

Introduction to Explainable AI



Helmholtz AI @ MALTAomics Summer School
13.09.2023

Who are we?

Helmholtz AI

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

WHAT IS OUR MISSION?



Maximise research impact by democratising access to AI

WHO ARE WE?



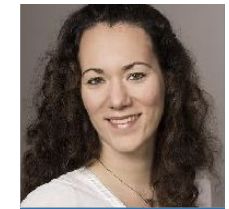
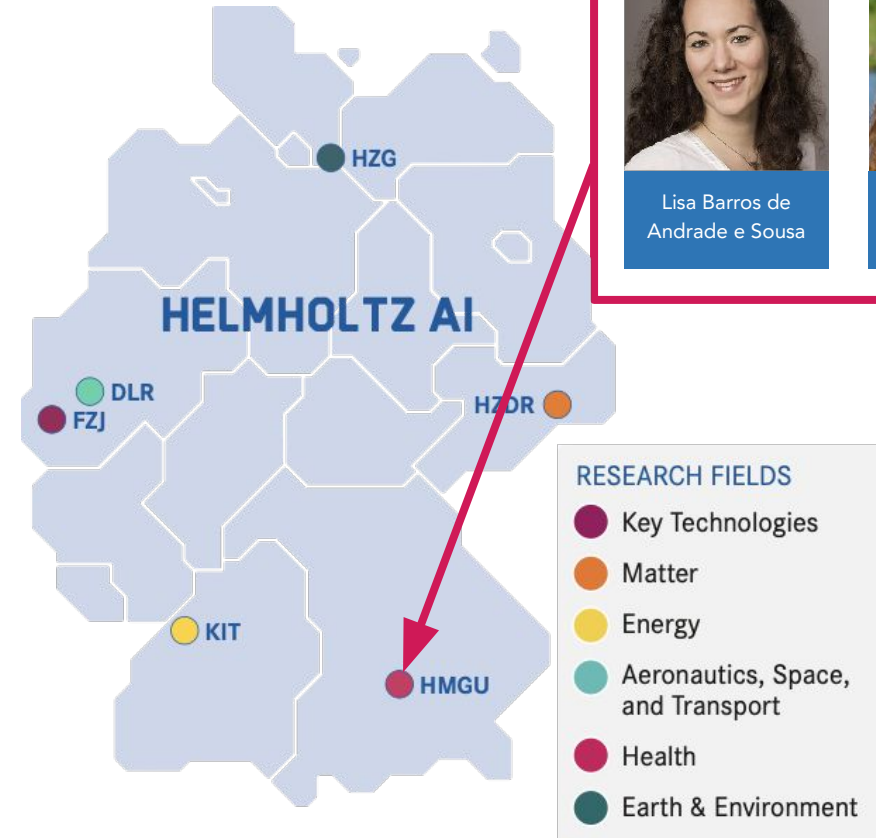
Interdisciplinary platform for innovative research in AI



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



Aims to reach international leadership in applied AI



Lisa Barros de
Andrade e Sousa



Donatella Cea



What is your field of study?

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

Outline

Schedule



You can ask
questions anytime!

10.00 - 10.20	Introduction on XAI
10.20 - 10.50	XAI Model-Agnostic Methods: "Permutation Feature Importance"
10.50 - 11.00	Break
11.00 - 11.40	XAI Model-Agnostic Methods: "SHAP"
11:40 - 11.50	Break
11.50 - 12.20	XAI Model-specific Methods: "Forest Guided Clustering"
12.20 - 12.50	"XAI for Random Forests"
12.50 - 13.00	Wrap-up and conclusions

Introduction

Terminology

Explainability or Interpretability?

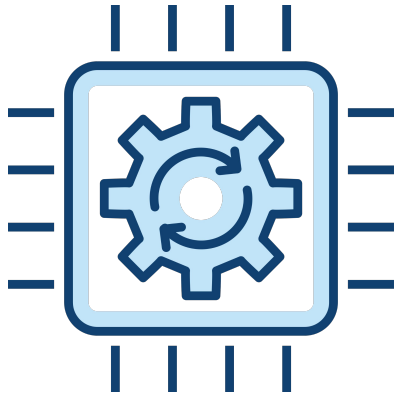


Introduction

Terminology

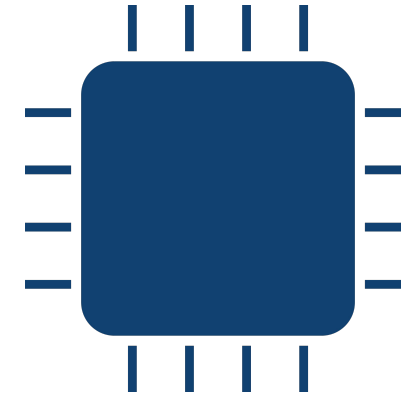
Interpretability

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.

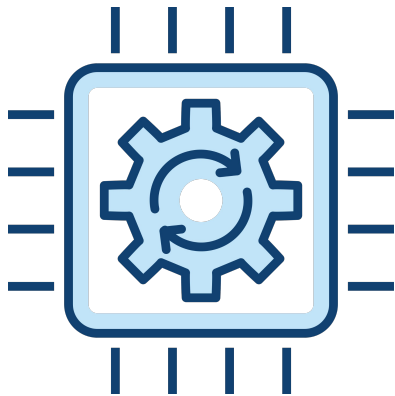


Introduction

Terminology

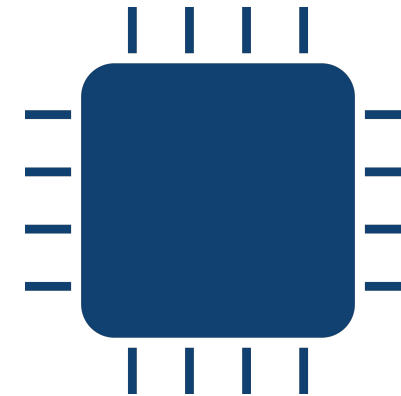
Interpretability

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.

