



From Explanations to Trust

Human-centred approach AI research

Senka Krivić, PhD
Assistant Professor
Faculty of Electrical Engineering
University of Sarajevo

AI in media

The Wall Street Journal homepage. Headlines include: Qualcomm Won't Revive NXP Deal, After White House Flags China Concession; Trump: China to 'Reduce and Remove' Tariffs on American Cars; GlaxoSmithKline to Acquire Tesaro for \$4.16 Billion; Smaller Firms Finding Big Problems in China; Amazon Cashierless for Biggest; and CIO JOURNAL.

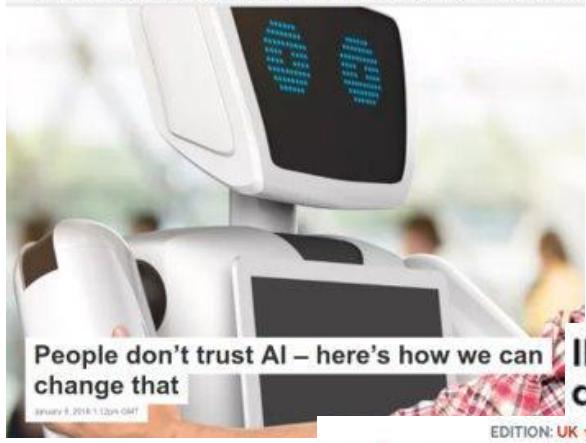
Tech Giants Launch New AI Tools as Worries Mount About Explainability

About 60% of 5,000 executives in ITM countries expressed concern with AI's 'black box'

THE CONVERSATION

Academic rigor. Journalistic fire.

Arts + Culture Business + Economy Cities Education Environment + Energy Health + Medicine Politics +



COMPUTERWORLD

HEALTH IT, CONVERGENCE

How do you make doctors trust machines in an AI-driven clinical world?

During a panel at the MedCity INVEST Twin Cities conference leaders from the payer, provider and investor spaces spoke about how to actually drive adoption of AI tools in the clinical system.

By KEVIN TRUONG

IBM Researchers propose transparency docs for AI services

other technologies and industries, artificial intelligence will need to adopt supplier's declaration of conformity agreements to build trust. How was that model built exactly?

ZDNet



Analytics & Optimization GDPR and Other Regulations Demand Explainable AI

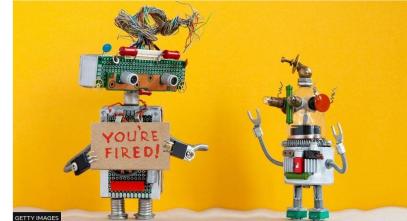
Save



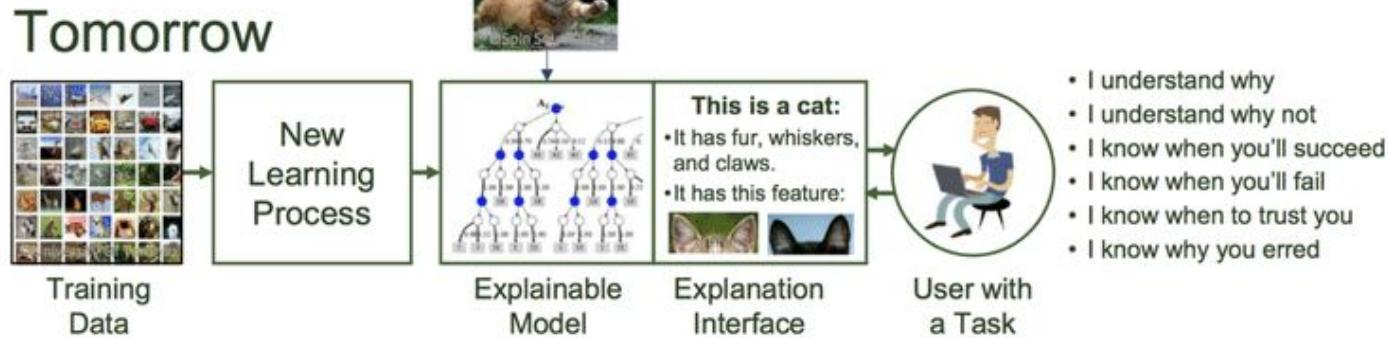
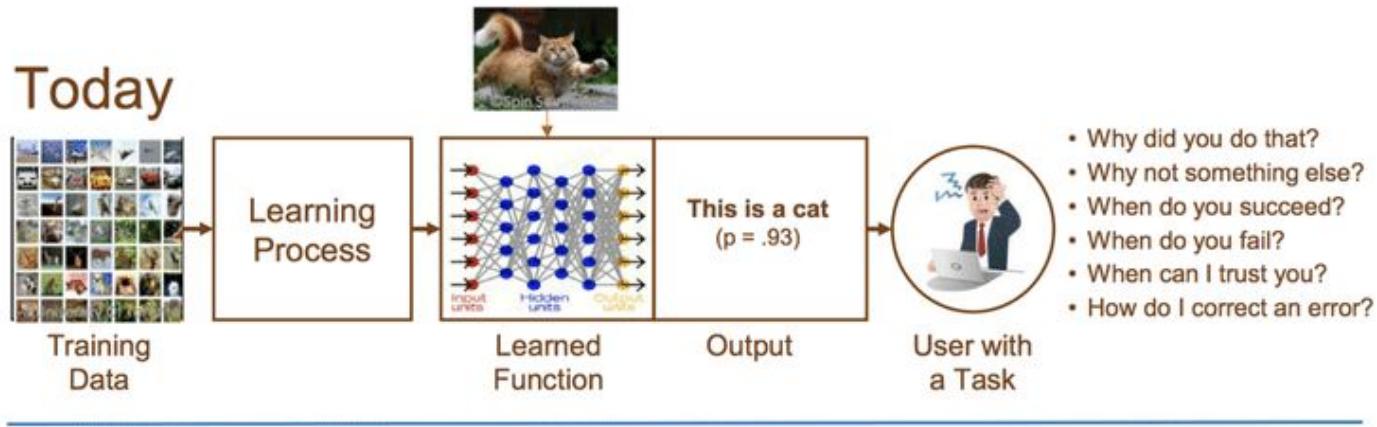
AI at work: Staff 'hired and fired by algorithm'

0 25 March 2021 | 0 Comments

3

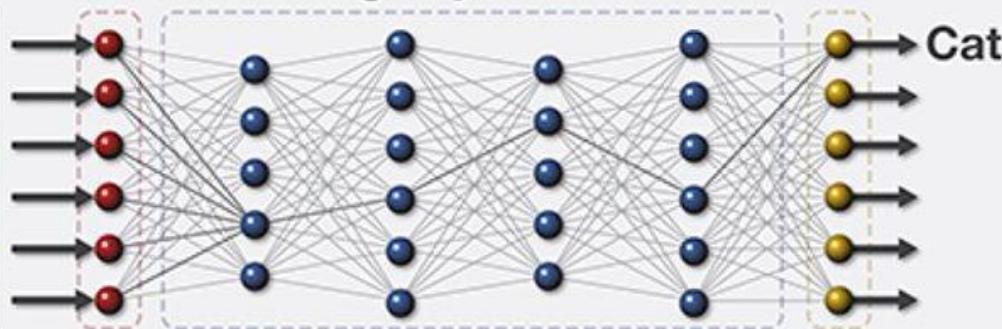


Explainable AI



XAI Explanation

Machine Learning System



This is a cat.

