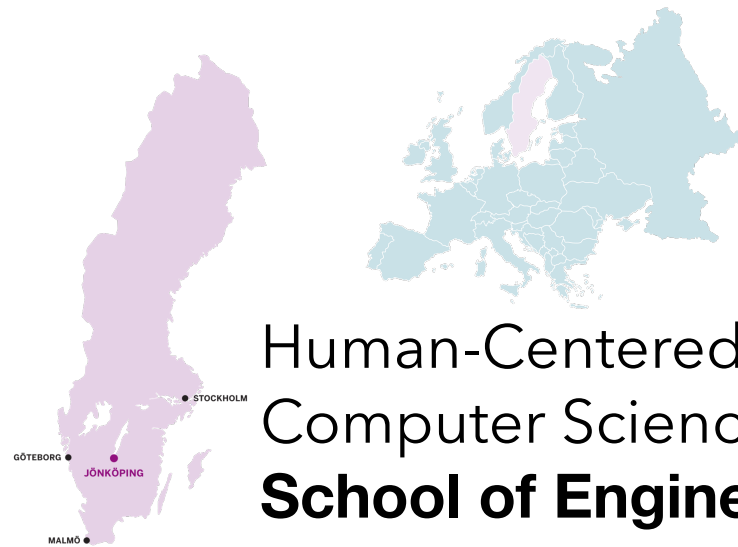
The image features a silhouette of a human figure on the left and a small, wheeled robot on the right. They are standing on a grassy horizon, holding hands. The background is a dramatic sky at sunset or sunrise, with the sun low on the horizon, casting a warm orange glow. The sky is filled with soft, white clouds. The overall mood is one of partnership and exploration.

Visualization-Empowered Human-in-the-Loop Artificial Intelligence

Explainable AI: what are explanations from AI-systems good for, anyway?

Maria Riveiro



Human-Centered Technology Computer Science and Informatics **School of Engineering**



Outline

- Human-Machine Collaboration.
Explainable AI
- Empirical Studies
- Future



Visualization-Empowered Human-in-the-Loop AI and Explainable AI

- **Visualization and Explainable AI provide transparency**

- Visualization as a powerful mediator
- Explanations (visual, textual, sound, multimodal)
encourage user engagement and facilitate informed
decision-making
- Support human-in-the-loop interaction
- Strengthen user trust and confidence in AI-generated
outcomes



FOCUS PERIOD LINKÖPING UNIVERSITY (CAMPUS NORRKÖPING) 2025

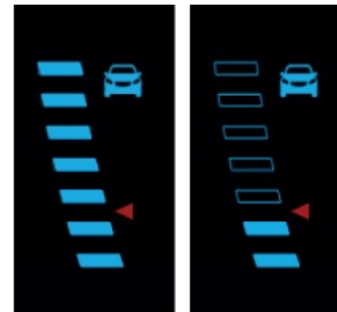
Visualization-Empowered Human-in-the-Loop Artificial Intelligence

APRIL 28 – MAY 30, 2025

In the past years, experts in visualization, human-computer interaction and related fields have substantially contributed to the topic, for instance, by the development of visualization approaches to open the typically closed black box design of popular machine learning methods. However, the rapid developments in AI/ML potentially trigger a fundamental change in our understanding of the capabilities and applicability of the models as they are now also able to “interact” with the general population. What are the implications in terms of trust into the analytical results and potential biases that may occur? How should visualization research react and adapt to increase trust and call our attention to critical biases to avoid them?



Autonomous driving



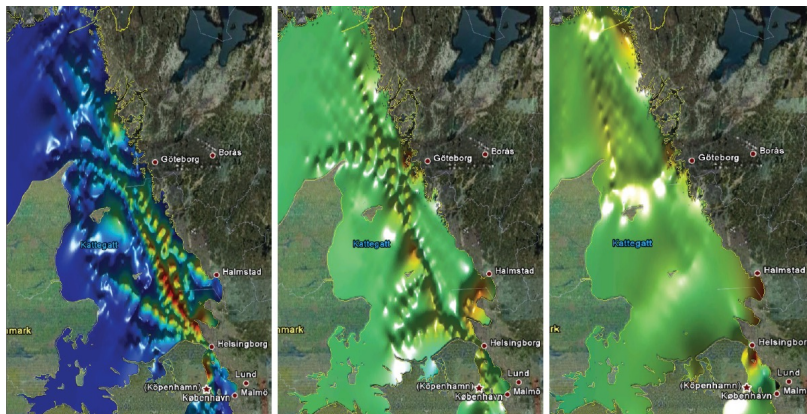
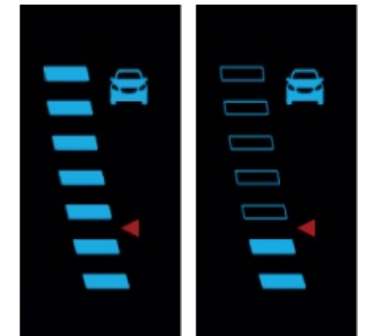
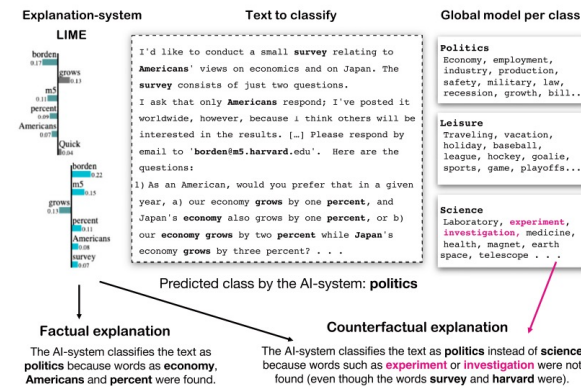
Results:

- Faster take-over
- More look away
- Better trust calibration

Helldin, T., Falkman, G., Riveiro, M., & Davidsson, S. (2013, October). Presenting system uncertainty in automotive UIs for supporting trust calibration in autonomous driving. In *Proceedings of the 5th international conference on automotive user interfaces and interactive vehicular applications* (pp. 210-217).

