

# Revisiting **Intelligence Augmentation**: Investigating and Mitigating the Risks of AI to Human Intelligence

**Q. Vera Liao**

Computer Science and Engineering  
University of Michigan



**M** UNIVERSITY OF  
MICHIGAN

# AI v.s. IA (Intelligence Augmentation)

*“Machines will be capable, within twenty years, of doing any work that a man can do”*

— Herbert Simon, 1960

# AI v.s. IA (Intelligence Augmentation)

*“Machines will be capable, within twenty years, of doing any work that a man can do”*

— Herbert Simon, 1960

This report covers the first phase of a program **aimed at developing means to augment the human intellect**. These "means" can include many things—all of which appear to be but extensions of means developed and used in the past to help man apply his native sensory, mental, and motor capabilities—and we consider the whole system of a human and his augmentation means as a proper field of search for practical possibilities. It is a very important system to our society, and like most systems its performance can best be improved by considering the whole as a set of interacting components rather than by considering the components in isolation.

*Augmenting Human Intellect: A Conceptual Framework*

—Douglas Engelbart, 1962

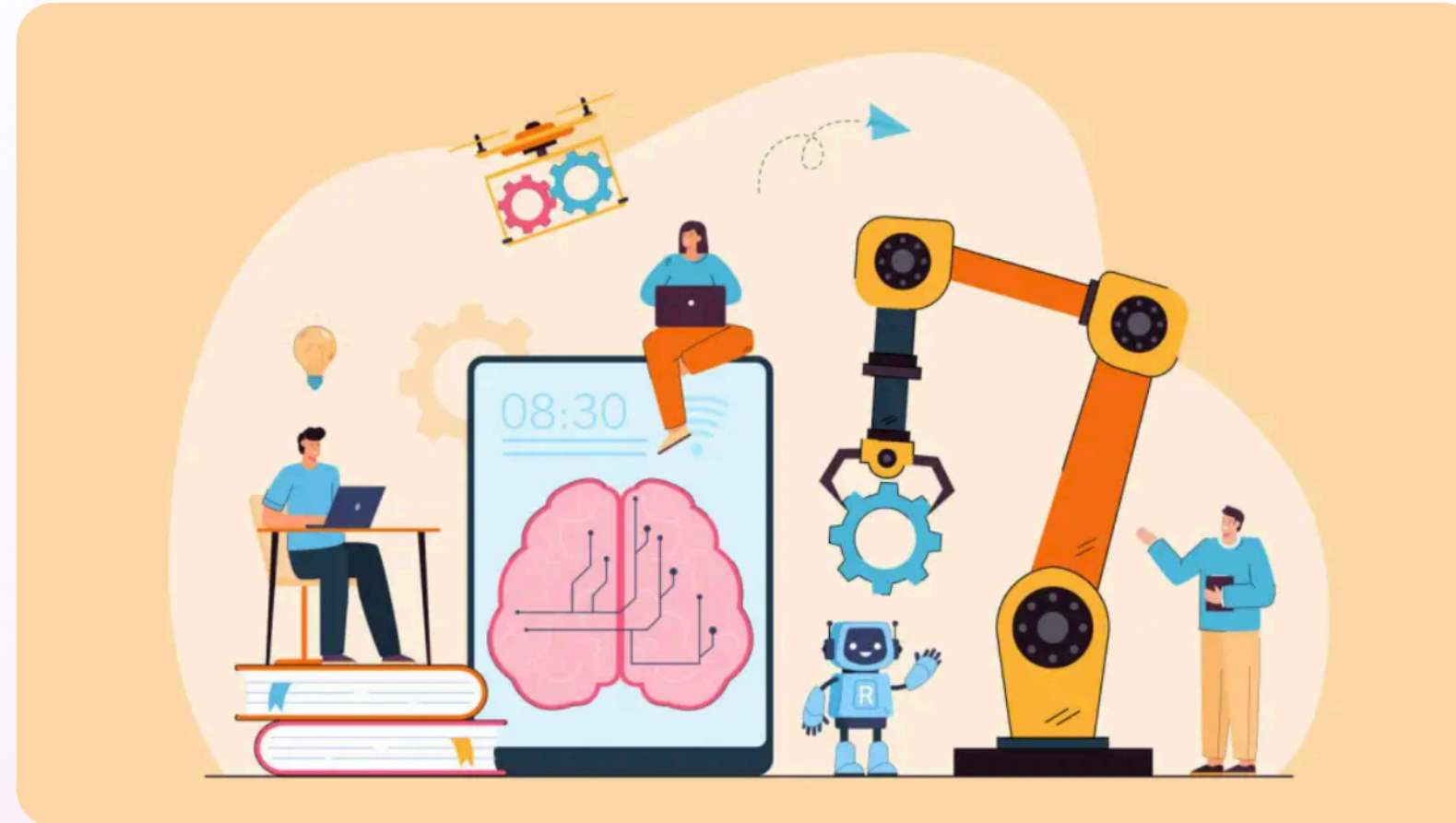
**AI augmenting, not replacing, people**

# AI as Assistant/Copilot/Collaborator...

## Human-AI Collaboration: How to Work Together

Prepare your digital workforce to thrive in an AI-driven world.