Current Explanation

This is a cat:

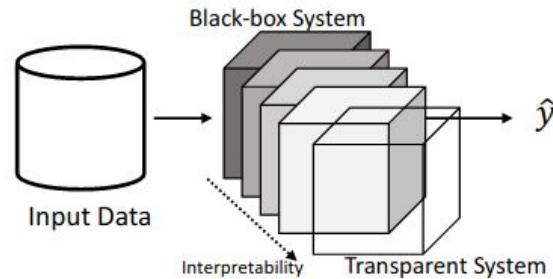
- It has fur, whiskers, and claws.
- It has this feature:



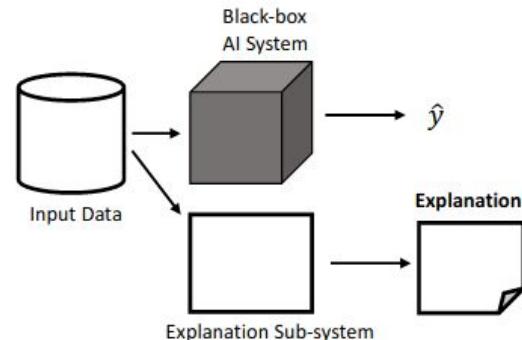
XAI Explanation

XAI Systems

Transparent-by-design systems



Post-hoc Explanation (black-box explanation) systems



[Mittelstadt et al. 2018]

Our work in XAI

- Human-Centred AI Lab in Sarajevo, UNSA
- Explainable Robotics
- Explainable and trusted AI Industry - Semiconductor industry use case
- AI-assisted learning



Ajla Karajko



Mubina
Kamberović



Vahidin Hasić



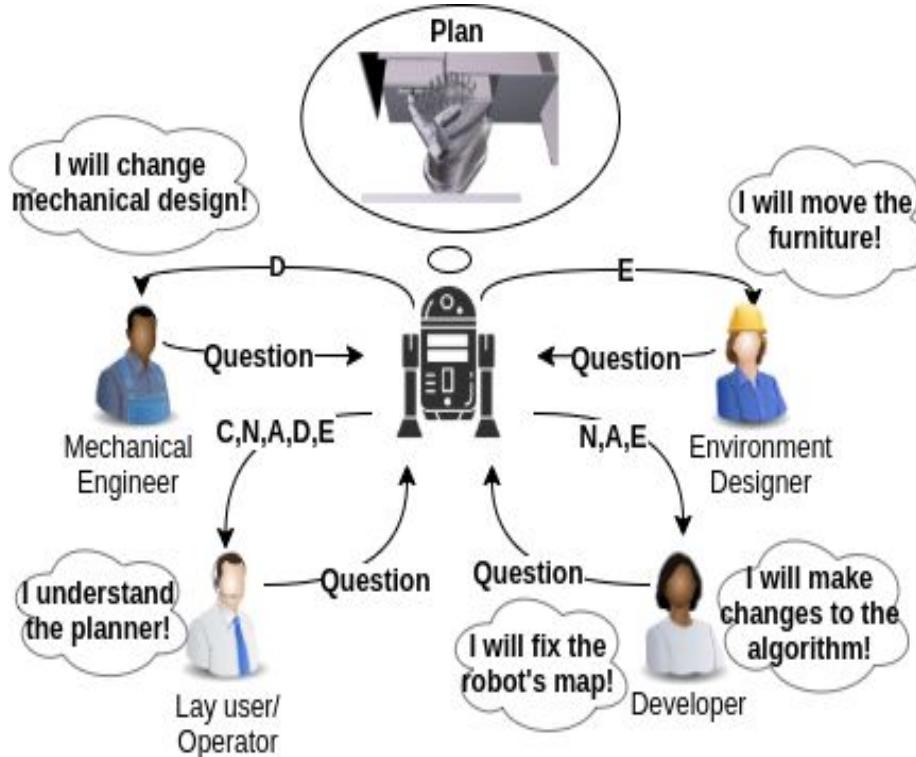
Amar Halilović



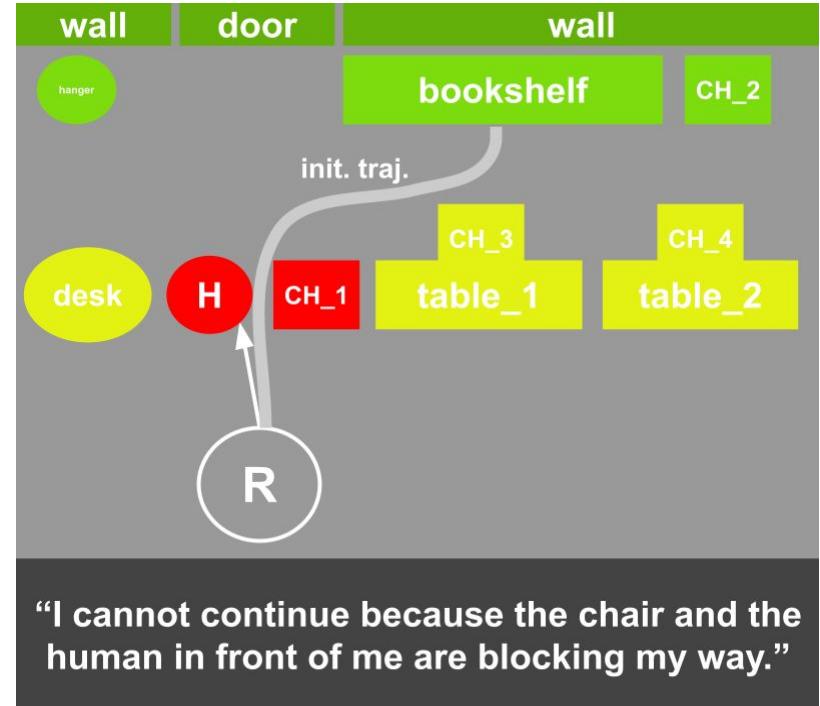
Amina Mević



Explainable AI and Robotics



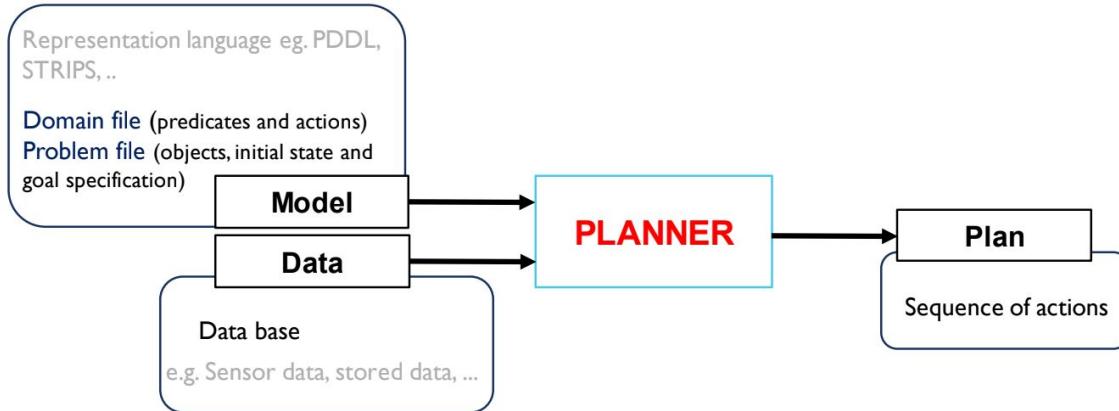
Explanation of Robot Navigation



“I cannot continue because the chair and the human in front of me are blocking my way.”

Planning of Explanations

AUTOMATED PLANNING AND SCHEDULING



- We propose a way of planning of explanations along with other robot actions.⁹
- The explanation occurrence can vary based on parameters such as user preference or other task priorities.



I am navigating.
Everything is
fine.

FAILURE!



I detected failure. I
cannot proceed.
What should I do?

**HUMAN IS
NEARBY!**



Human is nearby.
Maybe he can help
me. Let me explain
him/her what
happened.