EXPLAINABLE AI & VISUALIZATION

- Explainable methods
- Visualization
- Human-Centered AI





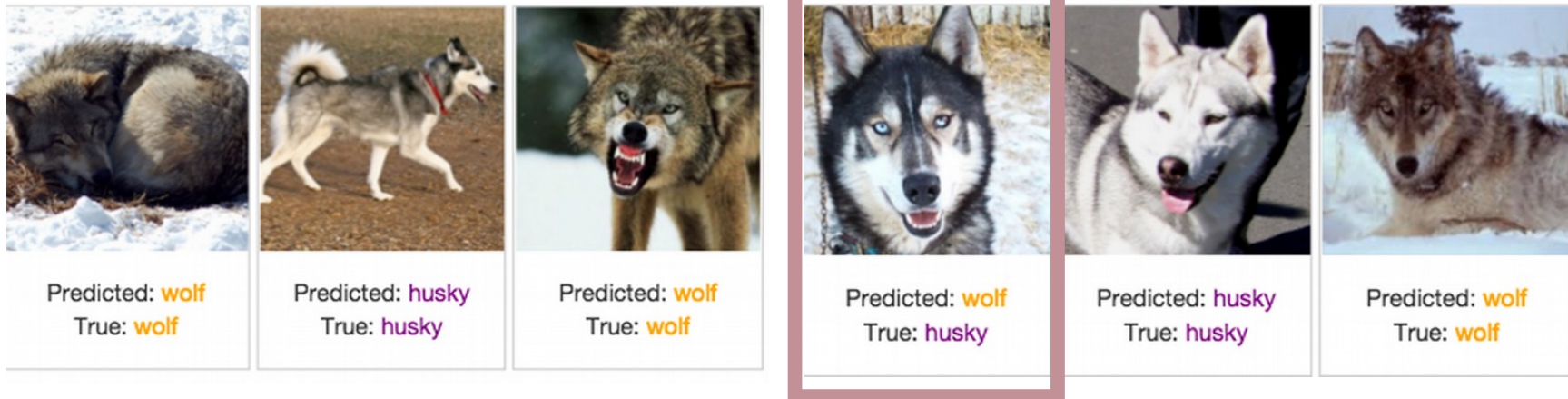
Human-machine collaboration

1. AI-systems need to support humans in understanding them
2. AI-systems need to be able to understand humans

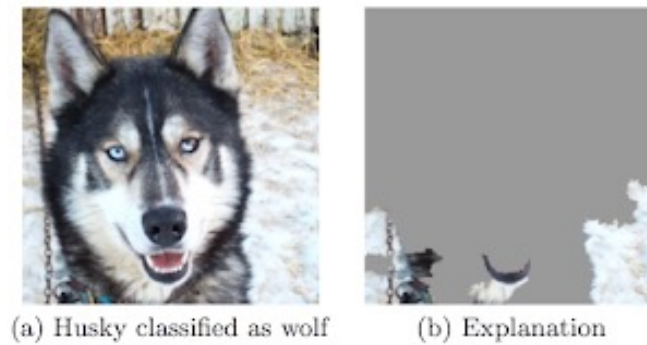
Transparency and explainability

- AI systems often operate as 'black boxes', lacking transparency.
- This makes it difficult for users to understand and trust the system.
- Explainable AI (XAI) seeks to provide clarity and justification for AI decisions.





Explainable AI



Ribeiro, M. T., Singh, S., & Guestrin, C. (2016, August). "Why should i trust you?" Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining* (pp. 1135-1144).

AI used in many application areas (automation - human support, LoA)



Empirical studies

The brains...

People train robots
by demonstration,
and robots train
other robots.
Reinforcement
learning



AI – human
communication
decision-making



What is my
neural network
actually learning?
Visualising
concepts in NN



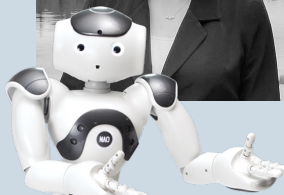
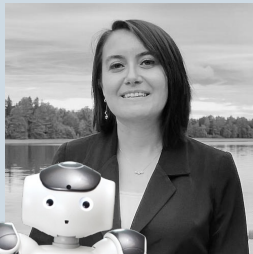
Mental models,
theory of mind,
expectations



Human Robot
Interaction

Expectations

PostDoc



Detection of
defects in
manufacturing



Automatic product
recognition.

Computer Vision,
fine-grained
techniques and
data generation



ML for
forecasting
and planning

AI adoption



What are explanations from AI-systems good for, anyway?

- Explanations lead to **positive results** (better understanding, trust, higher confidence in own decisions, satisfaction, performance, better mental models) but also....

negative effects or trade-offs

- ... revealing limitations led to negative heuristics, under reliance
- ... persuasion (follow advice even if it is incorrect, advice-taking) and overreliance
- ... unnecessary explanations lead to higher cognitive load, information overload, more time
- ... confusion
- ... perceived accuracy is more important than explainability, no explanations needed

Trust. The relationship between explainability and trust is difficult to comprehend... (trust calibration)

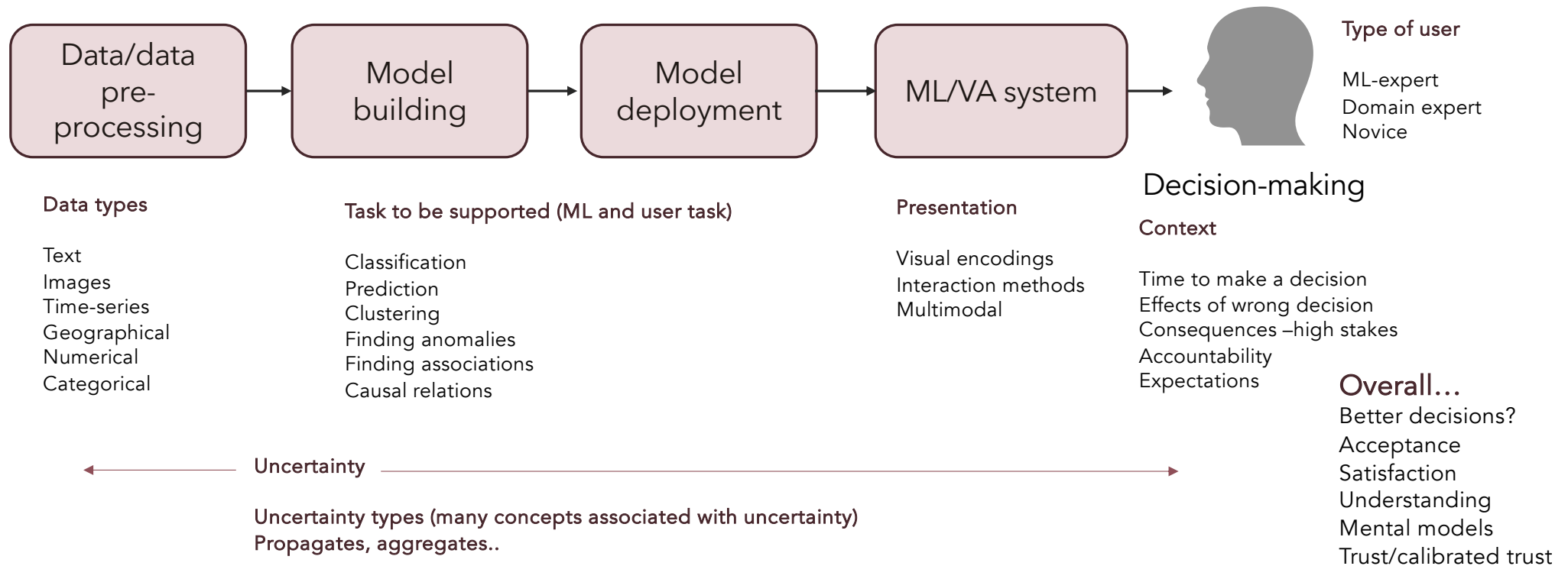
Explaining something to someone is a **complex cognitive process**... (not only XAI-ML)