[Choi Chow, Product Marketing Manager](#)

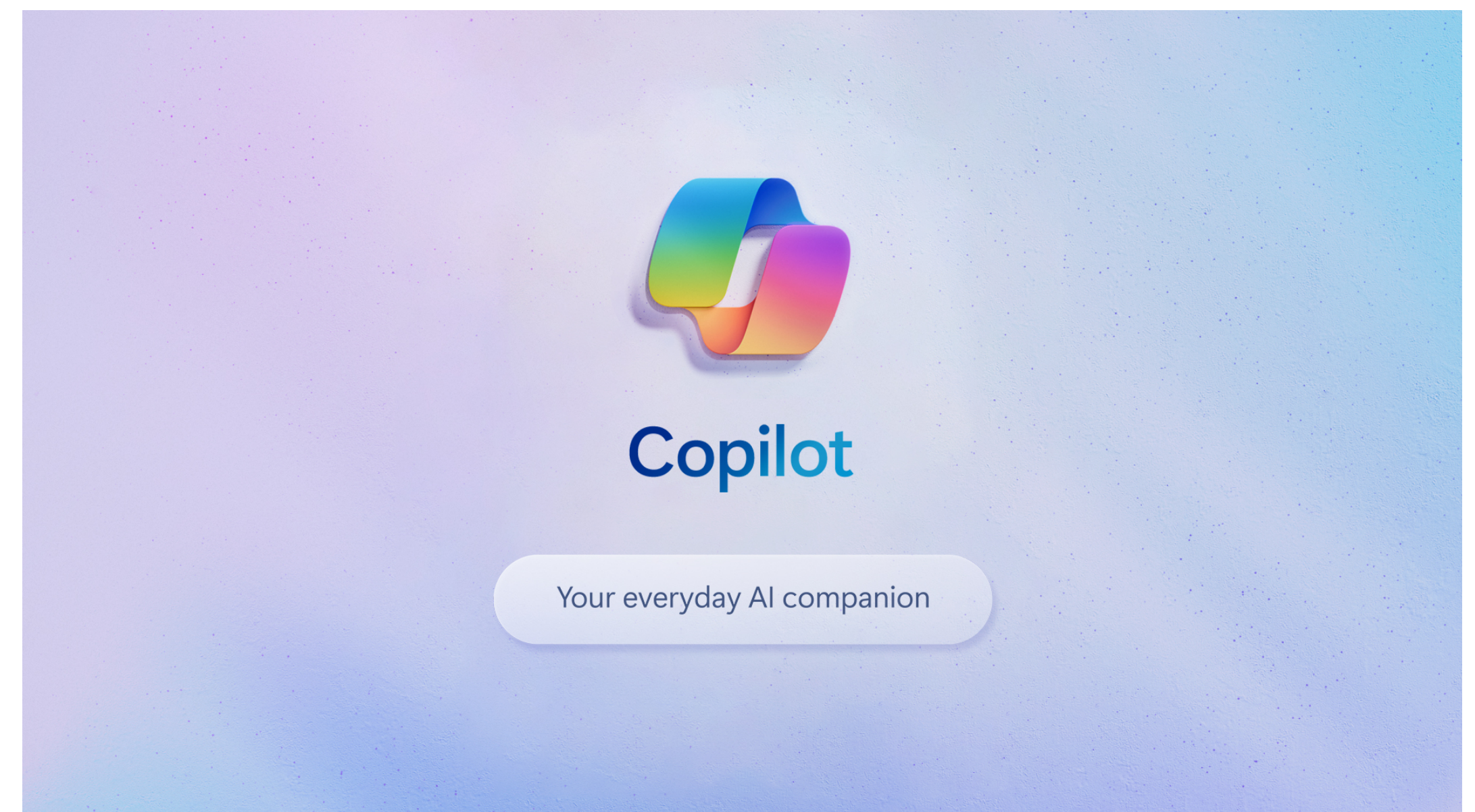
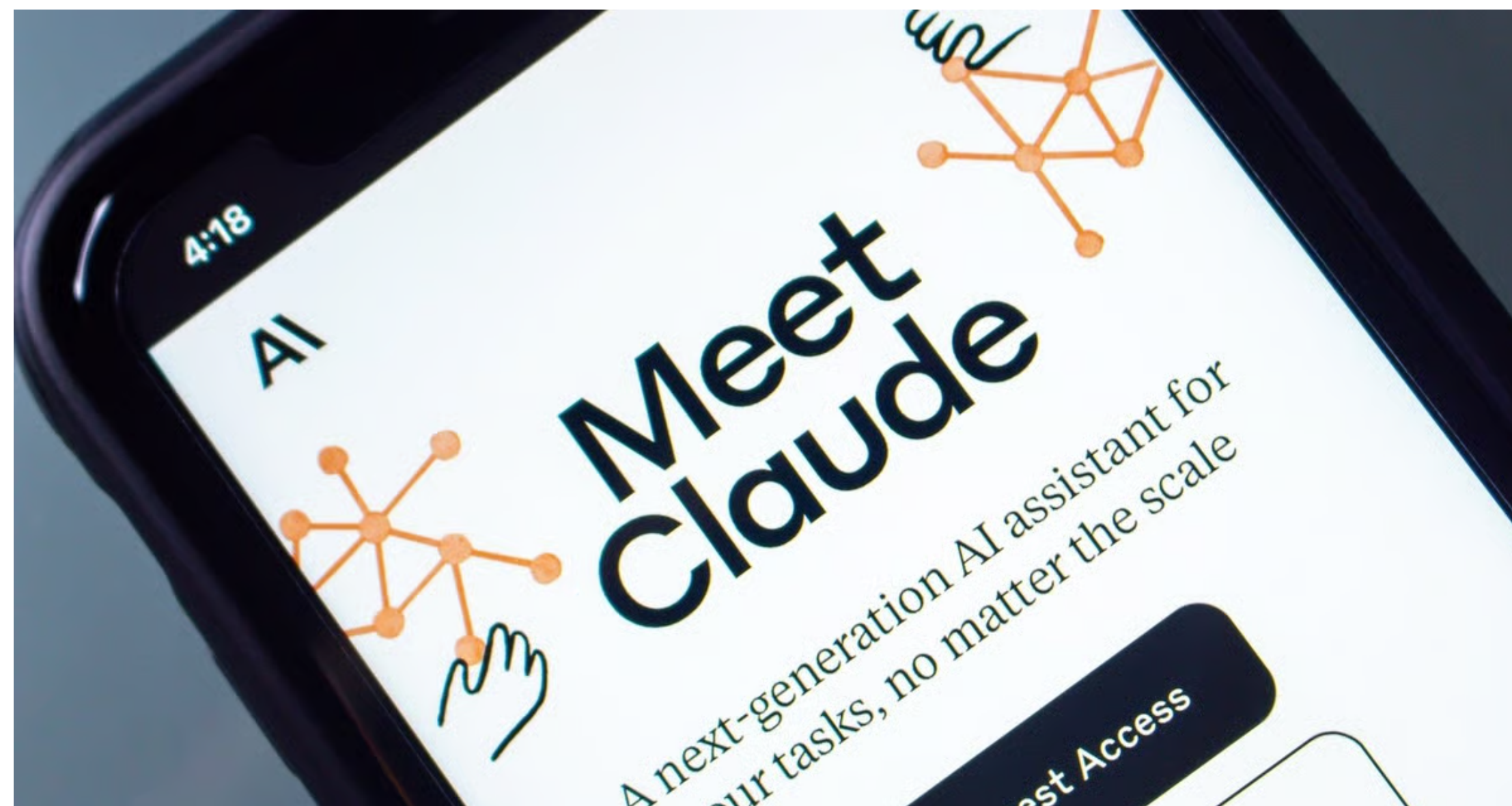


