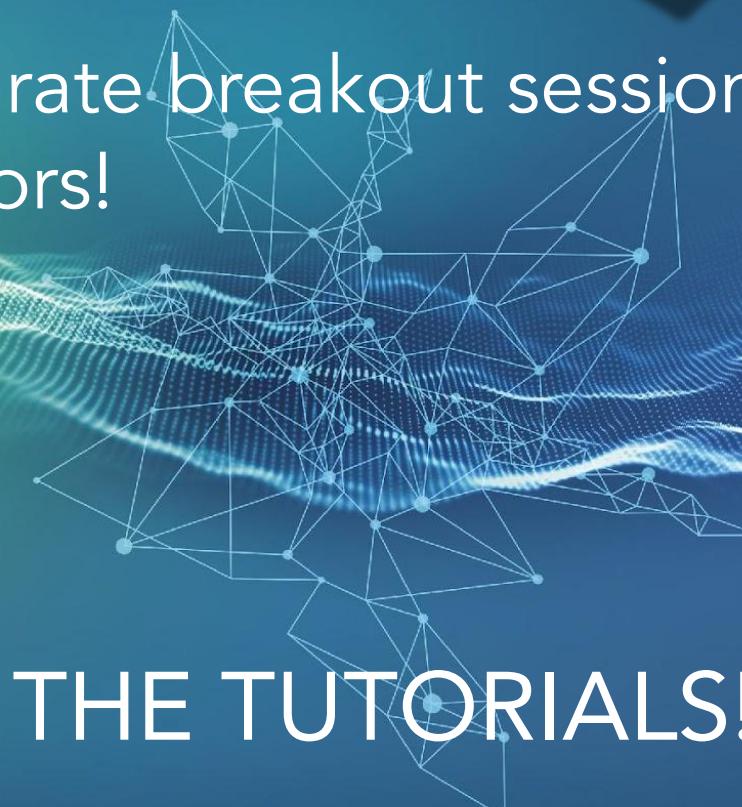
Overview on post-hoc methods



We will move you now into separate breakout sessions with your
tutors!



HAVE FUN WITH THE TUTORIALS!



Who are we?

Helmholtz AI

WHAT IS OUR MISSION?



Maximise research impact by
democratising access to AI

WHO ARE WE?



Interdisciplinary platform for
innovative research in AI



Compiles develops and fosters
applied AI methods nationwide
across all Helmholtz Centers



Aims to reach international
leadership in applied AI

If you have questions on Helmholtz AI, contact us at:
consultant-helmholtz.ai@helmholtz-muenchen.de