XAI – a solution looking for a problem?

Expert systems (80s)

XAI design space is complex ...

WHY, WHAT, WHEN, WHERE & HOW?



Empirical studies

Expectations

- Riveiro, M., and Thill, S. (2021). "That's (not) the output I expected!" On the role of end user expectations in creating explanations of AI systems. Artificial Intelligence, Elsevier, ISSN 0004-3702, E-ISSN 1872-7921, Vol. 298, article id 103507
- Riveiro, M., and Thill, S. (2022). The challenges of providing explanations of AI systems when they do not behave like users expect. New York: Association for Computing Machinery (ACM), UMAP '22: 30th ACM Conference on User Modeling, Adaptation and Personalization, Barcelona, Spain, July 4-7, 2022

Task difficulty

- Ingesson, E. and Riveiro, M (2025 to be submitted) When Do We (Not) Want Explanations? A Study on Explanation Demand in Human-AI Decision-Making.

Human robot interaction - explaining errors

- Akalin, N. and Riveiro, M (2025 ROMAN). Let Me Explain Why I didn't Take the Action You Wanted! Comparing Different Modalities for Explanations in Human-Robot Interaction

Chatbot Alba mental health support

- Holmberg, L., Sikström, S. and Riveiro, M. Speaking or Writing? Do Response Times Influence Anthropomorphism Differently for ADHD and Neurotypical Users in a Mental Health Chatbot? (2025 under review) Conversational User Interfaces.

How empirical studies inform theory

- Riveiro, M and Thill, S. (2025 under review) The diversity of empirical research on explainable artificial intelligence and implications for theory building. ACM Transactions on Interactive Intelligent Systems.

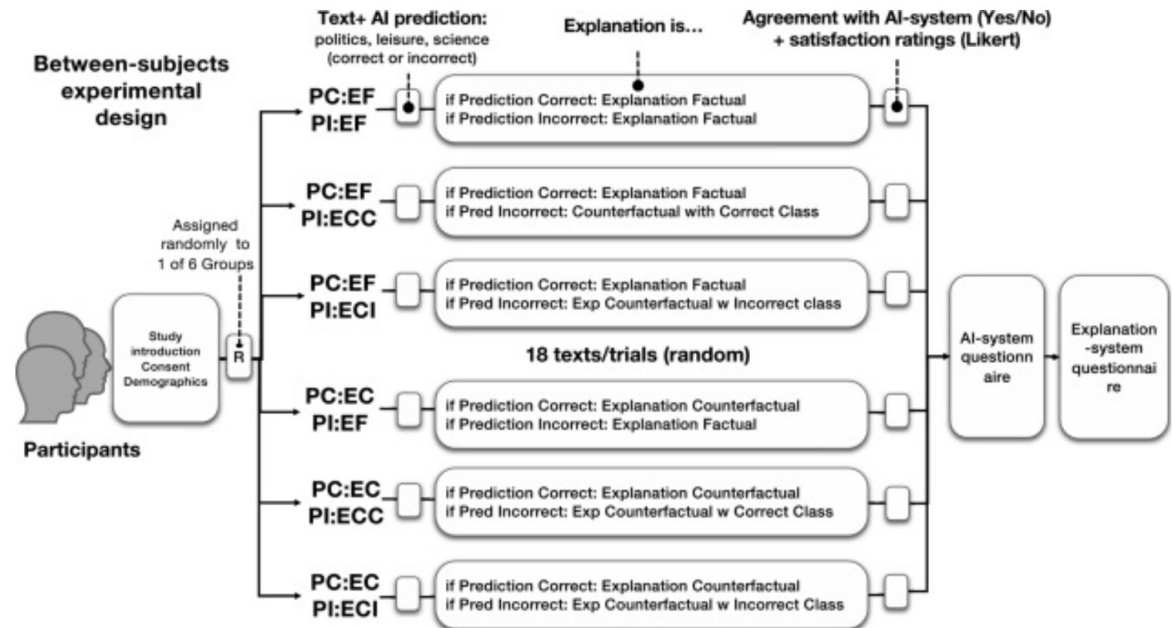


Expectations

- In human-human interactions, explanations are often needed when an event is unexpected (Why?), and we need to explain the unexpected fact in relation to an implicit expected foil
- Do expectations play a role in when and what? Do they modulate the content of explanations? If we don't consider them, do we risk that you are not getting the explanation you are looking for?
- So.... what do we want to see in the explanations when we don't agree with the system/when it does something that we don't expect?

Measures/metrics

- System understanding
- Explanation satisfaction, completeness
- Performance
- Perceived need for interaction

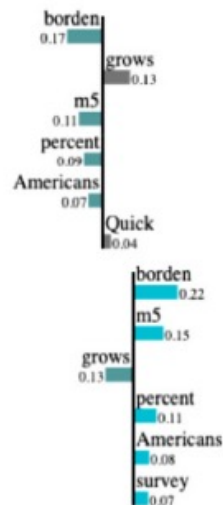


Example. Prediction Correct (PC): politics. Prediction Incorrect (PI): science or leisure.

- Explanation Factual (EF): The AI-system classifies the text as *politics* because words as *economy*, *Americans* and *percent* were found.
- Explanation Counterfactual (EC): The AI-system classifies the text as *politics* instead of *science* because words such as *experiment* or *investigation* were not found (even though the words *survey* and *harvard* were).
- Explanation Counterfactual with Correct Class (ECC): The AI-system classifies the text as *science/leisure* instead of *politics* because words such as *financial* or *growth* were not found (even though the words *American* and *percent* were).
- Explanation Counterfactual with Incorrect Class (ECI): The AI-system classifies the text as *leisure* instead of *science* because words such as *experiment* or *investigation* were not found (even though the words *survey* and *harvard* were).