## Go from answers to action with Claude Cowork

Cowork brings Claude Code's agentic capabilities to your desktop. Give Claude access to your local files, set a task, and step away. Come back to completed work. Now in research preview.

[Download Claude](#)

[Talk to sales](#)



# Intelligence Augmentation in Education

For teachers

Experience the best AI for teachers.

Built for educators by educators, with your needs in mind. Khanmigo simplifies your workflow while keeping your work and your student data private and secure.

## Safe, teacher-guided AI for students

Brisk's student-facing AI tools empower educators to give every student personalized, interactive support right where they're learning — all with teachers in the driver's seat.

For Teachers

An AI Toolkit for the Work Teachers Do Every Day.

## Enhance teaching with research-backed tools

Built on K-12 research and best practices, Solara includes ready-made prompts and tools to create lesson plans, rubrics, grade-level texts, and more—saving time for meaningful student engagement. Easily configure the tools to fit your district's needs.

AI FOR TEACHERS

Save time.  
Spark creativity.  
Personalize learning.



Teachers at the center

Our human-in-the-loop approach to AI ensures that teachers always have the final say, keeping educator expertise at the heart of every decision.



OFFICE OF  
Educational Technology

# Artificial Intelligence and the Future of Teaching and Learning

Insights and Recommendations

May 2023



## Perspective: Intelligence Augmentation

*"Augmented intelligence is a design pattern for a human-centered partnership model of people and artificial intelligence (AI) working together to enhance cognitive performance, including learning, decision making, and new experiences." <sup>16</sup>*

Foundation #1 (above) keeps humans in the loop and positions AI systems and tools to support human reasoning. "Intelligence Augmentation" (IA)<sup>17</sup> centers "intelligence" and "decision making" in humans but recognizes that people sometimes are overburdened and benefit from assistive tools. AI may help teachers make better decisions because computers notice patterns that teachers can miss. For example, when a teacher and student agree that the student needs reminders, an AI system may provide reminders in whatever form a student likes without adding to the teacher's workload. Intelligence Automation (IA) uses the same basic capabilities of AI, employing associations in data to notice patterns, and, through automation, takes actions based on those patterns. However, IA squarely focuses on helping people in human activities of teaching and learning, whereas AI tends to focus attention on what computers can do.