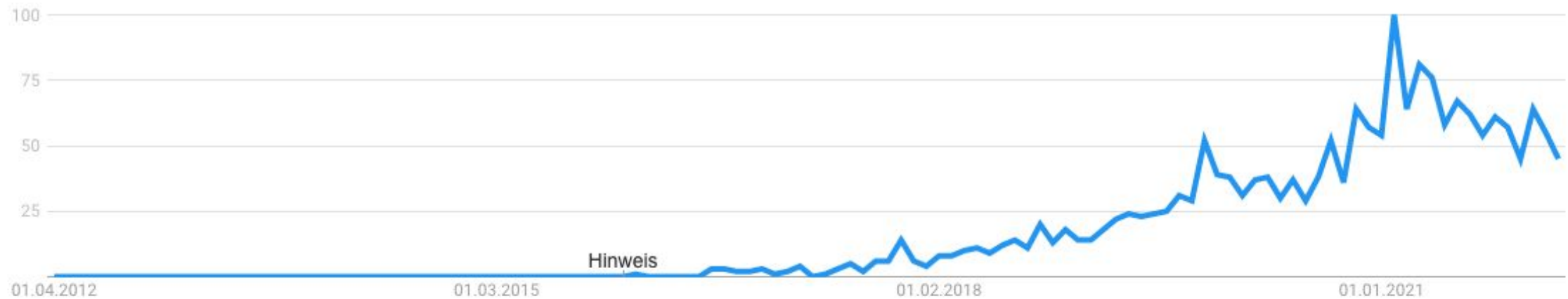


In this course, we will focus only on **eXplainable Artificial Intelligence (XAI)**.

Introduction

Why is explainability important?

Google Trends Popularity Index of the term *Explainable AI* over the last ten years (2012–2022)





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

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

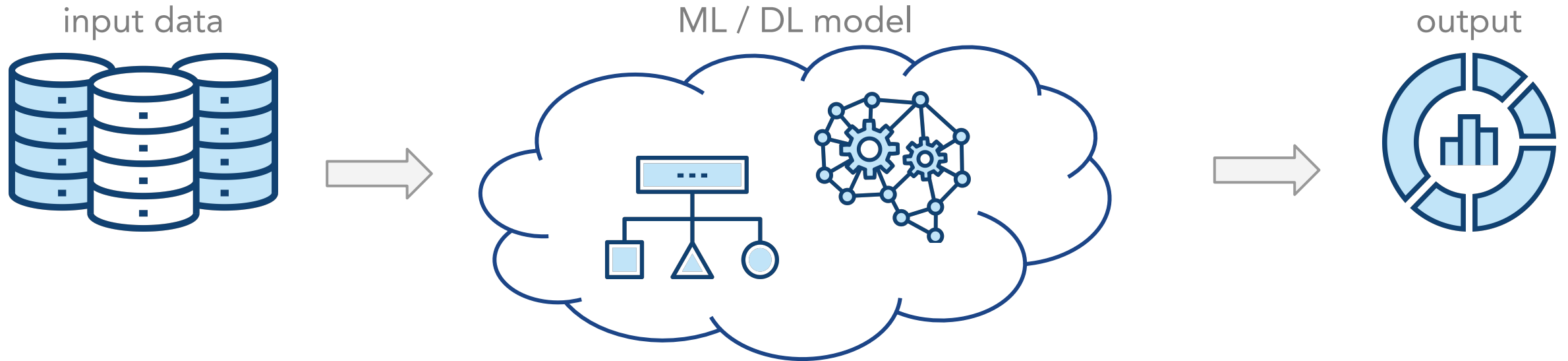
Impact story

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.