Explanation-system

LIME



Text to classify

I'd like to conduct a small **survey** relating to **Americans'** views on economics and on Japan. The **survey** consists of just two questions. I ask that only **Americans** respond; I've posted it worldwide, however, because I think others will be interested in the results. [...] Please respond by email to '**borden@m5.harvard.edu**'. Here are the questions:

1) As an American, would you prefer that in a given year, a) our economy **grows** by one **percent**, and Japan's **economy** also grows by one **percent**, or b) our **economy grows** by two **percent** while **Japan's** economy **grows** by three percent? . . .

Global model per class

Politics

Economy, employment, industry, production, safety, military, law, recession, growth, bill..

Leisure

Traveling, vacation, holiday, baseball, league, hockey, goalie, sports, game, playoffs...

Science

Laboratory, **experiment**, **investigation**, medicine, health, magnet, earth space, telescope . . .

Predicted class by the AI-system: **politics**

Factual explanation

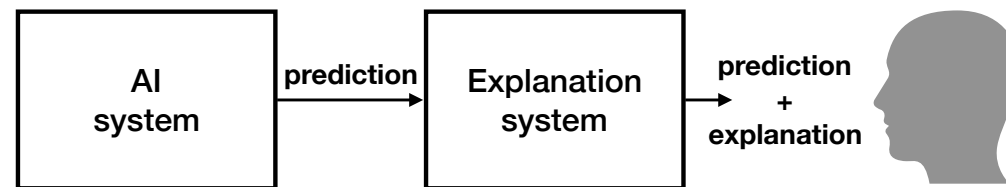
The AI-system classifies the text as **politics** because words as **economy**, **Americans** and **percent** were found.

Counterfactual explanation

The AI-system classifies the text as **politics** instead of **science** because words such as **experiment** or **investigation** were not found (even though the words **survey** and **harvard** were).

Role of expectations in explanations

- Do expectations determine explanation content?
- Are counterfactuals preferred when outcomes from AI-system are unexpected?



- Factual and counterfactual explanations
 - ✓ • H1: Factual explanations are appropriate for correct predictions because the system output is in line with the expected output.
 - ✗ • H2: Counterfactual explanations that contain the expected foil are appropriate when the system prediction is incorrect

So.... what do we want to see in the explanations
when the system does something that we don't
expect?

Method

STUDY I
(multiple-choice)

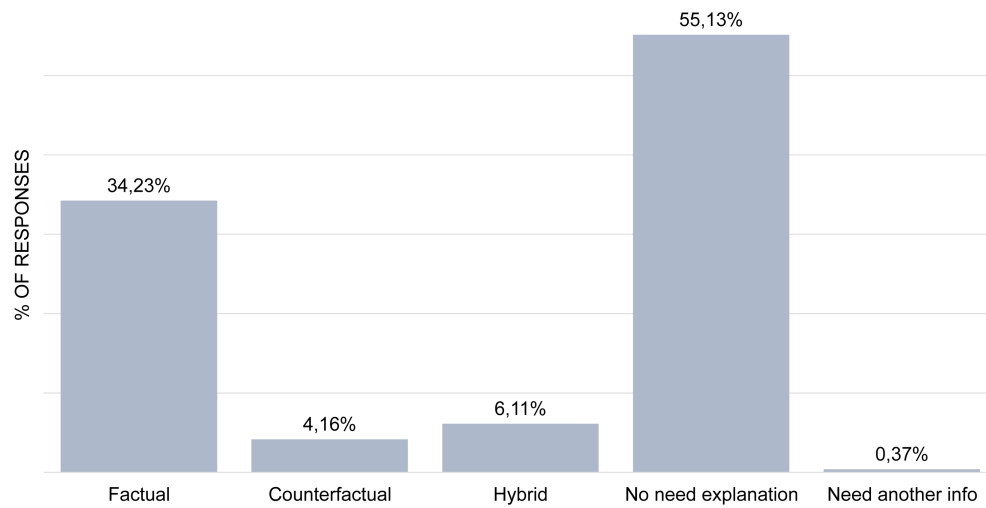
STUDY II
(open questions)

- We presented participants with various scenarios involving a text classifier and then asked them to indicate their preferred explanation for each scenario
- One group of participants chose the type of explanation from a multiple-choice questionnaire (Study I), the other had to answer using free text (Study II)

162
15 (+ 2)

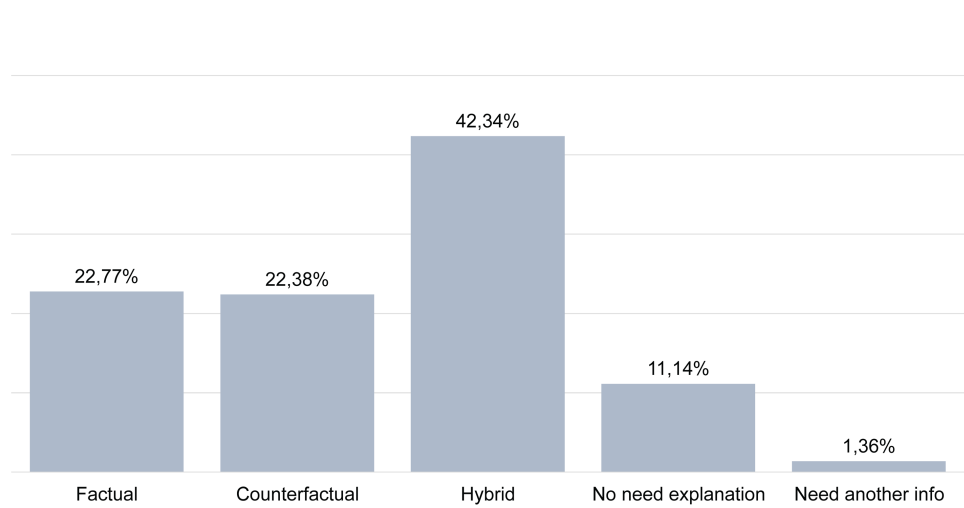
STUDY I

Preferred explanations when expectations are matched (correct output) and consistent behavior - 162 participants



Matched expectations

Preferred explanations when expectations do not match (incorrect output) and consistent behavior - 162 participants



