

Tópicos selectos de computación

María del Rocío Ochoa Montiel

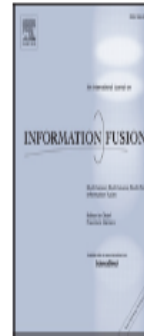
Contacto: ma.rocio.ochoa@gmail.com



Contents lists available at ScienceDirect

Information Fusion

journal homepage: www.elsevier.com/locate/inffus



Deep Learning and Security

DARPA's Explainable Artificial Intelligence Program

David Gunning, David W. Aha

Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI



44 AI MAGAZINE

Copyright © 2019, Association for the Advancement of Artificial Intelligence. All rights reserved. ISSN 0738-460

Alejandro Barredo Arrieta^a, Natalia Díaz-Rodríguez^b, Javier Del Ser^{a,c,d,*}, Adrien Bennetot^{b,e,f},
Siham Tabik^g, Alberto Barbado^h, Salvador Garcia^g, Sergio Gil-Lopez^a, Daniel Molina^g,
Richard Benjamins^h, Raja Chatila^f, Francisco Herrera^g

^aTECNALIA, Derio 48160, Spain

^bENSTA, Institute Polytechnique Paris and INRIA Flowers Team, Palaiseau, France

^cUniversity of the Basque Country (UPV/EHU), Bilbao 48013, Spain

^dBasque Center for Applied Mathematics (BCAM), Bilbao 48009, Bizkaia, Spain

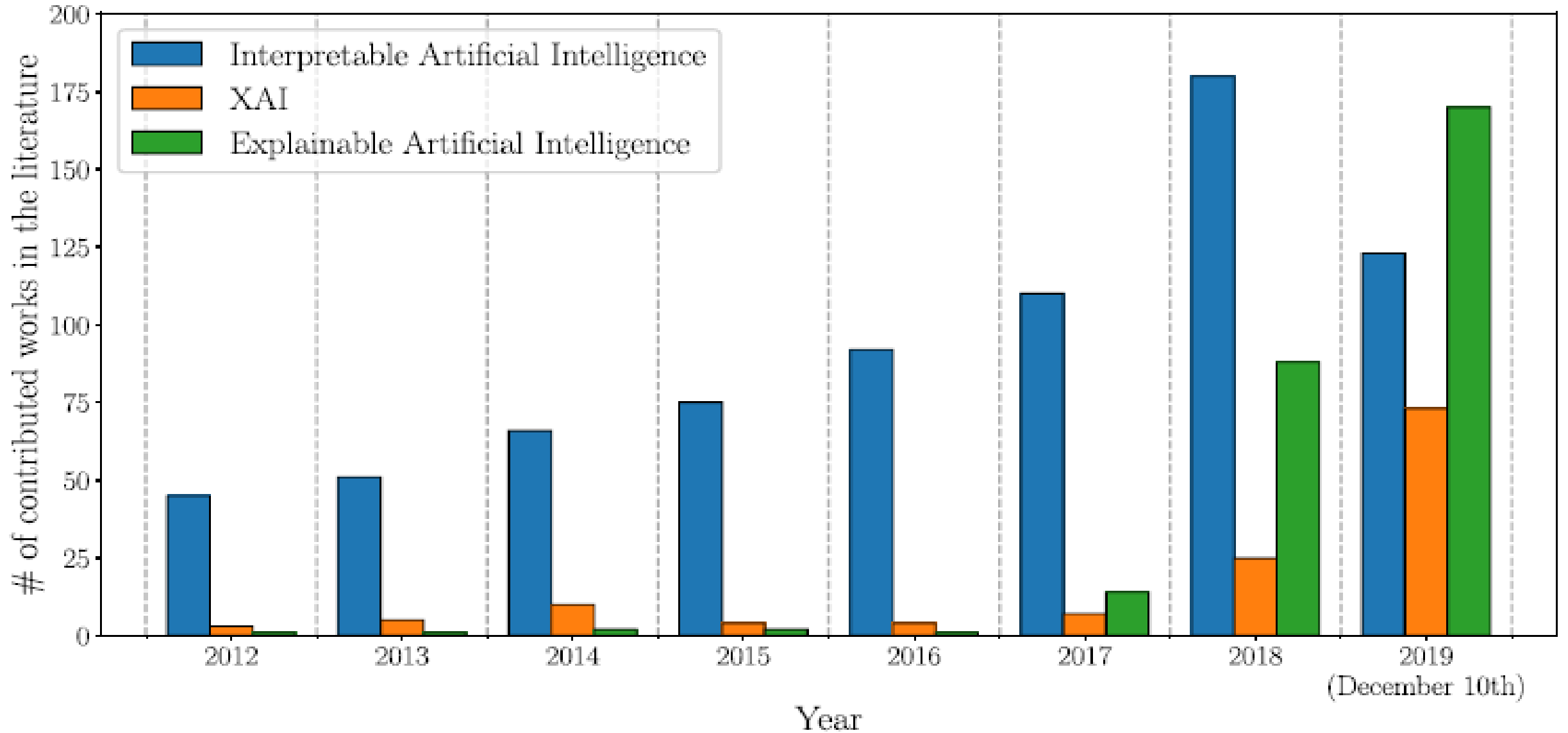
^eSegula Technologies, Parc d'activité de Pissaloup, Trappes, France