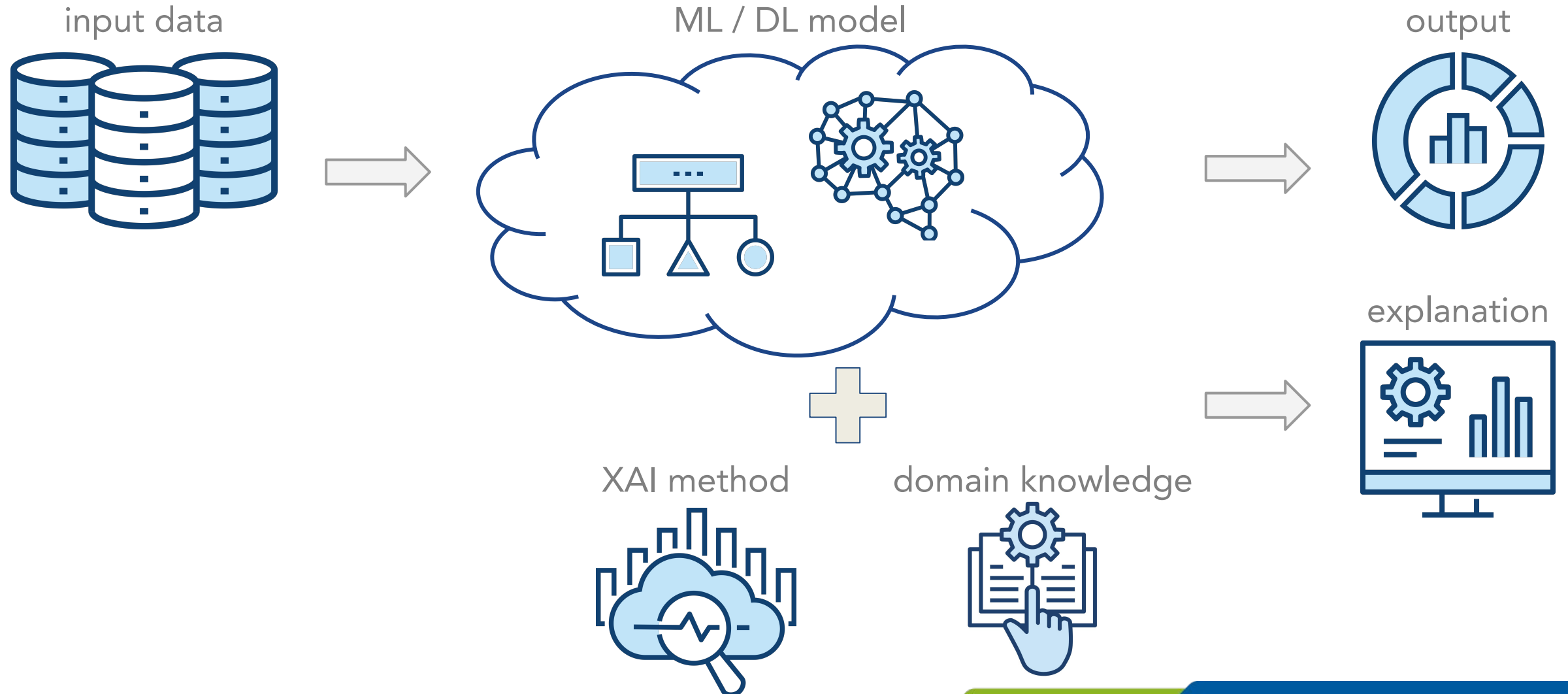
Introduction

XAI in your ML workflow



Introduction

XAI in your ML workflow



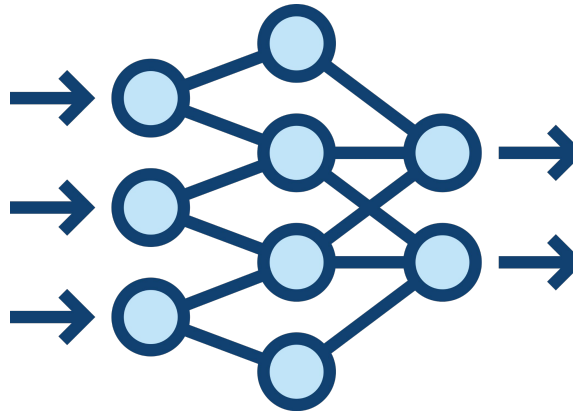
Introduction

XAI in your ML workflow

input data