Mismatched expectations

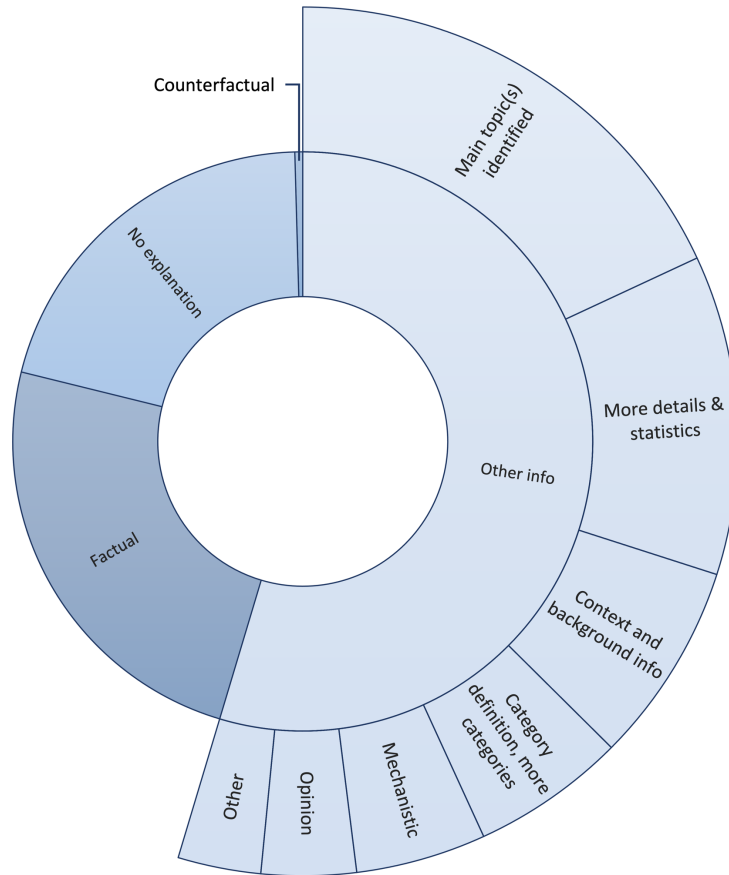
Table 2. Main categories found under class “other information” in response to “What kind of information you would like to see in the explanations from the AI system in order to understand this prediction?” when the AI system output matched user expectations. n=126.

Category	Description	Count	Example evidence from responses
Main topic(s) identified	A summary of the main themes, overall topic found in the text	41	<i>“Identify theme and central ideas”.</i> <i>“About UN to establish a police force for Haiti”.</i> <i>“Politics, colonial, occupation”.</i>
More details and statistics	More details, more evidence, statistics	27	<i>“More information in regards to the smaller details of the text”.</i> <i>“I’d like to know the ingredients of the medicine.</i> <i>Perhaps their specific function.”</i> <i>“Statistics.”</i> <i>“Also, the percentage of key words that were mapped to the category (assuming that the AI is running on some sort of key word matching logic).”</i>
Context and background info	Context, background info on article, where this text comes from, reference	17	<i>“The source of the article and the names of the people having the conversation.”</i> <i>“Research and references, something that can actually prove this is the case.”</i> <i>“It seems like it’s missing context.”</i>
Category definition, more categories	Unclear how the AI system defined each category, more categories needed	13	<i>“Explain how Sports is Leisure.”</i> <i>“Sub topic of politics.”</i> <i>“What kind of politics.”</i>
Mechanistic	How the AI system reached that conclusion, reasoning	11	<i>“How it came to this conclusion.”</i> <i>“Showing the reasoning behind the decisions if able to do so.”</i>
Opinion about the text	Participants expressed their opinions about topics in text	8	<i>“It should be said how the UN predicts policies.”</i> <i>“I would love to see the breakdown of how taxes are being spent .”</i>
Other/Not relevant	Comments not related to explanations	7	<i>“Don’t understand the text .”</i>

Table 3. Main categories found under class “other information” in response to “What kind of information you would like to see in the explanations from the AI system in order to understand this prediction?” when the output from the AI system did not match user expectations. n=198.

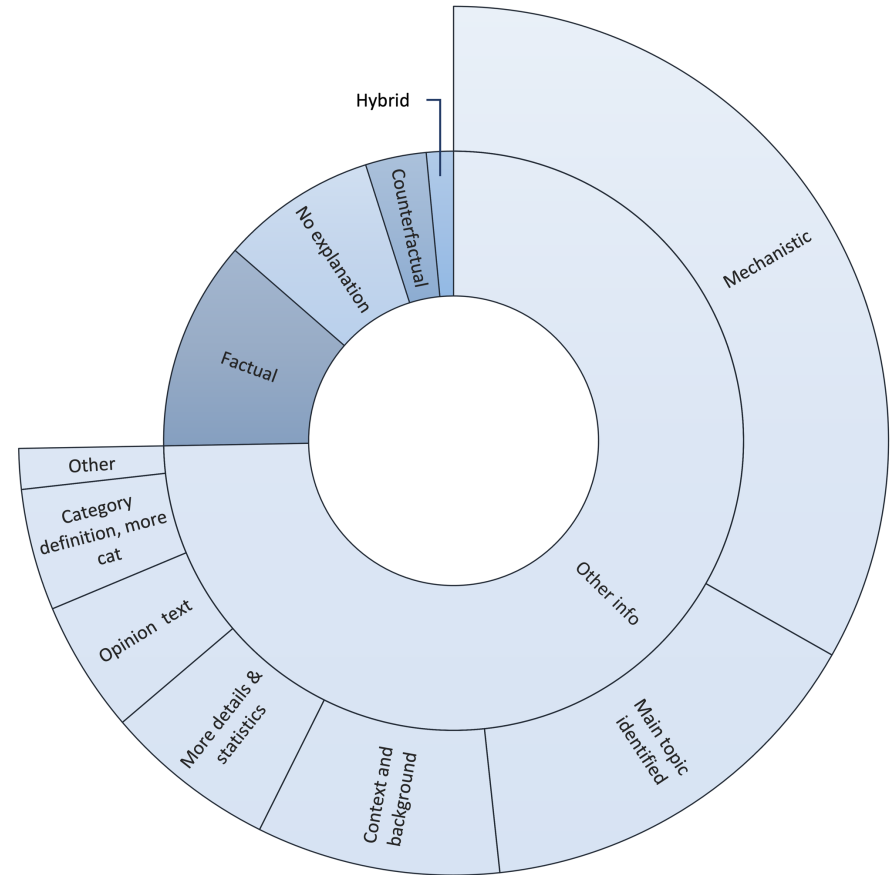
Category	Description	Count	Example evidence from responses
Mechanistic	How the AI system reached that conclusion, reasoning	88	<i>“I don’t understand what the AI system would think this is leisure. How exactly would a gun buyback program be of any leisure ?”</i> <i>“What got you to believe this was politics?”</i> <i>“I would like the AI system to list reason for suggesting the prediction as Politics.”</i>
Main topic(s) identified	A summary of the main themes, overall topic found in the text	40	<i>“The effects of gun buyback program and how it affects citizens.”</i> <i>“It’s about been tour round the world.”</i>
Context and background info	Context, background info on article, where this text comes from, reference	24	<i>“I would love to know information source about the fail-safe mechanism.”</i> <i>“I would like to see some reasons that justifies the prediction like the context in which the discussion was made.”</i> <i>“How the ecosystem in Utah works, and the climate of Utah.”</i>
More details and statistics	More evidence and statistics	17	<i>“I would love to know more about NOOP operation.”</i> <i>“More clearer information such as locations.”</i> <i>“Statistics.”</i>
Opinion about the text	Participants expressed their opinions about topics in text	13	<i>“Gun sport maybe leisure to some people but this is his opinion more than anything else.”</i>
Category definition, more categories	Unclear how the AI system defined each category, more categories needed	12	<i>“Language selected to recognise leisure, what is ‘leisure’ by definition.”</i> <i>“What leisure activity is referred to.”</i>
Other/Not relevant	Comments not related to explanations	4	<i>“Better paragraph structure.”</i>