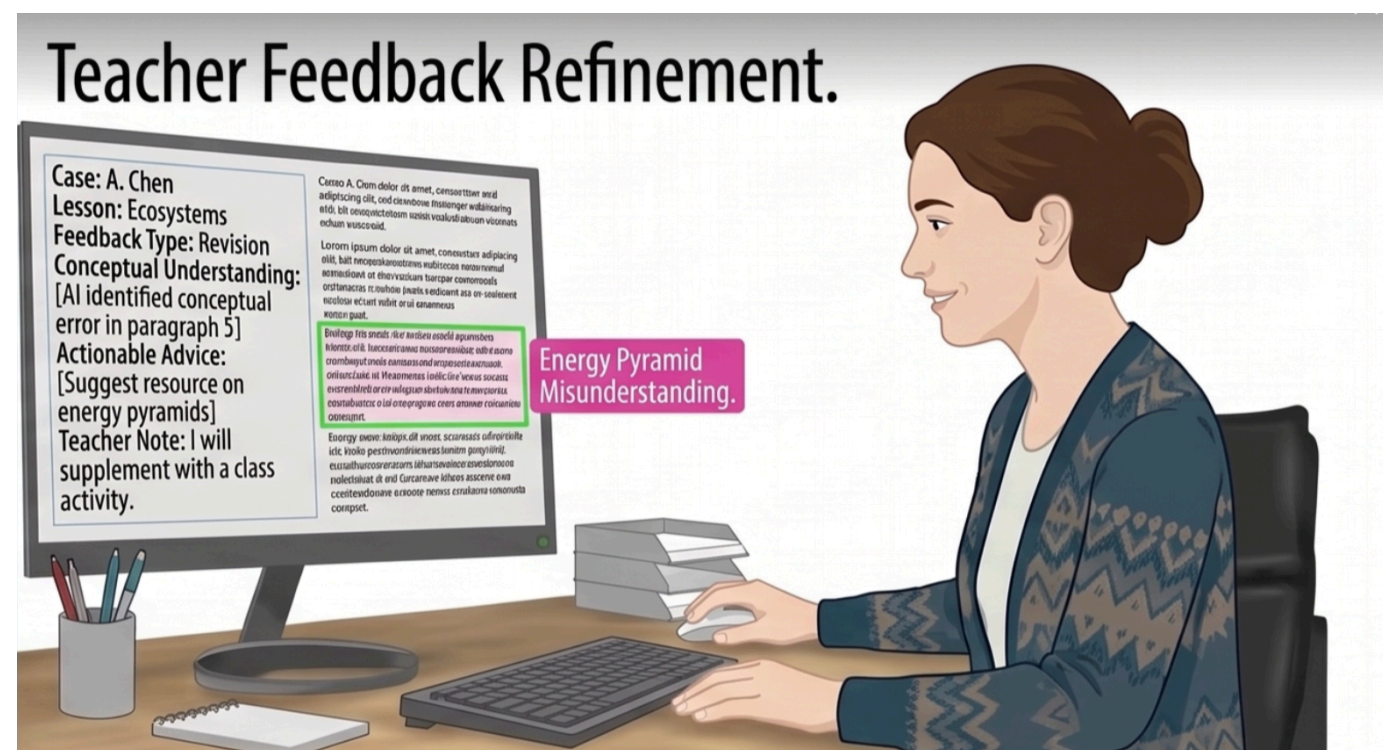
**What are the challenges in achieving IA?**

**How to teach IA for the “AI generation”?**

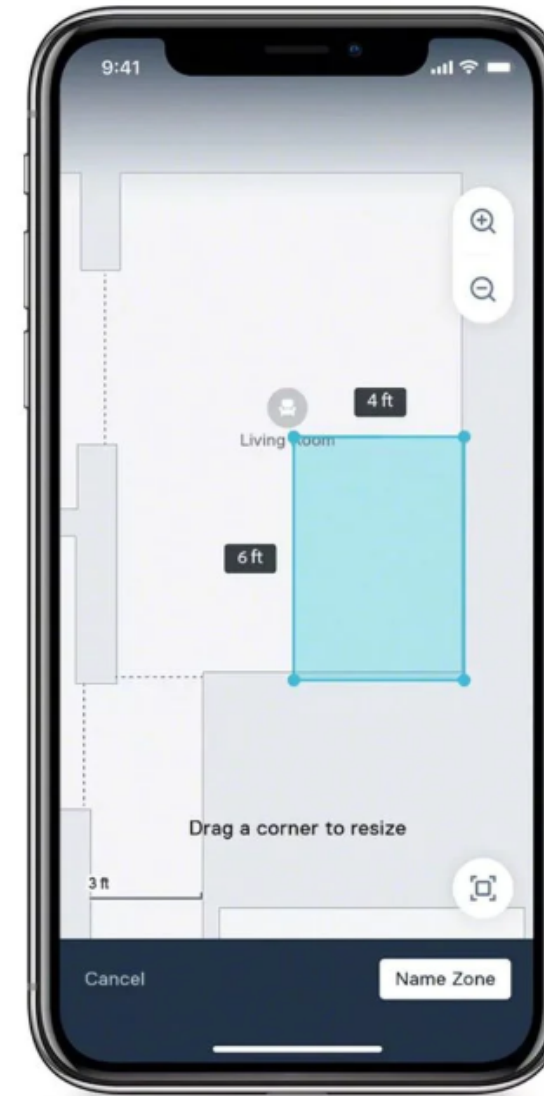


People have always developed IA tools to **offload** some cognitive process to make the whole tasks better, easier, faster...