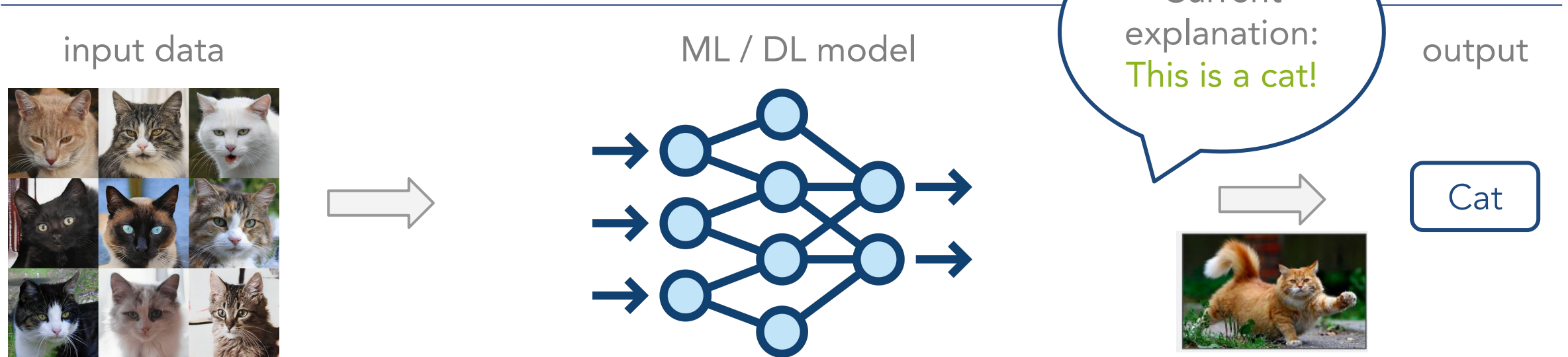
ML / DL model



output

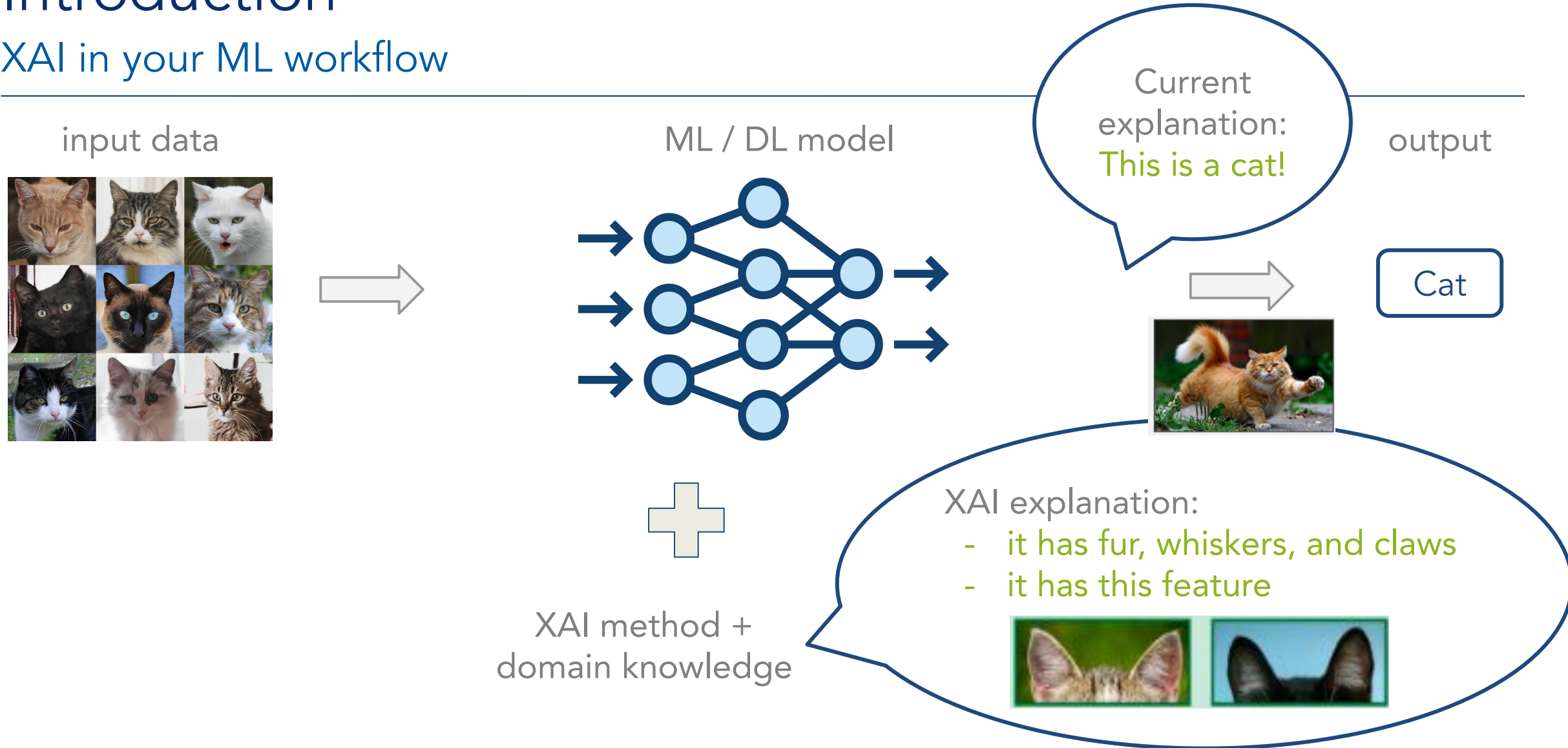
Introduction

XAI in your ML workflow



Introduction

XAI in your ML workflow



Introduction