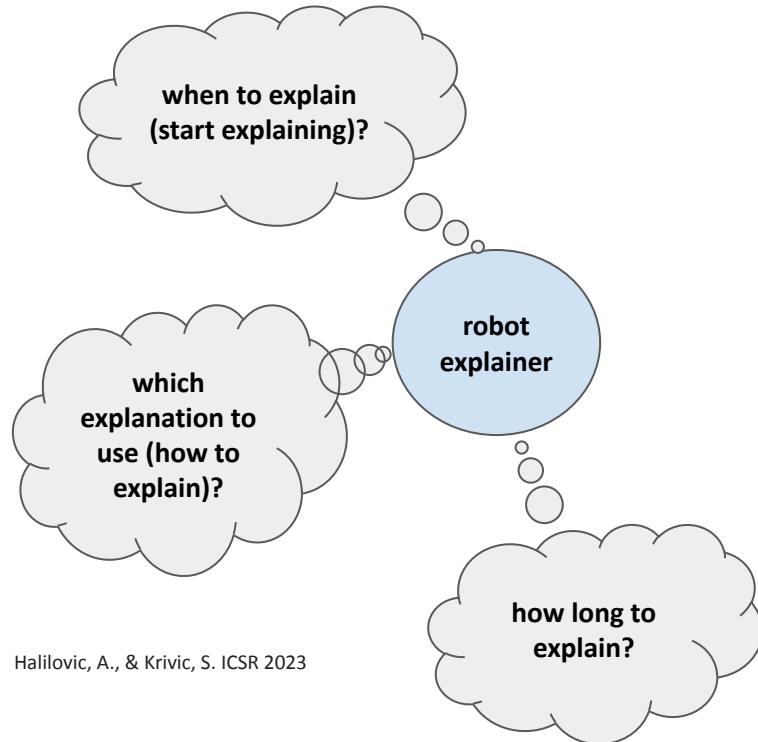
**I AM AN
INTROVERT**



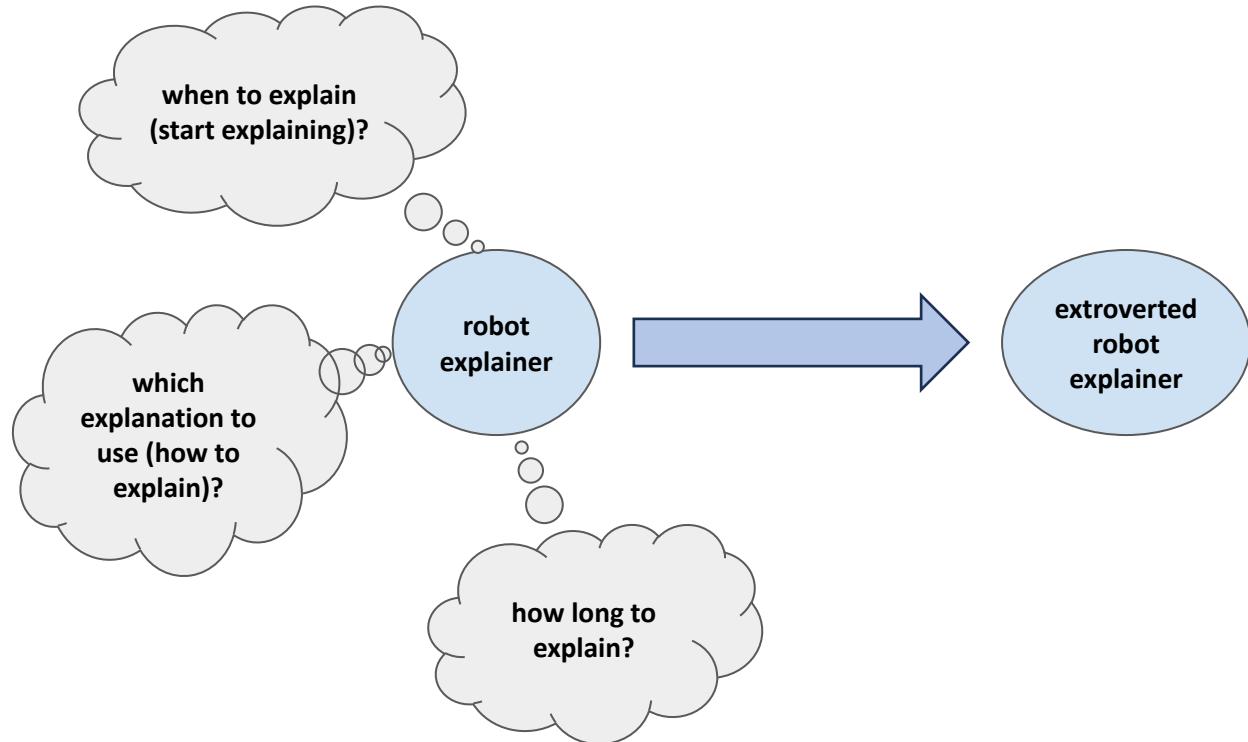
I am extroverted. I
will design my
explanation
accordingly.

**EXPLANA
TI
ON DESIGN
STARTED**

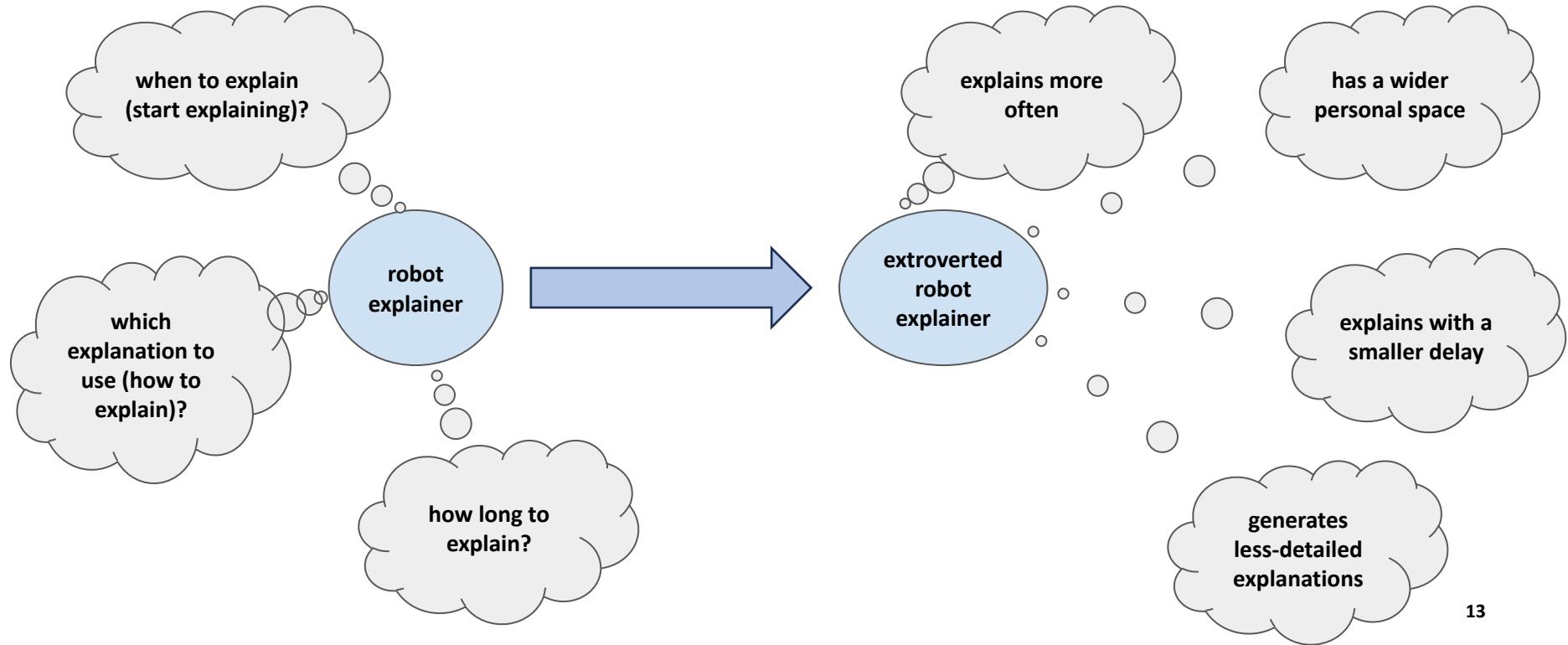
Influence of Robot Personality on Explanations