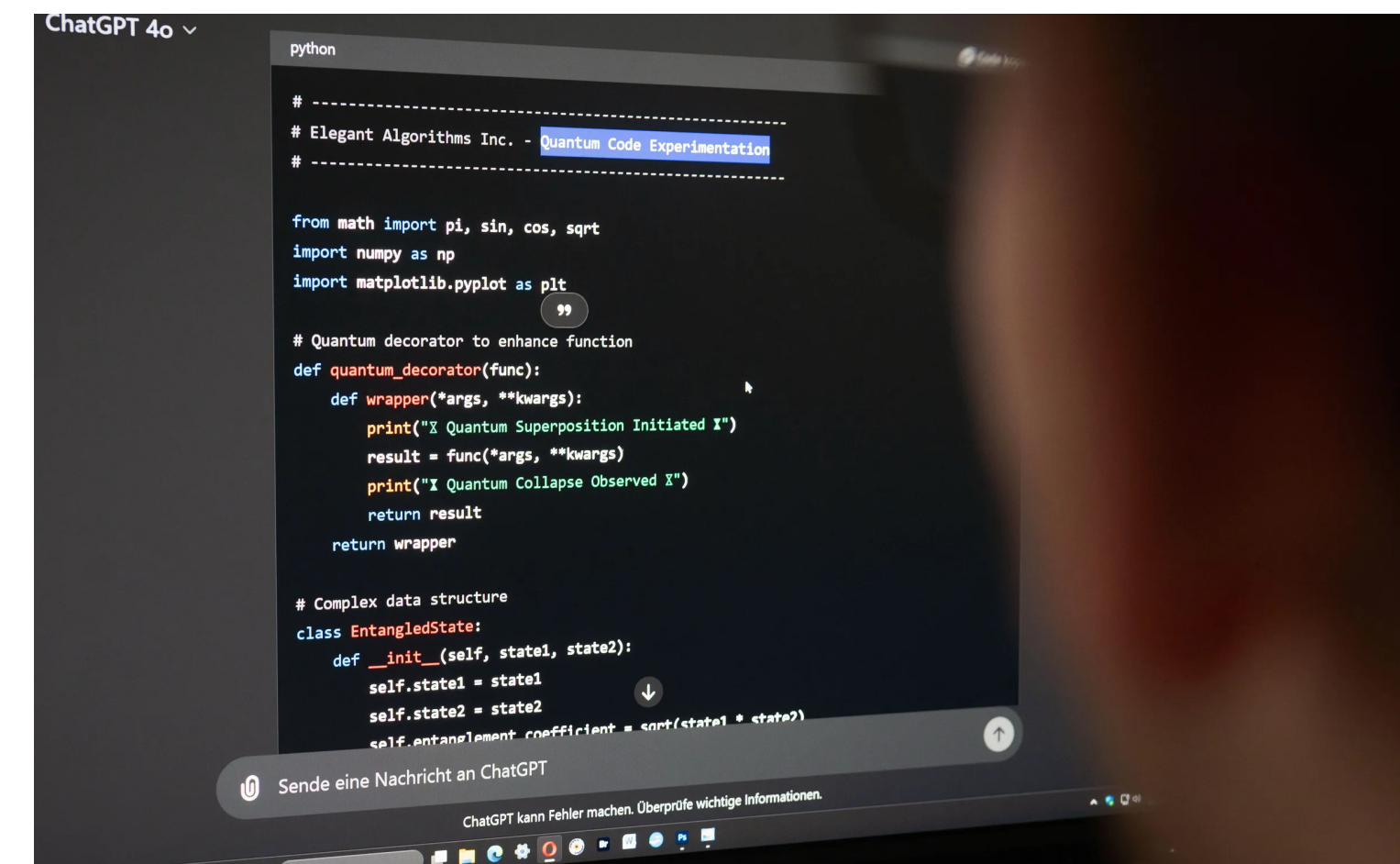
Increasing offloadability



AI suggests, human decides



Human plans, AI executes



AI plans, human fixes

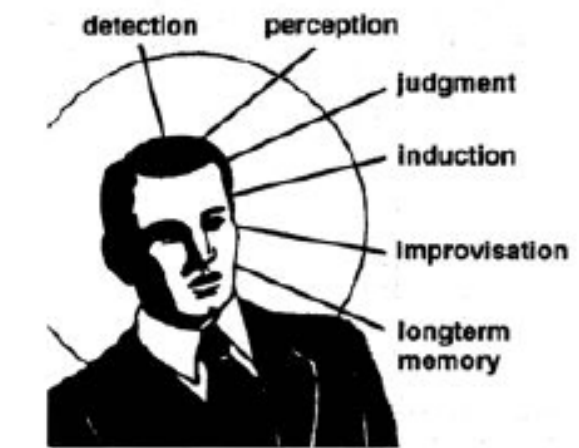
Increasingly capable and agentic AI

## MAN-COMPUTER SYMBIOSIS

Separable Functions of Men and Computers in the Anticipated Symbiotic Association... Men will set the goals and supply the motivations, of course, at least in the early years. They will formulate hypotheses... The information-processing equipment, for its part, will convert hypotheses into testable models and then test the models against data.

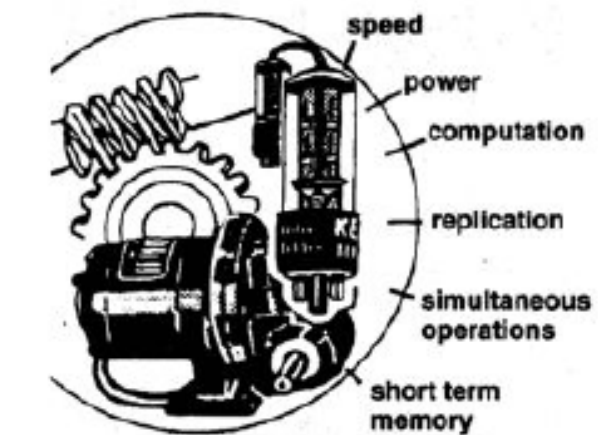
—J.C.R. Licklider

### HUMANS SURPASS MACHINES IN THE:



- Ability to detect small amounts of visual or acoustic energy
- Ability to perceive patterns of light or sound
- Ability to improvise and use flexible procedures
- Ability to store very large amounts of information for long periods and to recall relevant facts at the appropriate time
- Ability to reason inductively
- Ability to exercise judgment

### MACHINES SURPASS HUMANS IN THE:



- Ability to respond quickly to control signals, and to apply great force smoothly and precisely
- Ability to perform repetitive, routine tasks
- Ability to store information briefly and then to erase it completely
- Ability to reason deductively, including computational ability
- Ability to handle highly complex operations, i.e., to do many different things at once.