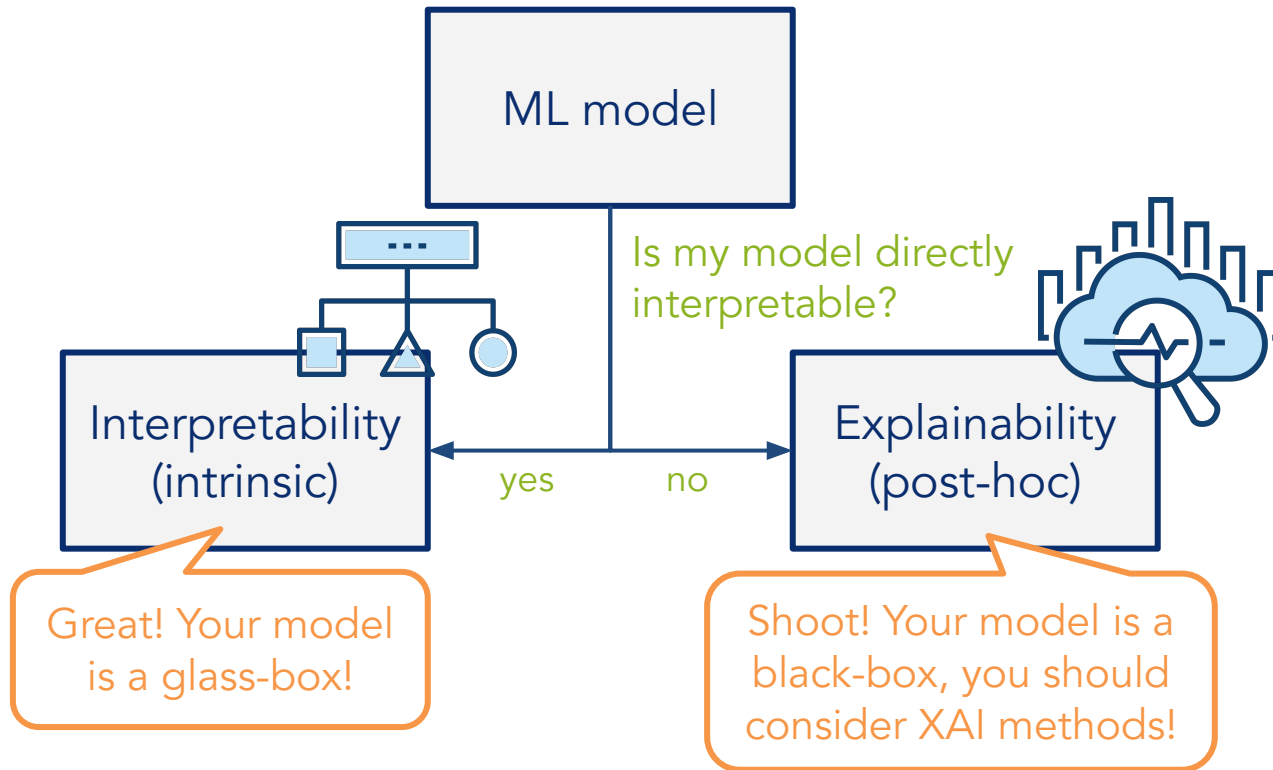
Taxonomy of XAI methods



ML model

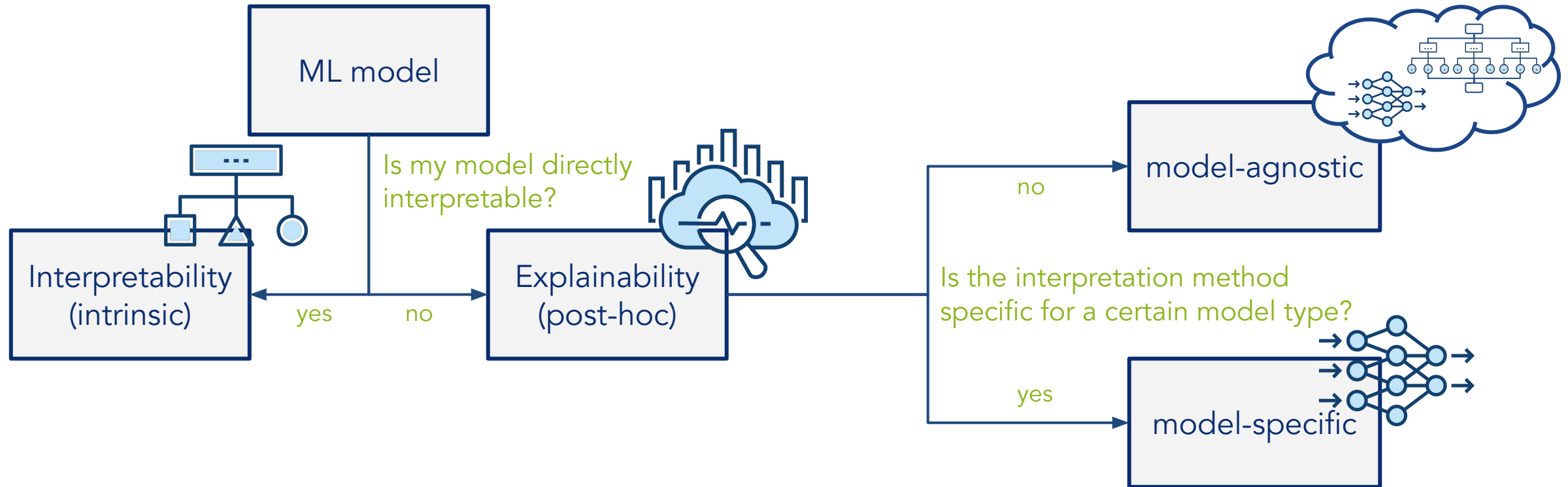
Introduction

Taxonomy of XAI methods



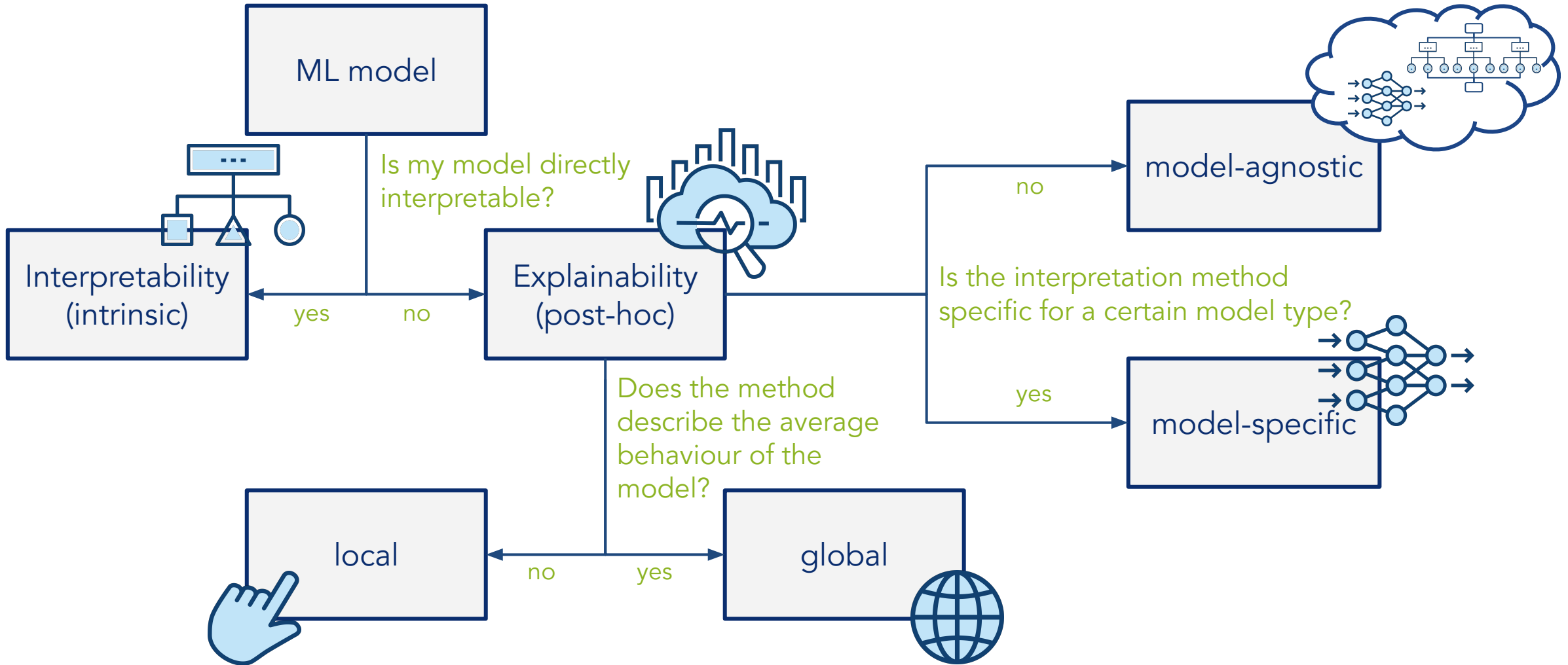
Introduction

Taxonomy of XAI methods