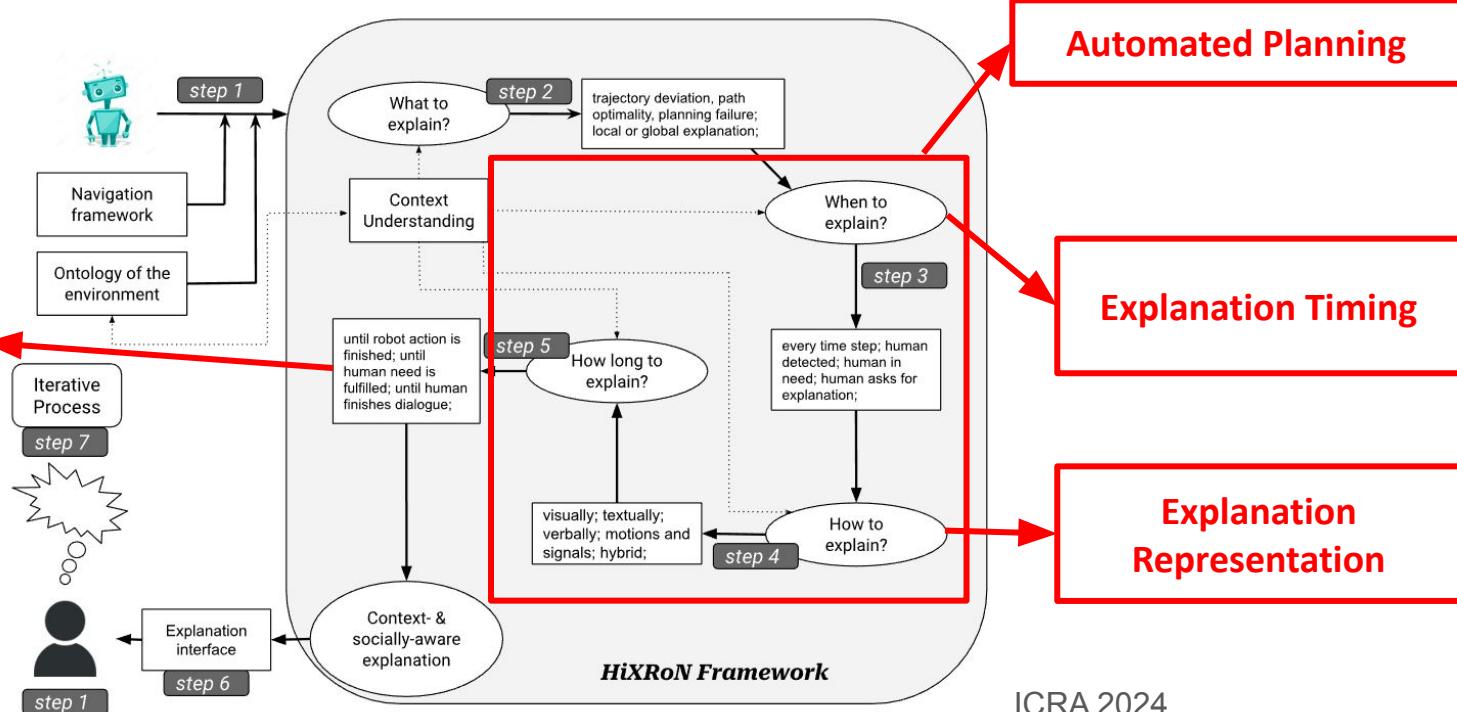
Influence of Robot Personality on Explanations



Influence of Robot Personality on Explanations



Explanation Generation



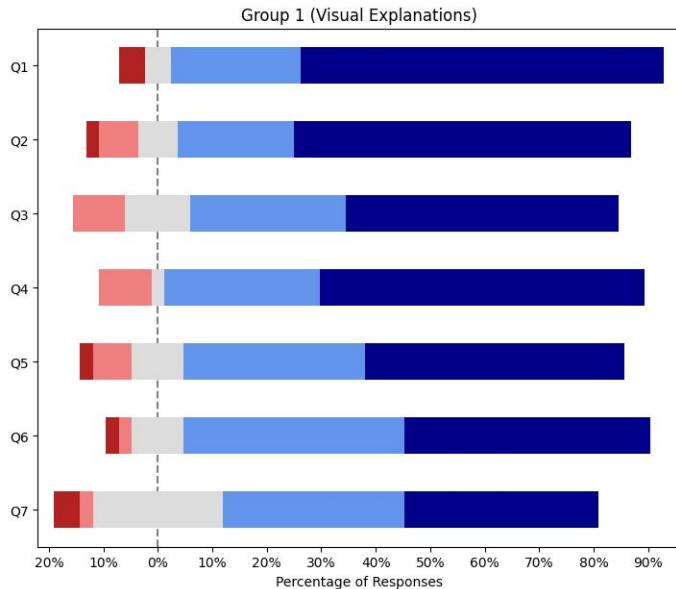
User Study on Explanation Representation

- Task: fetching a coffee from a coffee machine on a table
- Online study
- 84 participants (volunteers)
- Two groups:
 - Control group: visual explanations
 - Experimental group: visual-textual explanations
- User satisfaction scale used (Hoffman et al. 2018)*

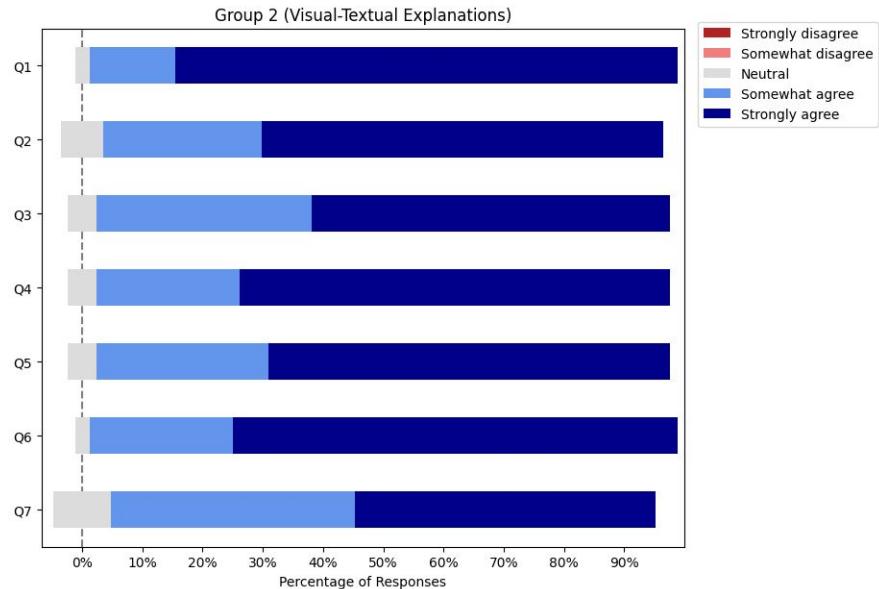


*Hoffman, R. R., Mueller, S. T., Klein, G., & Litman, J. (2018). Metrics for explainable AI: Challenges and prospects. arXiv preprint arXiv:1812.04608.

User study on explanation representation



The mean of all responses for Group 1 is 4.24 ($SD = 0.99$), corresponding to the overall attitude of “somewhat agree”.



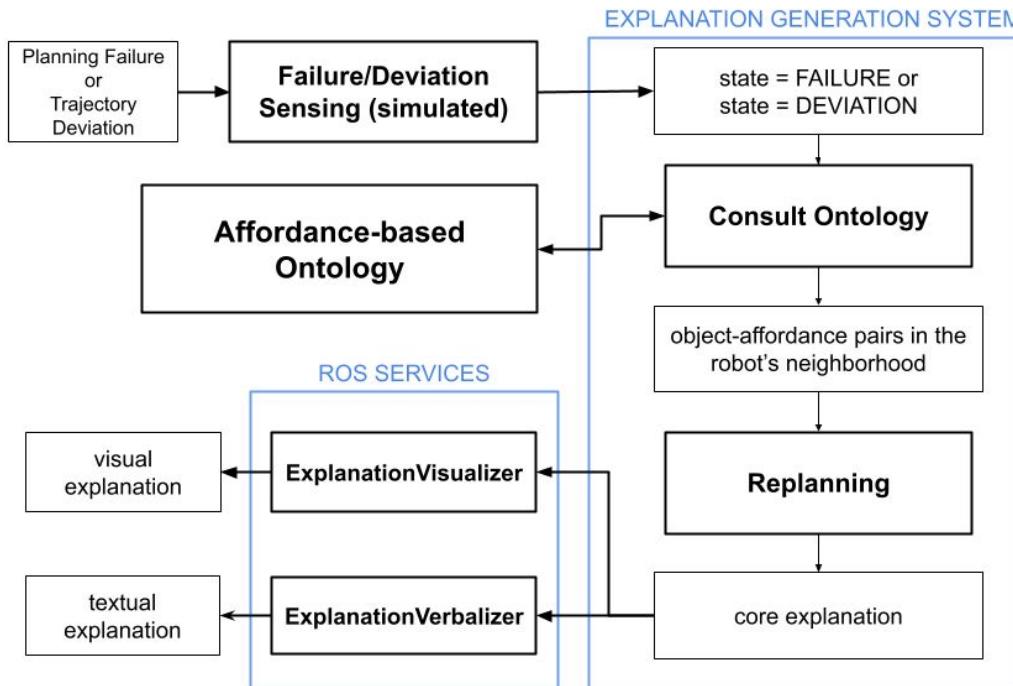
The mean of all responses for Group 2 is 4.62 ($SD = 0.58$), corresponding to the overall attitude of “strongly agree”.