Offloading “what machines are better at”  
to **complement** people



Generated with Pixtral Version 24.09 in the free version of mistral.ai;

**Oversight** (“executive power”) retained by people

## Article 14: Human Oversight

Humans oversight personnel should...

- ... prevent or minimise the risks to health, safety or fundamental rights
- ... properly monitor the AI system and understand its capacities and limitations
- ... remain aware of the possible tendency of automatically (over-)relying on the output produced by the AI system (automation bias)
- ... correctly interpret AI output and decide not to use / disregard AI
- ...to intervene in the operation of the high-risk AI system or interrupt the system through a 'stop' button or a similar procedure that allows the system to come to a halt in a safe state.
- ...

AI EU Act mandates oversight for high-risk AI systems

**Oversight** (“executive power”) retained by people

## Article 14: Human Oversight

Humans oversight personnel should...

- ... prevent or minimise the risks to health, safety or fundamental rights
- ... properly monitor the AI system and understand its capacities and limitations
- ... remain aware of the possible tendency of automatically (over-)reacting to data produced by the AI system (automation bias)
- ... correctly interpret AI output and decide not to use / disregard AI output
- ...to intervene in the operation of the high-risk AI system or interrupt the system through a 'stop' button or a similar procedure that allows the system to come to a halt in a safe state.
- ...

AI EU Act mandates oversight for high-risk AI systems

3. Education and vocational training:

- (a) AI systems intended to be used to determine access or admission or to assign natural persons to educational and vocational training institutions at all levels;
- (b) AI systems intended to be used to evaluate learning outcomes, including when those outcomes are used to steer the learning process of natural persons in educational and vocational training institutions at all levels;
- (c) AI systems intended to be used for the purpose of assessing the appropriate level of education that an individual will receive or will be able to access, in the context of or within educational and vocational training institutions at all levels;
- (d) AI systems intended to be used for monitoring and detecting prohibited behaviour of students during tests in the context of or within educational and vocational training institutions at all levels.

**Oversight** (“executive power”) retained by people

- **Overreliance on AI: challenges with complementarity and oversight**
- **New affordances of GenAI create new threat to human thinking: risks of offloading**

# **Research Thread 1: Investigating and Mitigating Overreliance on AI**

Challenges with complementarity and oversight

**1A: LLMs for Educational Decision Support** 📅

🕒 10:30 AM – 12:00 PM

📍 Kongesal 4

Concurrent Session

---

**5 Subsessions**

- **1A-01: Data-Driven Evaluation of LLM-Based from Programming Learning Content**  
🕒 10:30 AM – 12:00 PM
- **1A-02: Empowering Multimodal Learning A Comprehensive Platform for Simulation-based Assessment**  
🕒 10:30 AM – 12:00 PM
- **1A-03: Embedding-Based Rankings of Educ Learning Outcome Alignment: Benchmarking Performance**  
🕒 10:30 AM – 12:00 PM
- **1A-04: Investigating Self-regulated Learning AI-based Intelligent Tutoring System**  
🕒 10:30 AM – 12:00 PM
- **1A-05: AI Knows Best? The Paradox of Experience in Educational Tutoring Decision-Making Tasks**  
🕒 10:30 AM – 12:00 PM

**2C: Teaching and Learning Supported by Generative Models** 📅

🕒 1:00 PM – 2:30 PM

📍 Scandic Bergen City Dragefjellet

Concurrent Session

---

**5 Subsessions**

- **2C-01: Structuring versus Problematizing Learning in Diagnostic Reasoning**  
🕒 1:00 PM – 2:30 PM
- **2C-02: Modeling Longitudinal Student Performance Models**  
🕒 1:00 PM – 2:30 PM
- **2C-03: Prompting for Teachability: Designing Learning by Teaching Contexts**  
🕒 1:00 PM – 2:30 PM
- **2C-04: Evaluating Large Language Model Professional Vision: Prompting Strategies**  
🕒 1:00 PM – 2:30 PM
- **2C-05: The Blind Spots in Automated Feedback Writing**  
🕒 1:00 PM – 2:30 PM