Introduction

Taxonomy of XAI methods



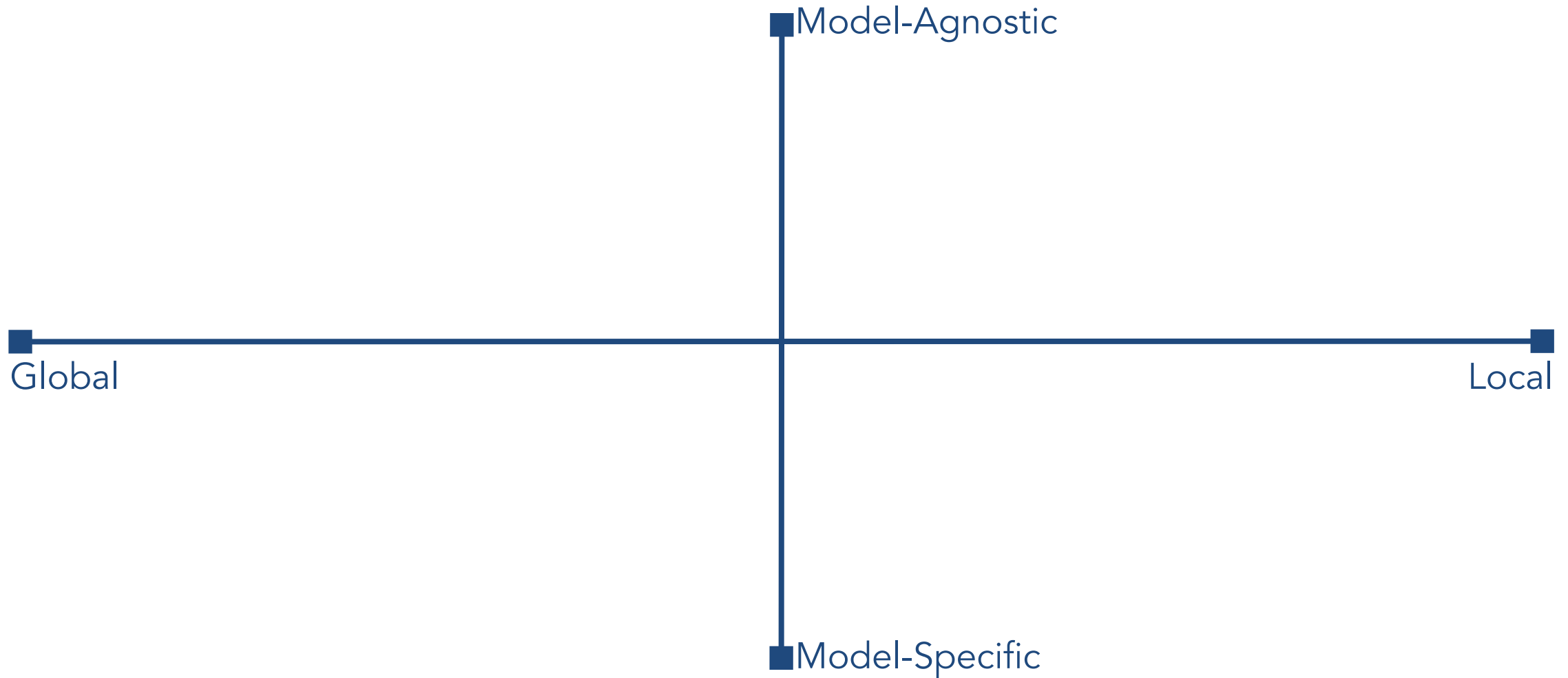


To understand what impact blood pressure has on the survival rate of patient John Doe in a Random Forest model, we need:

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

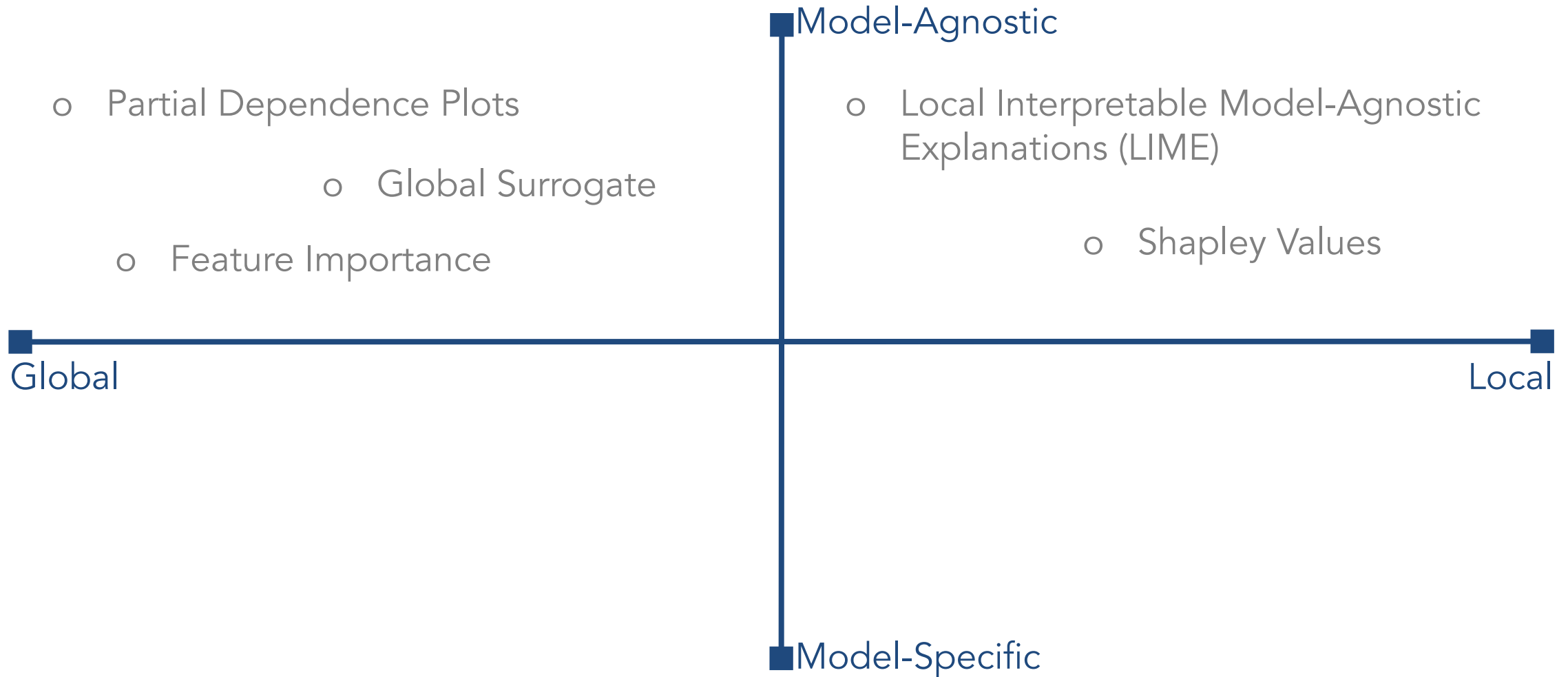
Introduction

Overview on post-hoc methods



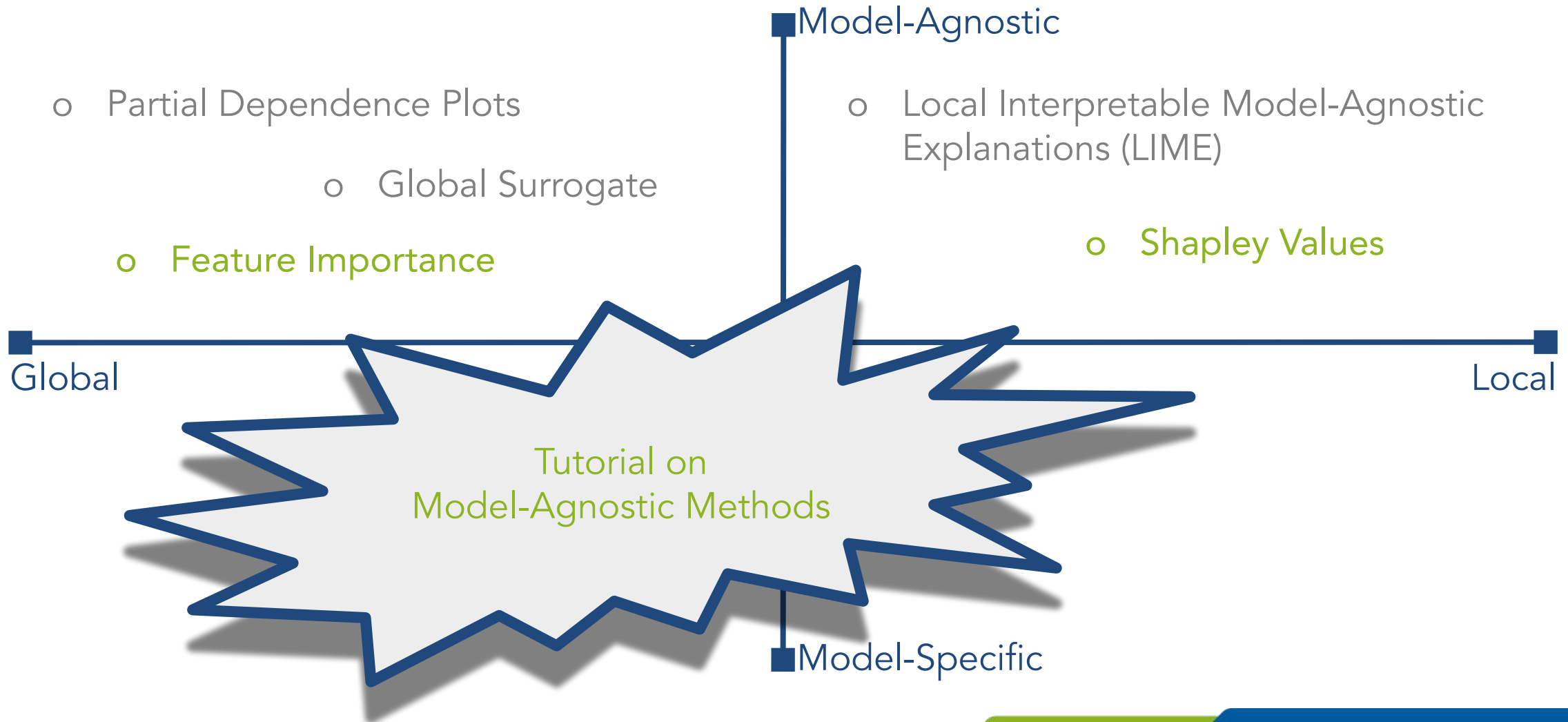
Introduction