Content of explanations when expectations are matched



Matched expectations

Content of explanations when expectations are **not** matched



Mismatched expectations

Conclusions from expectations

- For matched expectations, an explanation is often not required at all, while if one is, it is of the factual type
- Providing explanations when system output does not match user expectations is a challenging matter, primarily because there does not seem to be a unique strategy, although mechanistic explanations are requested more often than other types
- No one size fits all
- Overall, user expectations are a significant variable in determining the most suitable content of explanations (including whether an explanation is needed at all)

Akalin, N. and Riveiro, M (2025 ROMAN). **Let Me Explain Why I didn't Take the Action You Wanted!** Comparing Different Modalities for Explanations in Human-Robot Interaction

- Aim: Explore preferred modalities for robot explanations when robots decline user requests. We focus on explanations in scenarios where a user makes a request to the robot, but the robot does not perform the requested action for various reasons
- Method: User study assessing various explanation modalities (visual, auditory, gesture).

Results:

- Participants strongly preferred speech-based explanations for clarity, naturalness, and ease of understanding.
- Multimodal explanations (combining speech with lights, sounds, or gestures) were preferred for critical or urgent situations to ensure attention and clarity.
- User preferences for explanation modalities varied according to context, such as urgency, ambient noise, or user's distance from the robot.
- Overall, users found explanations clear and helpful, but desired improvements for more human-like or customizable voices.
- Users requested more interactive explanations, such as allowing follow-up questions.



Holmberg, L., Sikström, S. and Riveiro, M. Speaking or Writing? **Do Response Times Influence Anthropomorphism Differently for ADHD and Neurotypical Users in a Mental Health Chatbot?** (2025 under review) Conversational User Interfaces.

- Aim: Investigate how interaction modality (text vs. speech) affects anthropomorphism of chatbots among neurotypical and ADHD users.
- Method: User study with 107 participants interacting via text or speech.
- Results: Neurotypical users equally anthropomorphize chatbots in both conditions, while ADHD users anthropomorphize less in voice due to longer response times.

Chatbot Alba mental health support

AI CBT: Cognitive behavioral therapy with an AI-therapist

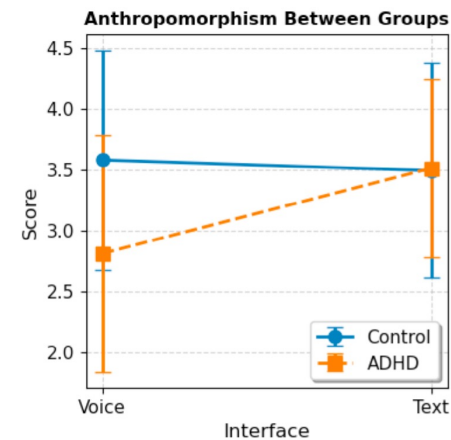


Welcome! Today I would like to talk to you about your mental health. I'd like to ask you to share what brings you in today or any primary concerns you've been experiencing lately. Please feel free to describe your symptoms or any distressing experiences in as much detail as you're comfortable with.

I have difficulties concentraring and do not get anything done

I'm sorry to hear that you're experiencing these difficulties. Can you tell me more about when you first noticed these concentration issues and how they've been affecting your daily life?

It has always been like this



User settings

Please define how you want to use the chatbot.

☒ Voice chat

Please use the button in the 'Voice recorder' section when you start and stop to speak

Choose a voice style:

allory

Voice recorder

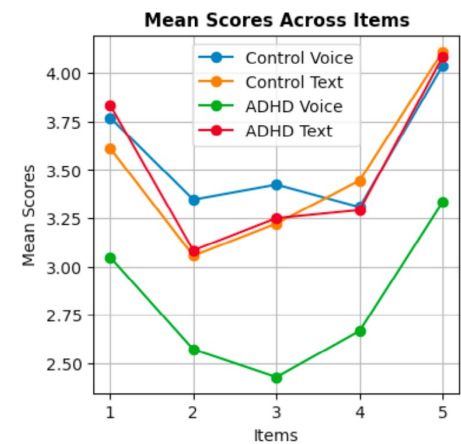
Start Recording

Welcome to an interview about your mental health!

Before starting the interview, you must provide informed consent through this form: <https://tinyurl.com/mr7abde>

Have you given your informed consent? Please answer yes or no.

☐ Waiting for the assistant's answer...