**3B: LLM-supported feedback** 📅

🕒 3:00 PM – 4:30 PM

📍 Dræggen 7

Concurrent Session

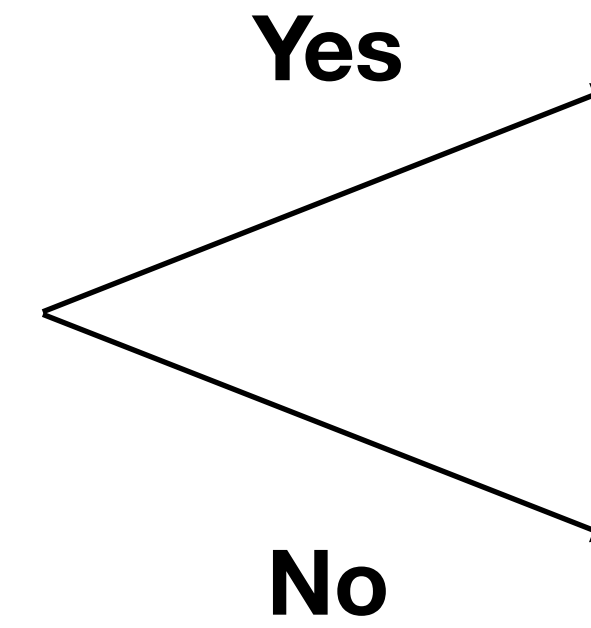
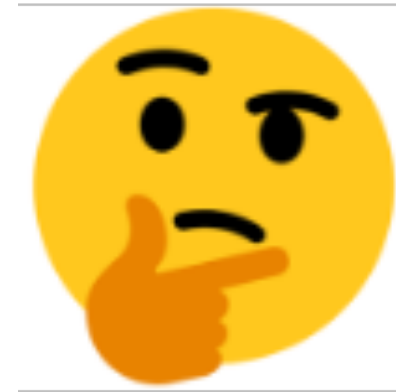
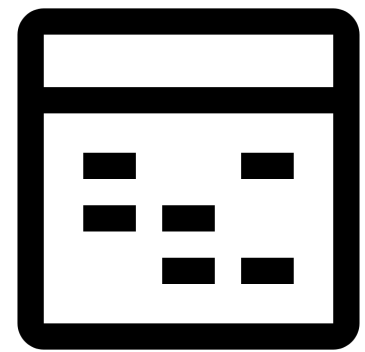
---

**5 Subsessions**

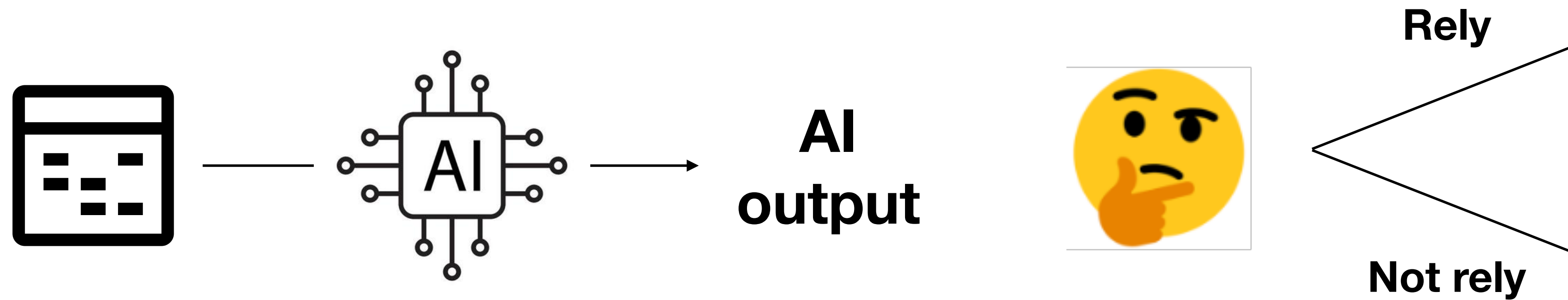
- **3B-01: LLM-based Multimodal Feedback Produces Equivalent Learning and Better Student Perceptions than Educator-written Feedback**  
🕒 3:00 PM – 4:30 PM
- **3B-02: MedSimAI: Simulation and Formative Feedback Generation to Enhance Deliberate Practice in Medical Education**  
🕒 3:00 PM – 4:30 PM
- **3B-03: From Data to Dialogue: Using AI to Scale Feedback-Informed Learning**  
🕒 3:00 PM – 4:30 PM
- **3B-04: When Cognitive Offloading Masks Authentic Ability: Using GenAI-Driven Metacognitive Feedback to Support Valid Game-Based Assessment**  
🕒 3:00 PM – 4:30 PM
- **3B-05: The Effects of AI Feedback on College Students' Reading and Writing Performance in an Intelligent Text Framework**  
🕒 3:00 PM – 4:30 PM