Motivation for Affordance-based Explanations

- Lack of environmental explanations
- Affordances for explanations
- Service robot navigation in indoor social settings
- Visual and textual explanation embodiments

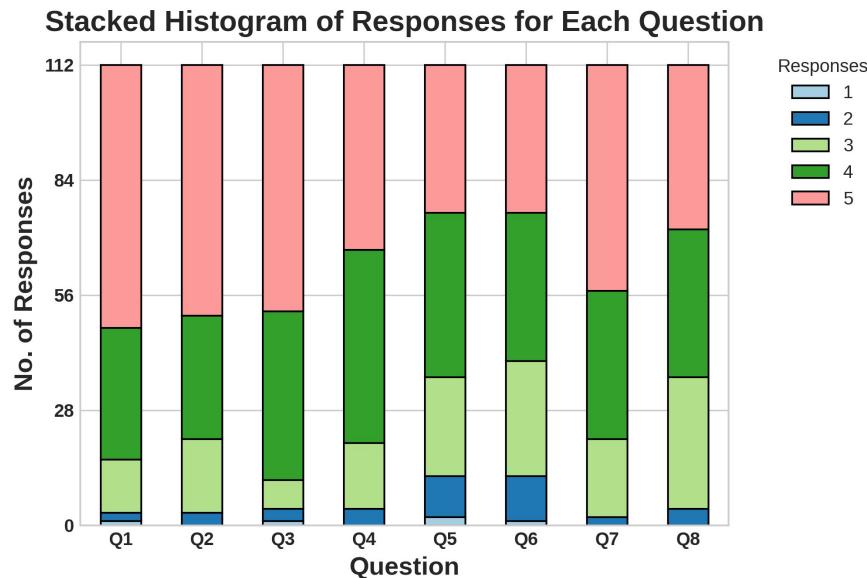


Generation of Affordance-Based Explanations



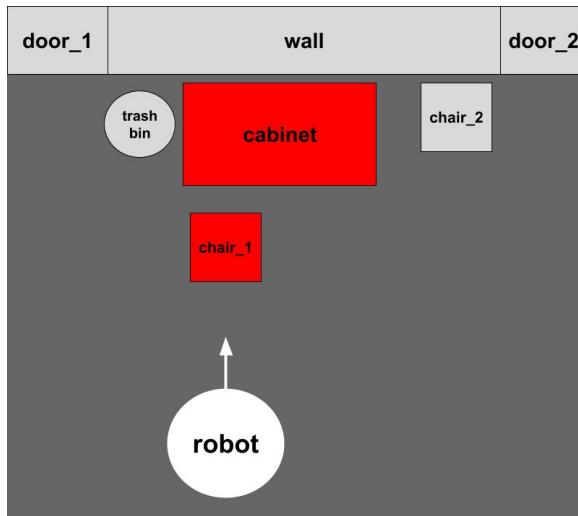
User study on explanation satisfaction

- Five different tasks (scenarios)
- Online study
- 112 participants (volunteers)
- One participant group
- User explanation satisfaction scale used (Hoffman et al. 2018)*

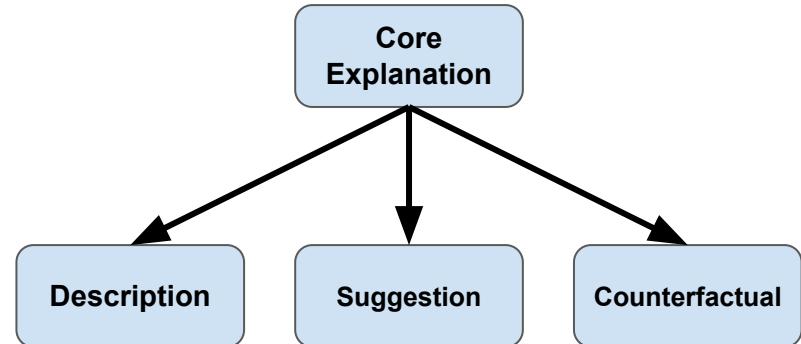


Explanation Visualization and Verbalization

Visualization

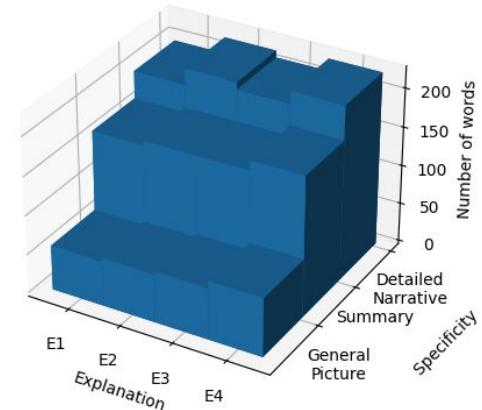
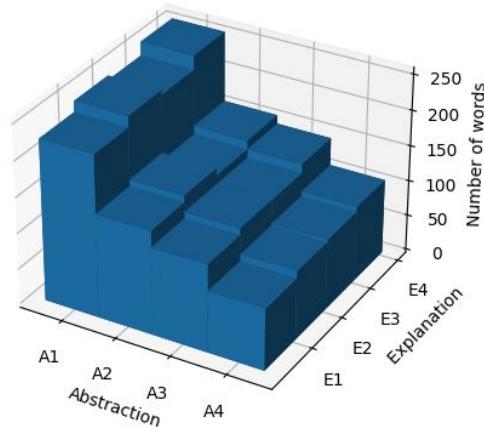


Verbalization



Verbalization experiments

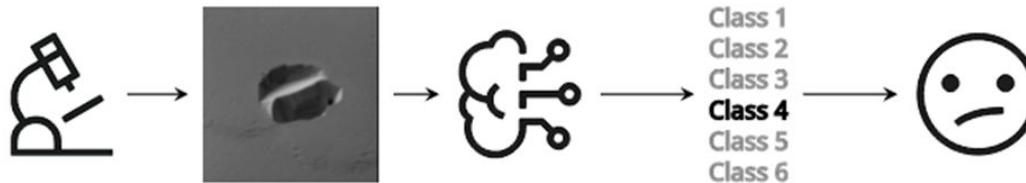
- Expanded the existing 3D verbalization framework* with affordance-based explanations
- Explanations do not increase the number of words in verbalizations



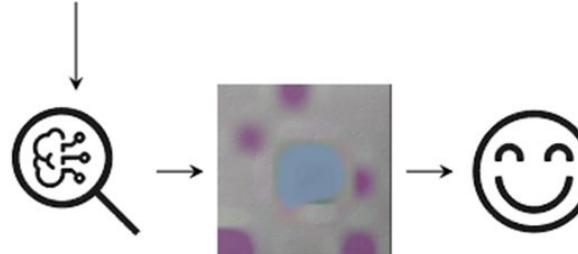
*Rosenthal, Stephanie, Sai P. Selvaraj, and Manuela M. Veloso. "Verbalization: Narration of Autonomous Robot Experience." IJCAI. Vol. 16. 2016.

*Canal, Gerard, et al. "PlanVerb: Domain-Independent Verbalization and Summary of Task Plans." Proceedings of the AAAI Conference on Artificial Intelligence. Vol. 36. No. 9. 2022.