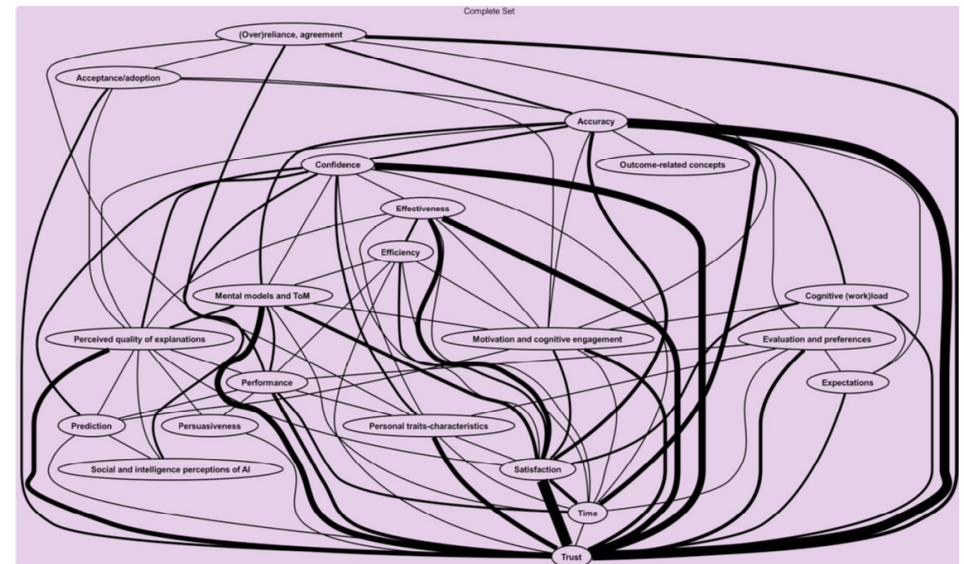
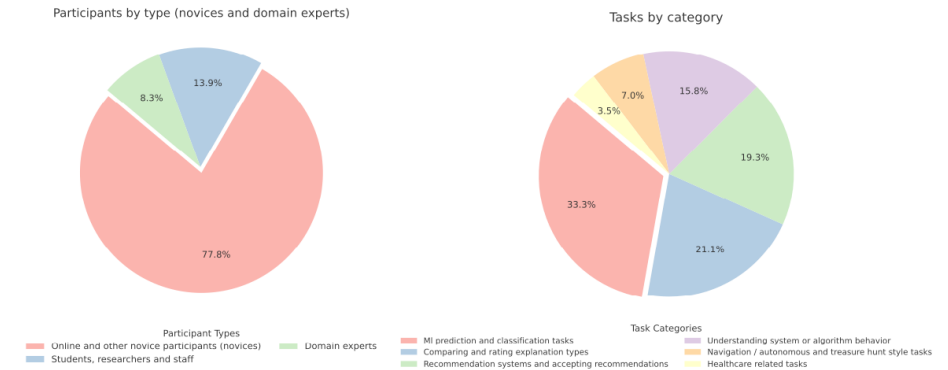


Riveiro, M and Thill, S. (2025 under review) **The diversity of empirical research on explainable artificial intelligence and implications for theory building.** ACM Transactions on Interactive Intelligent Systems.

- Aim: Overview and summary of empirical studies in XAI, analyzing contexts, purposes, and effects of explanations.
- Method: Reviewed 95 empirical studies evaluating explanations in human-AI interaction.

Results:

- Lack of ecological valid experiments.
- XAI research is highly diverse, lacking a common theoretical framework.
- Explanations serve multiple purposes (trust, understanding, decision support) but show varied effectiveness.
- ... principles for deriving a general theory of explanations from AI-systems?

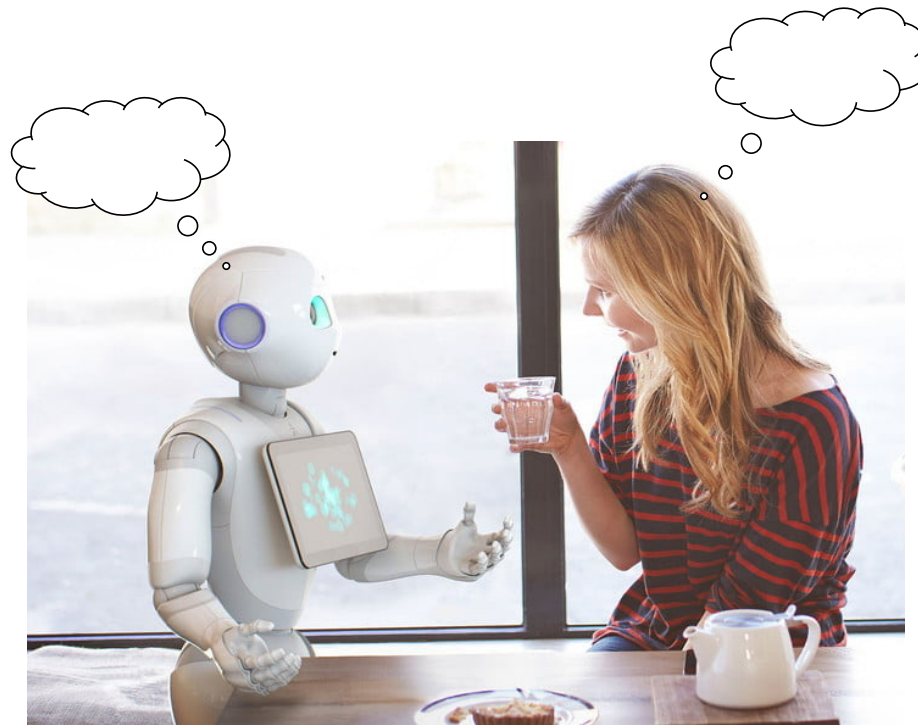




Human-machine collaboration

1. AI-systems need to support humans in understanding them
2. AI-systems need to be able to understand humans

AI-systems need to be able to understand humans



Adaptation

- Human-AI collaboration (my stand is that it mirrors human-human interactions)
- Understand users (needs, expectations, abilities, personality traits) and adapt interactions accordingly