^fInstitut des Systèmes Intelligents et de Robotique, Sorbonne Université, France

^gDaSCI Andalusian Institute of Data Science and Computational Intelligence, University of Granada, Granada 18071, Spain

^hTelefonica, Madrid 28050, Spain

Evolution of XAI (eXplainable AI)



Purposes of explainability in ML(machine learning) models

Who? Domain experts/users of the model (e.g. medical doctors, insurance agents) ?
Why? Trust the model itself, gain scientific knowledge



Who? Users affected by model decisions
Why? Understand their situation, verify fair decisions...



Who? Regulatory entities/agencies
Why? Certify model compliance with the legislation in force, audits, ...



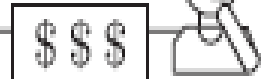
Target audience
in XAI



Who? Data scientists, developers, product owners...
Why? Ensure/improve product efficiency, research, new functionalities...



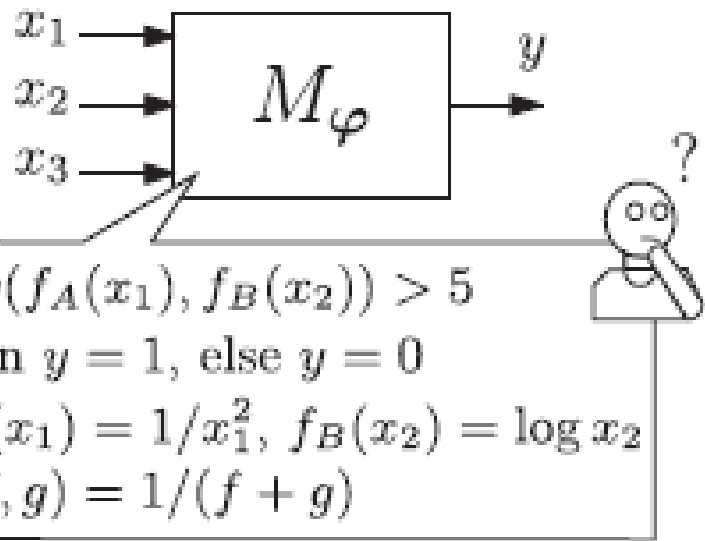
Who? Managers and executive board members
Why? Assess regulatory compliance, understand corporate AI applications...



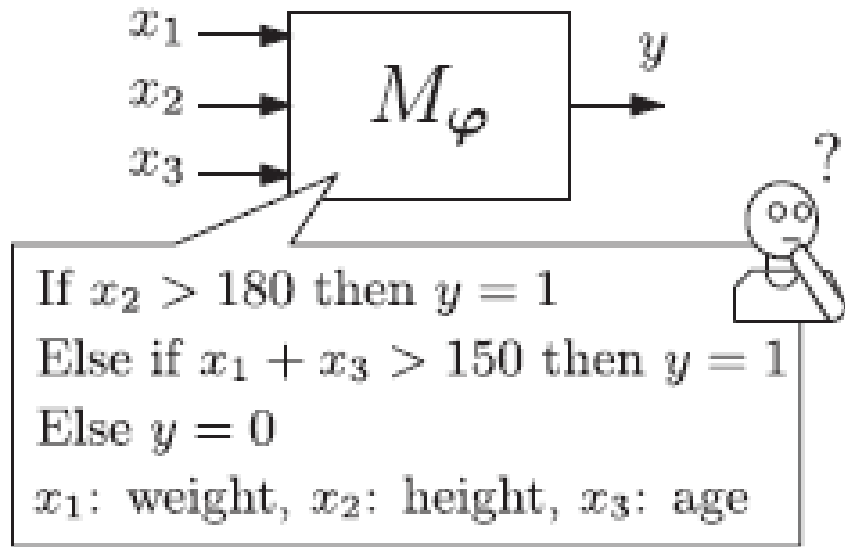
Goals pursued in the reviewed literature toward reaching explainability, and their main target audience