Explaining Scanning Electron Microscope (SEM) Defect Image Classification



SEM inspection	Defect image	Classification model	Predicted class	Unsatisfied expert
----------------	--------------	----------------------	-----------------	--------------------



XAI method	Explained prediction	Satisfied expert
------------	----------------------	------------------

ASE - Any Segment Explanation

- Traditional XAI relies on individual pixel-based explanations, which highlight significant pixels but may not align with human intuition
- Concept-based explanations are more human-understandable but typically depend on human-annotated concepts
- **ASE - black-box method that automatically extracts concepts and assigns importances based on concept insertion and deletion**

Prediction output: goldenfish

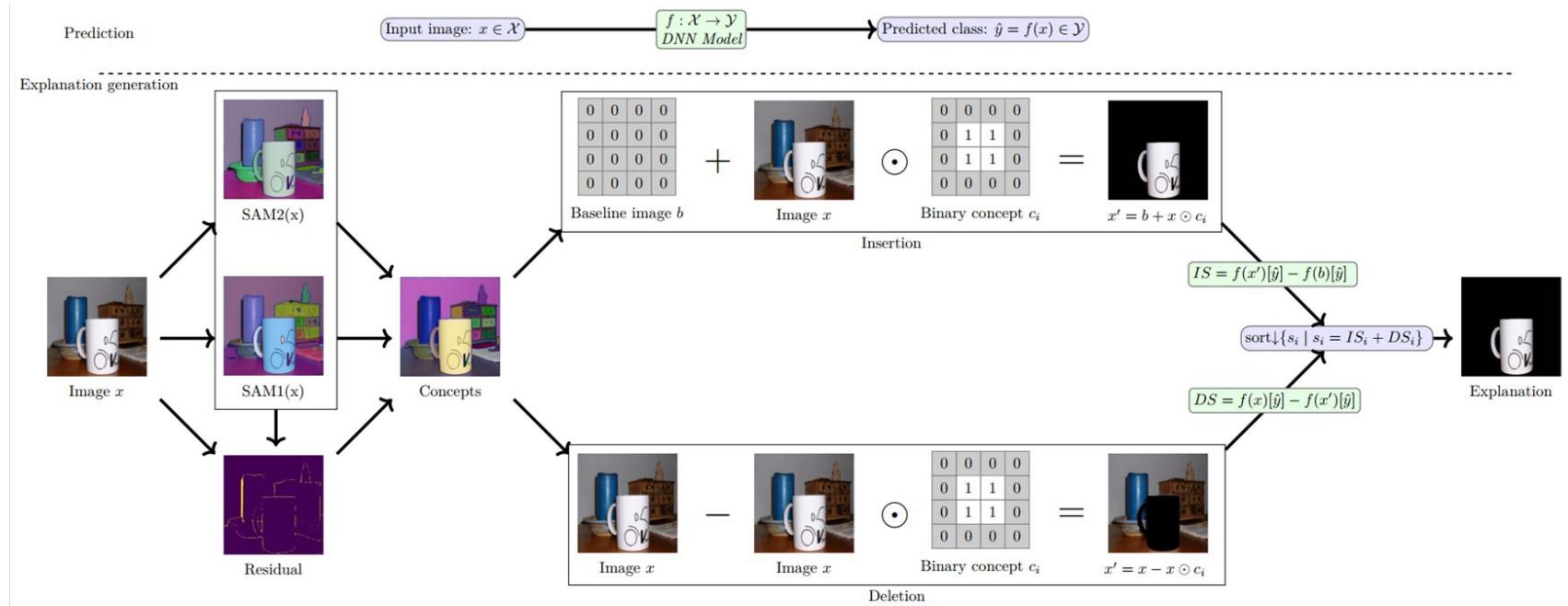


Original image

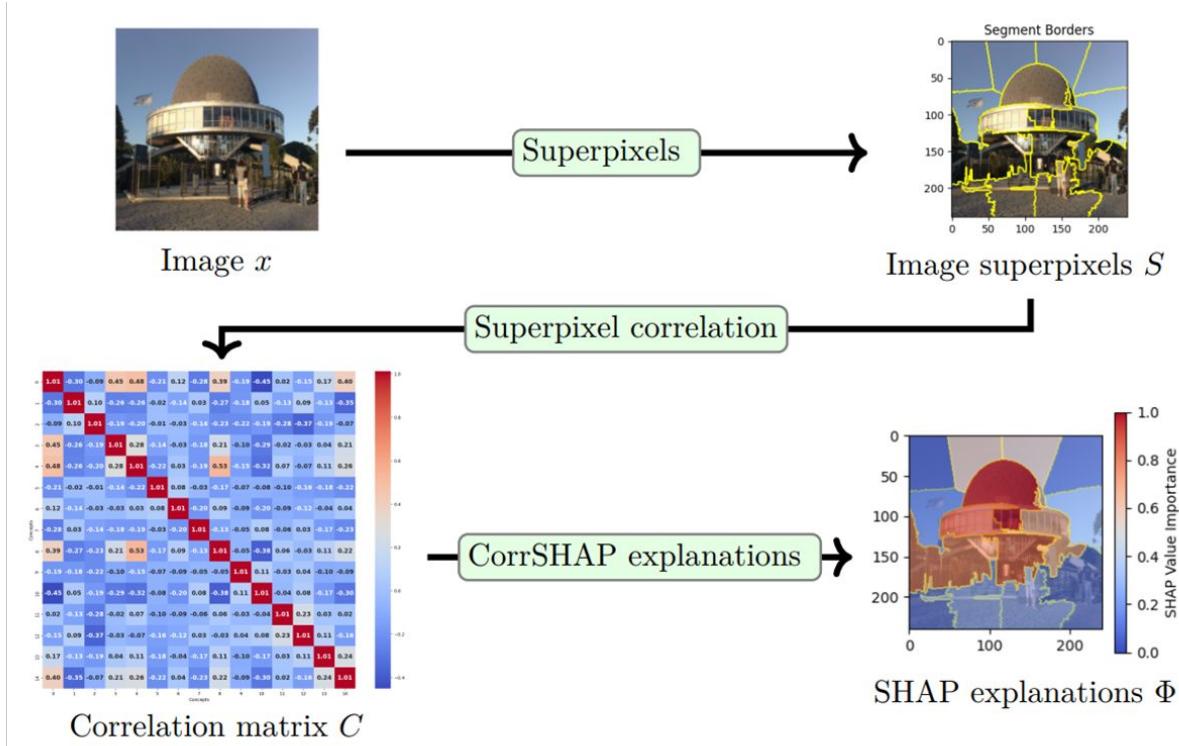


ASE explanation

ASE - Any Segment Explanation



CorrSHAP

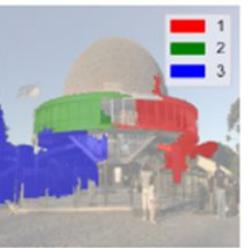


CorrSHAP - Results

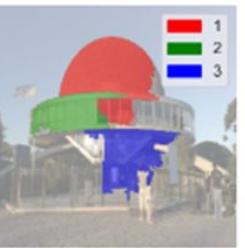
Original image



MCSHAP



CorrSHAP

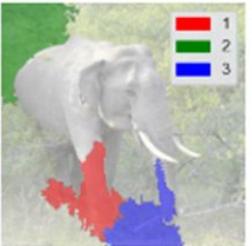


Model prediction: apairy

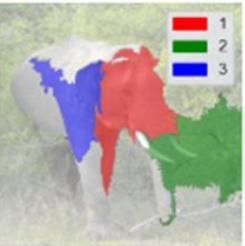
Original image



MCSHAP



CorrSHAP



Model prediction: elephant

Area Under the Curve (AUC) Insertion ↑

Model	Superpixels	CorrSHAP 1	CorrSHAP 2	CorrSHAP 3	MCSHAP
MobileNet-v2	Quickshift	80.4	80.37	80.29	80.89
	SLIC	78.12	78.13	78.15	77.79
ResNet-18	Quickshift	60.21	60.23	60.21	61.60
	SLIC	54.66	54.65	54.65	55.41
ResNet-50	Quickshift	82.63	82.61	82.63	82.27
	SLIC	81.20	81.20	81.20	80.66
ViT-b16	Quickshift	80.83	80.83	80.74	81.84
	SLIC	76.36	76.33	76.41	76.82