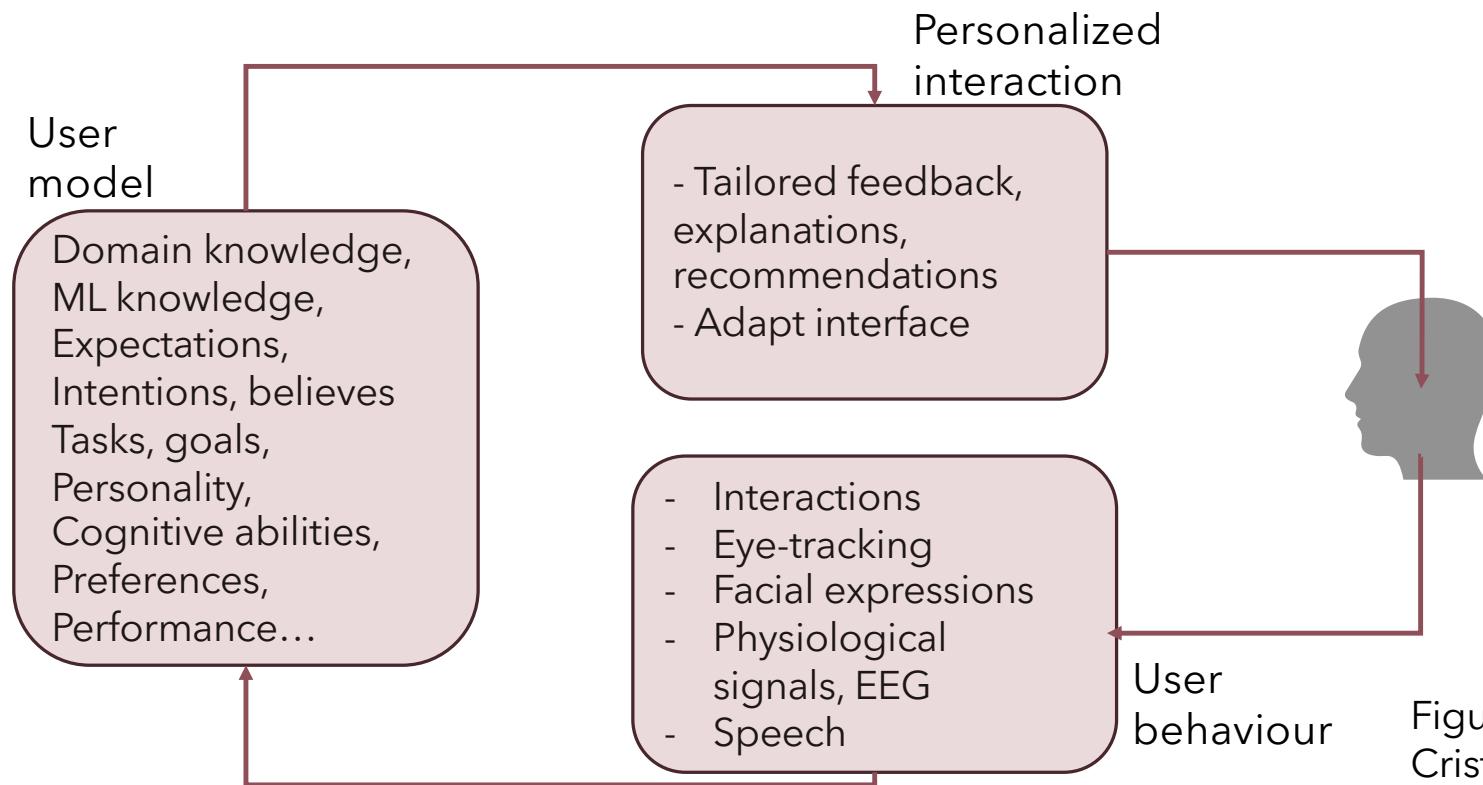
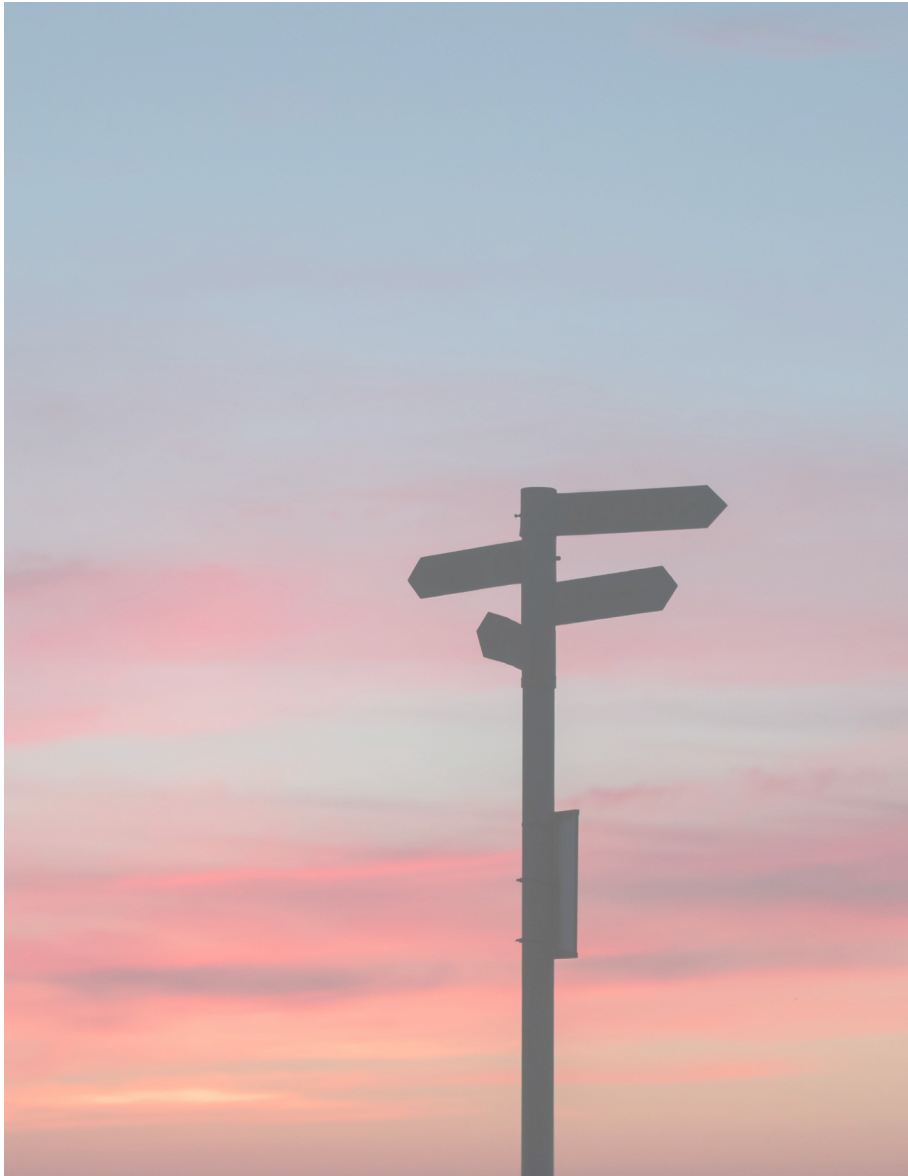


Figure inspired by
Cristina Conati's work



Future

- User models, mental models
- Expectations
- Theory of Mind
- Intention recognition
- Visualization, presentation and interaction modalities
- Adaptation: engagement, curiosity, knowledge
- Use cases!

Thanks for listening!



JÖNKÖPING UNIVERSITY

Human-Centered Technology (HCT)



HUMAN-COMPUTER COMMUNICATION



EDUCATION & HEALTHCARE



SMART INDUSTRIES



INFORMATION VISUALIZATION