XAI Goal	Main target audience (Fig. 2)
Trustworthiness	Domain experts, users of the model affected by decisions
Causality	Domain experts, managers and executive board members, regulatory entities/agencies
Transferability	Domain experts, data scientists
Informativeness	All
Confidence	Domain experts, developers, managers, regulatory entities/agencies
Fairness	Users affected by model decisions, regulatory entities/agencies
Accessibility	Product owners, managers, users affected by model decisions
Interactivity	Domain experts, users affected by model decisions
Privacy awareness	Users affected by model decisions, regulatory entities/agencies

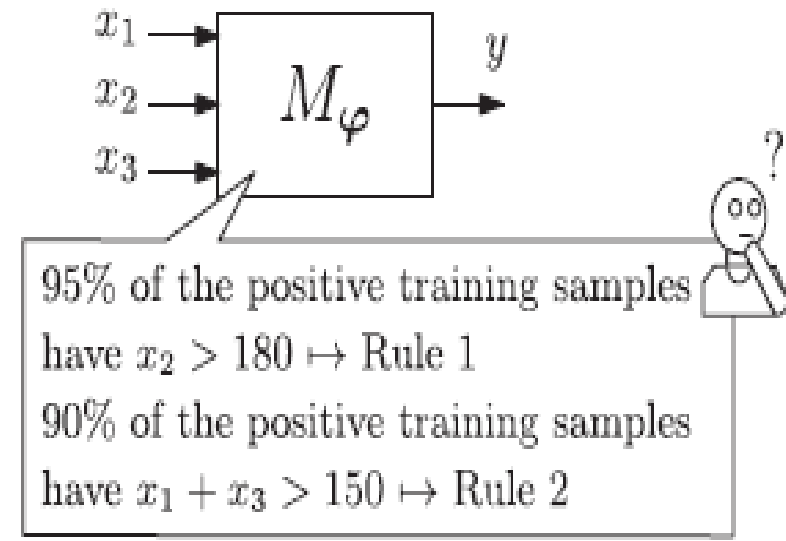
Levels of transparency characterizing a ML model



a) simulatability



b) decomposability

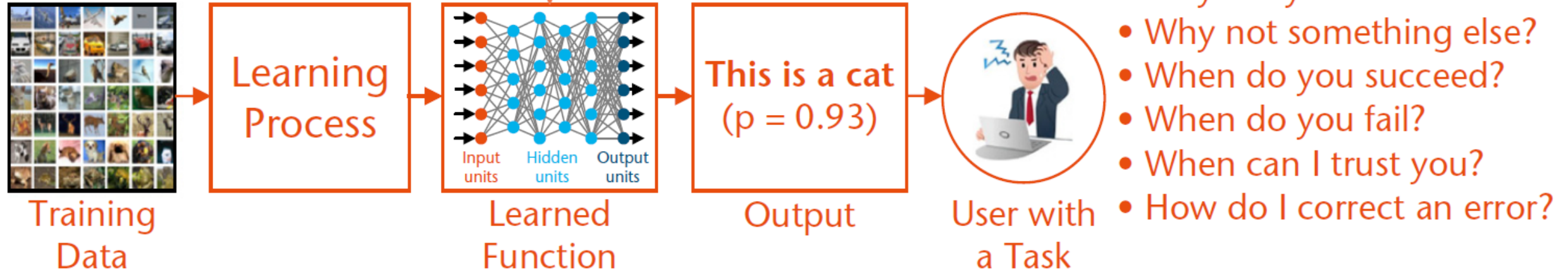


c) Algorithmic transparency

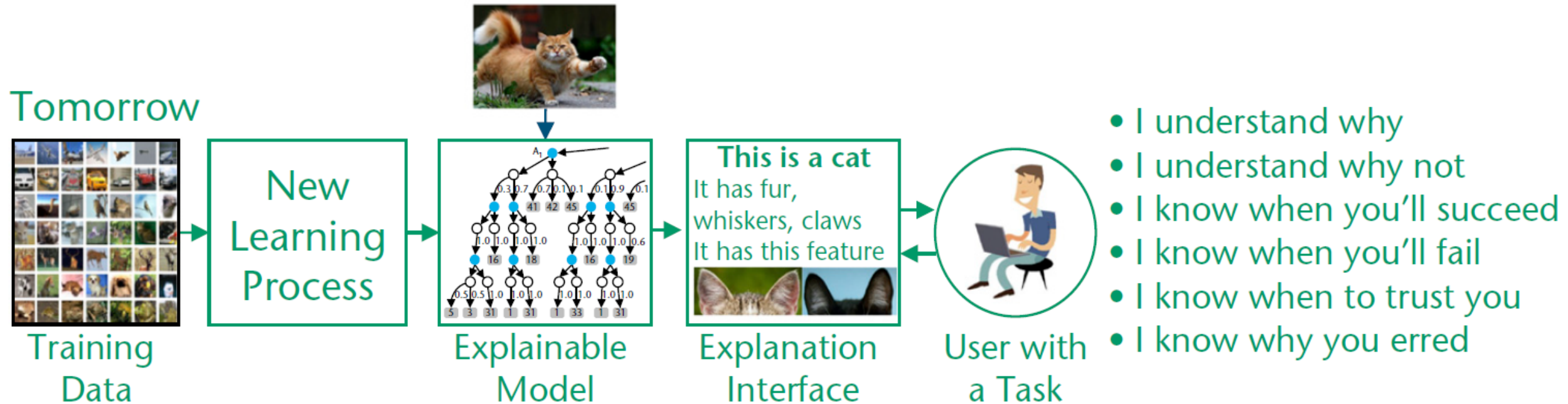
Targets for explainability may include:

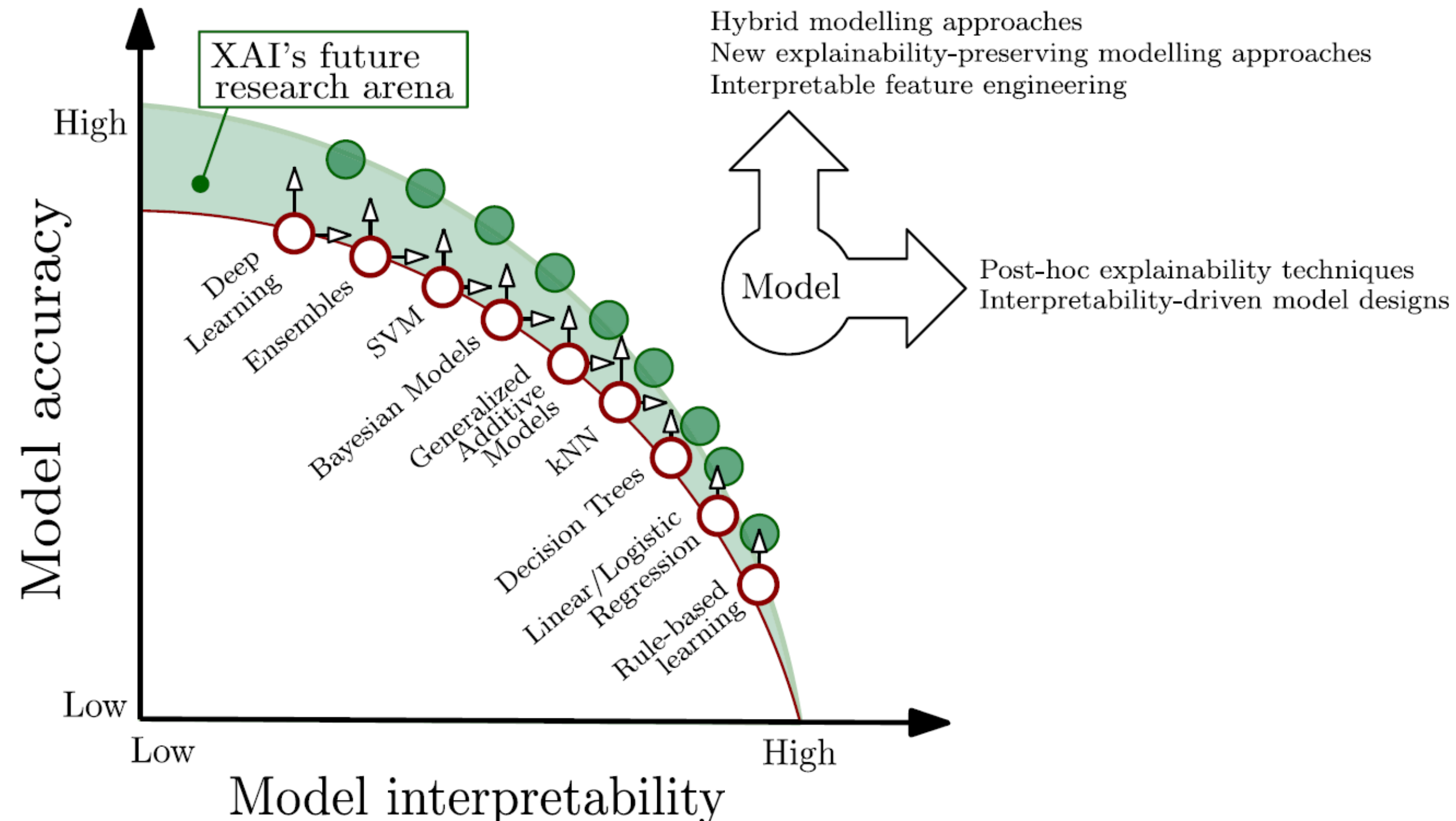
- a given example,
- the output classes, or
- the dataset itself

Today



Tomorrow





Post-hoc explainability approaches for a ML model