Area Under the Curve (AUC) Deletion ↓

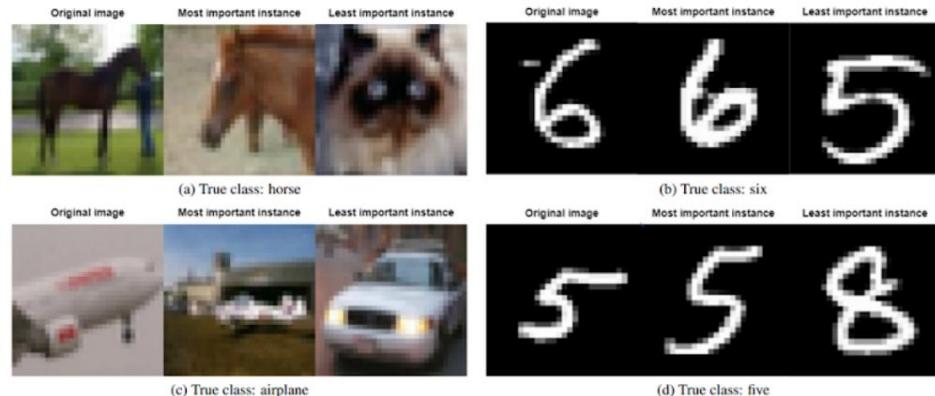
Model	Quickshift	20.14	20.14	20.16	20.93
MobileNet-v2	SLIC	19.48	19.48	19.48	21.40
	Quickshift	8.25	8.25	8.25	8.03
ResNet-18	SLIC	9.06	9.06	9.06	9.01
	Quickshift	22.79	22.79	22.79	24.16
ResNet-50	SLIC	22.36	22.36	22.39	23.79
	Quickshift	17.20	17.17	17.13	16.86
ViT-b16	SLIC	20.53	20.53	20.53	20.48

Execution Time (seconds) ↓

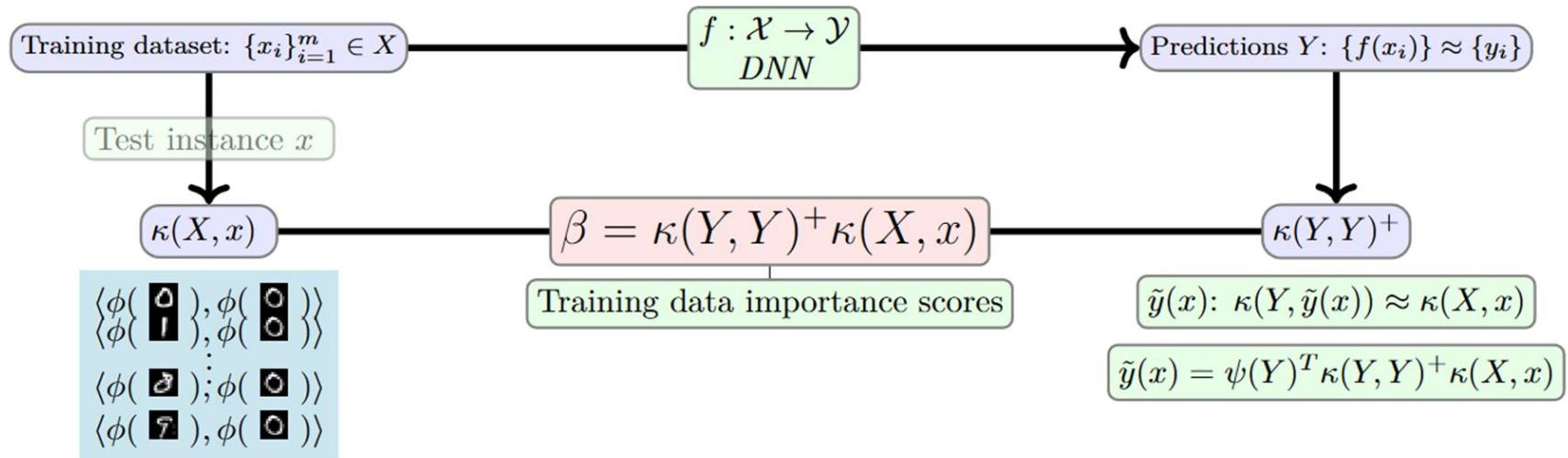
Model	Quickshift	0.42	0.54	1.15	16.13
MobileNet-v2	SLIC	0.74	0.89	1.86	14.98
	Quickshift	0.43	0.50	1.09	7.85
ResNet-18	SLIC	0.66	0.51	1.91	13.82
	Quickshift	0.46	0.58	2.76	25.16
ResNet-50	SLIC	0.78	0.93	5.93	36.49
	Quickshift	0.40	0.52	5.41	7.26
ViT-b16	SLIC	0.73	1.01	10.53	18.00

Kernel Sample-Based Explanations

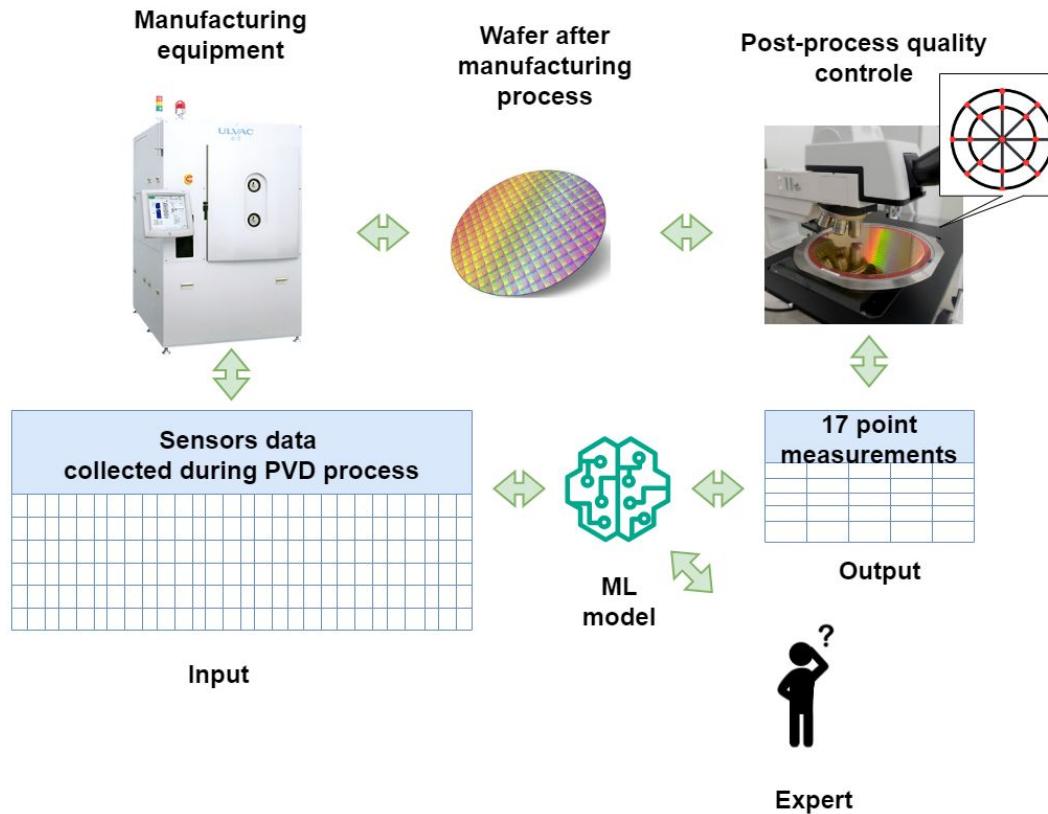
- Knowing the influence of individual training instances on the predictions has a lot of applications:
 - detecting mislabeled data
 - identifying data leakage
 - analyzing memorization effects
 - optimizing training dataset
 - control theory
 - active learning
 - system identification
- **Kernel Sample Based Explanations (K-SBE) - novel model-agnostic approach built upon kernel functions**