Overview on post-hoc methods



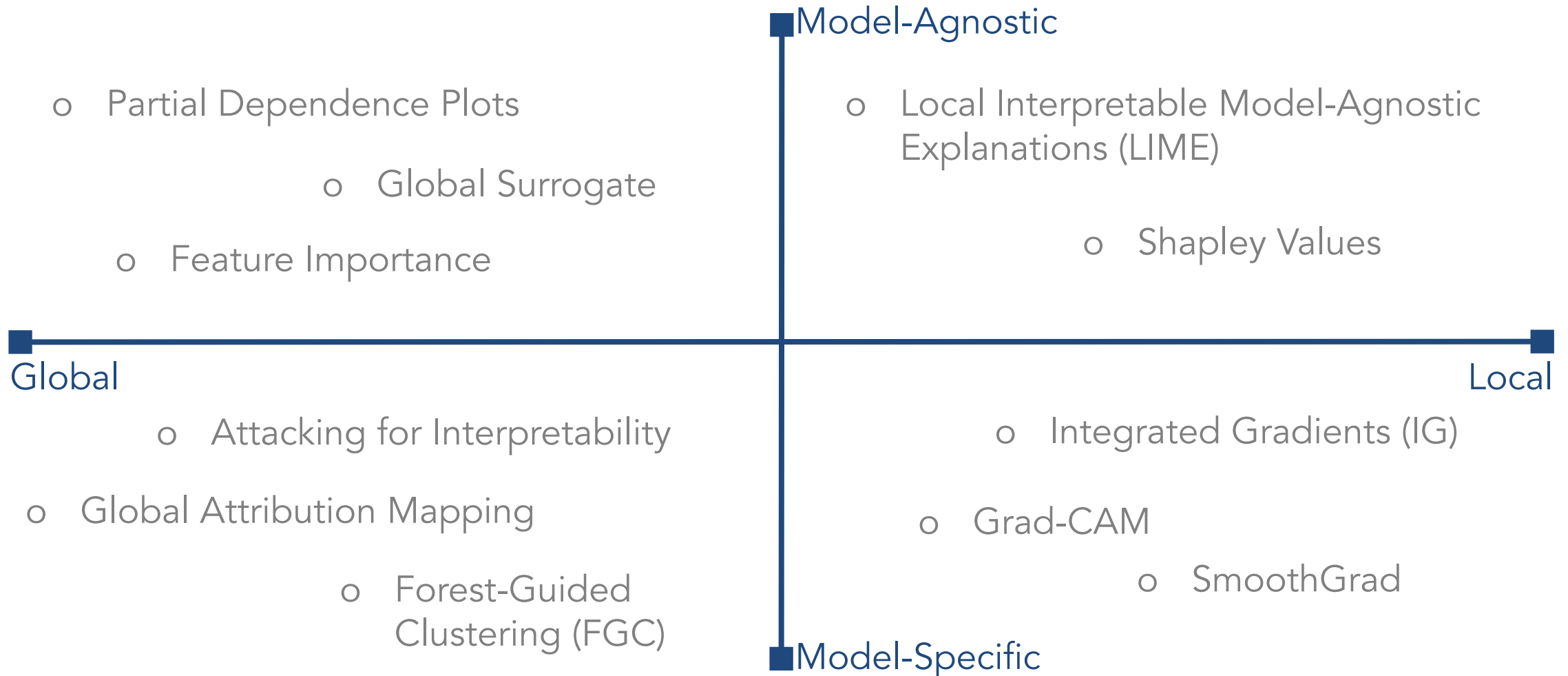
Introduction

Overview on post-hoc methods



Introduction

Overview on post-hoc methods



Introduction

Overview on post-hoc methods