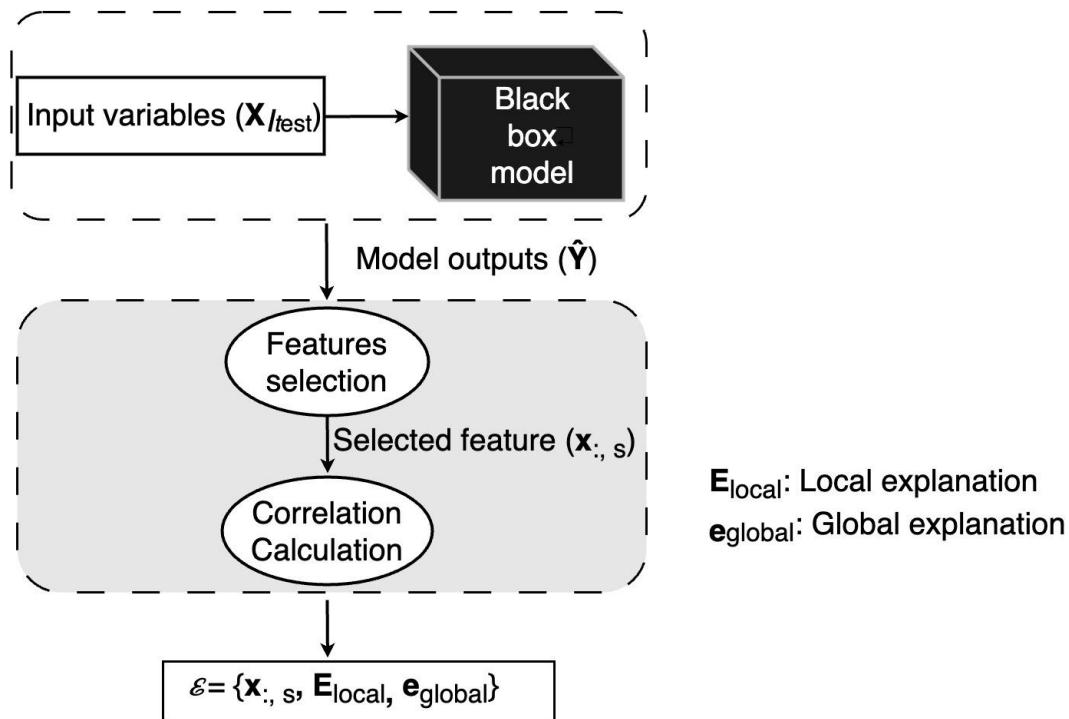
Kernel Sample-Based Explanations



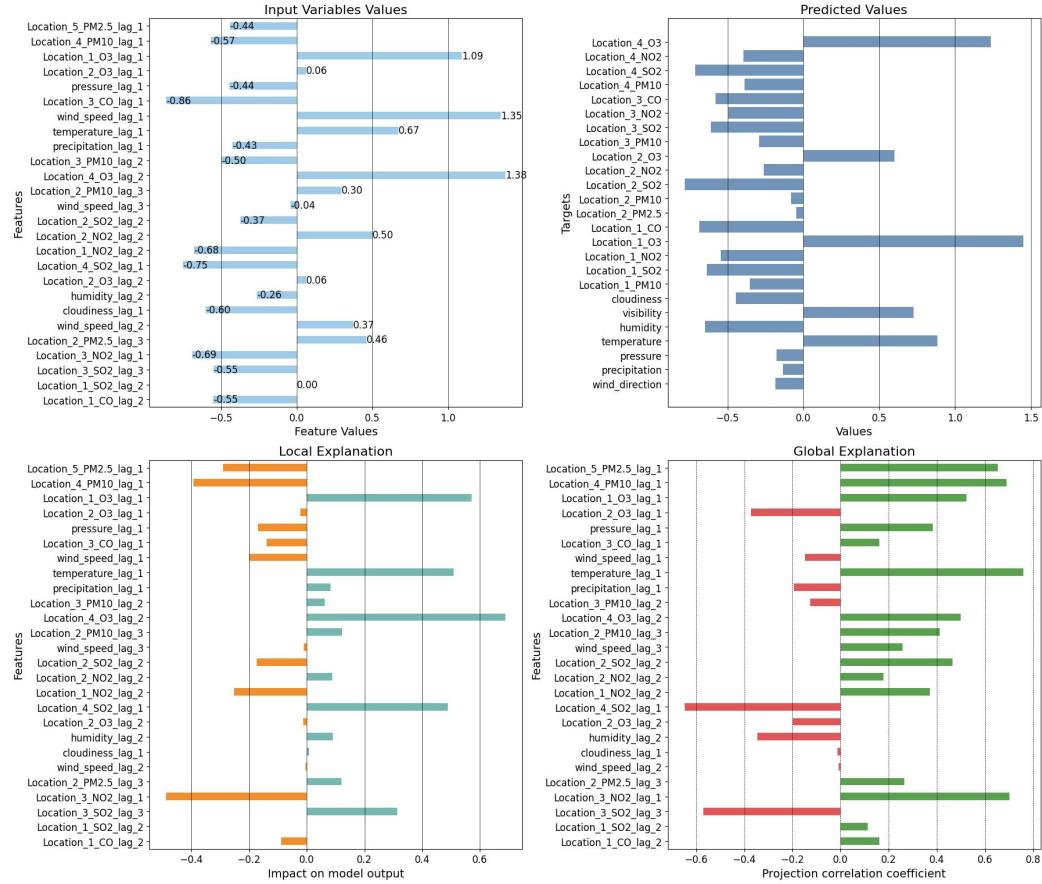
Virtual Metrology (VM) in Semiconductor Industry



Explanations based on Projective Operator for Multi-Output Models



Example of MIMO Explanations



eXProj explanation for one input sample:

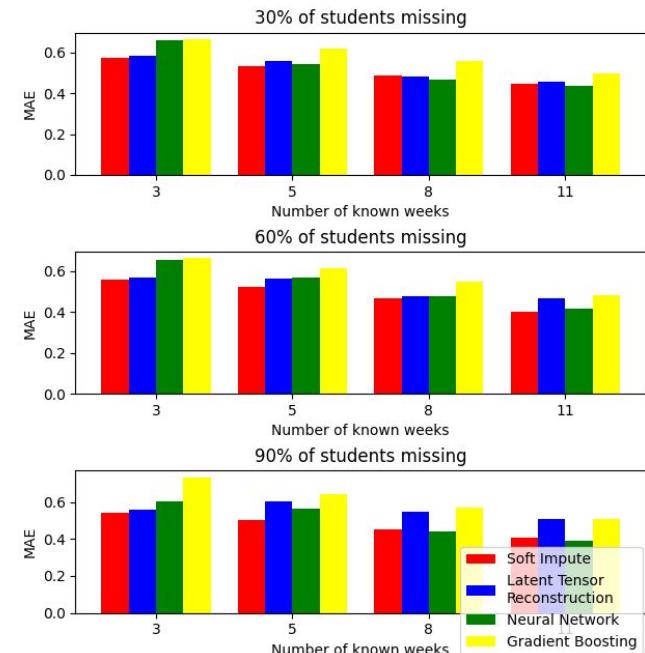
- Local explanation indicates the impact of the feature in on the prediction.
- Global explanation of the model indicate the influence of features on the predicted outputs.

AI-Assisted Learning

Objective

- Create AI systems that adapt to student's individual needs and knowledge state in order to maximize learning gains

Predicting mistakes that students make during learning to program using collaborative filtering techniques



Study on Investigating AI in Programming Education

Demographics

Age, Gender, Nationality, First time enrolled, Study department

Previous Experience

Programming proficiency in C, Programming proficiency, English proficiency, Frequency of AI use for work, Frequency of AI use for learning

Performance Scores

Total points, Grade

Weekly Usage and Feedback

Frequency of use for work, Frequency of use for learning, This week material comprehension, Lab difficulty, Last task difficulty

General Attitudes toward AI Scores

Positive attitude, Negative attitude

Weekly Motives Scores

Expectancy value, Attainment Value, Utility Value, Interest Value, Cost Value

Learning Style Scores

Visual, Auditory, Kinesthetic

Personality Traits Scores

Extraversion, Agreeableness, Conscientiousness, Negative Emotionality, Open-Mindedness

Thank you for your attention!



Email: senka.krivic@etf.unsa.ba