# AI-Assisted Decision-Making: AI Suggests, Human Decides

# Decision-Making



# AI Assisted Decision-Making

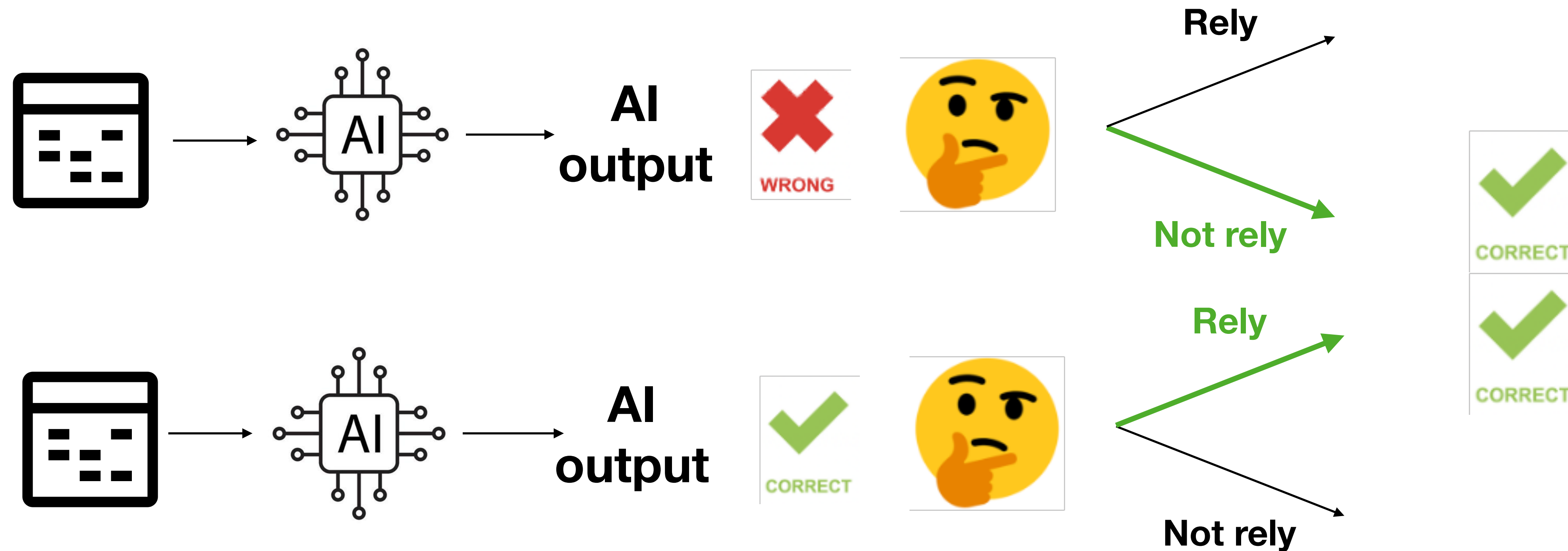


# AI Assisted Decision-Making



**Human oversight** to catch AI errors

# AI Assisted Decision-Making

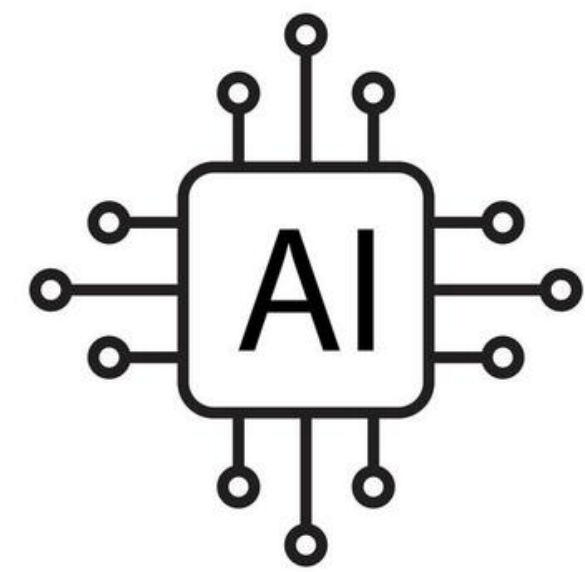


**Human-AI complementarity** to make more correct decisions than either human or AI would have done alone



**Not Rely**

**Rely**

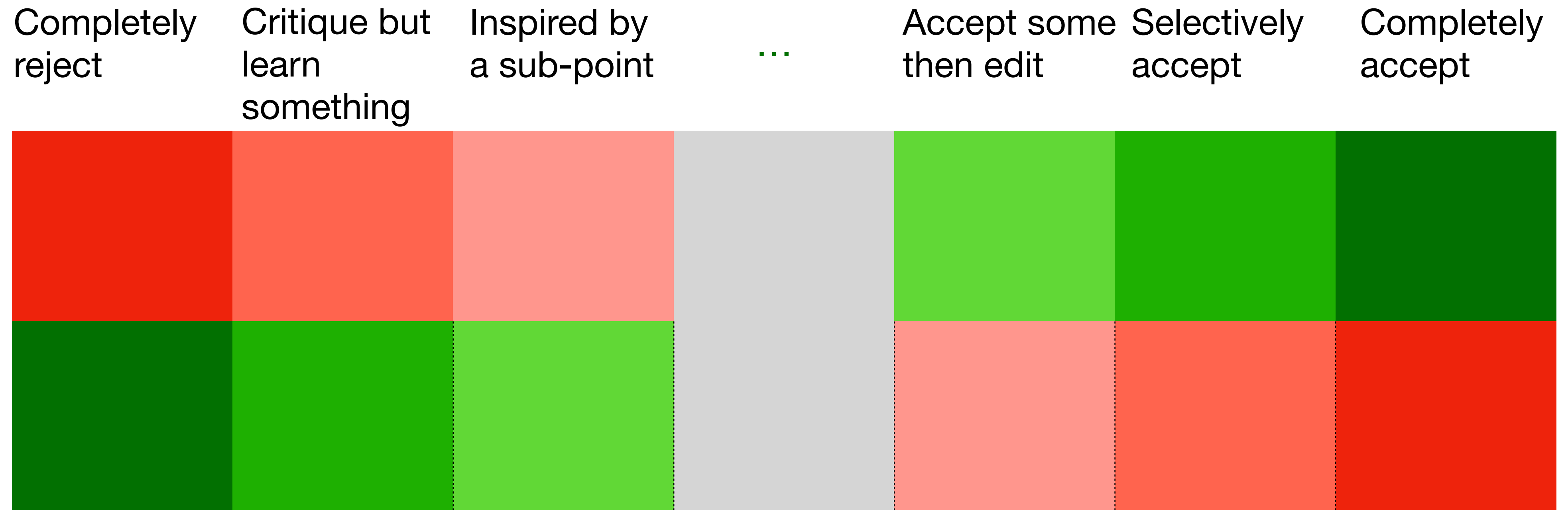


**AI Correct**

**AI Incorrect**

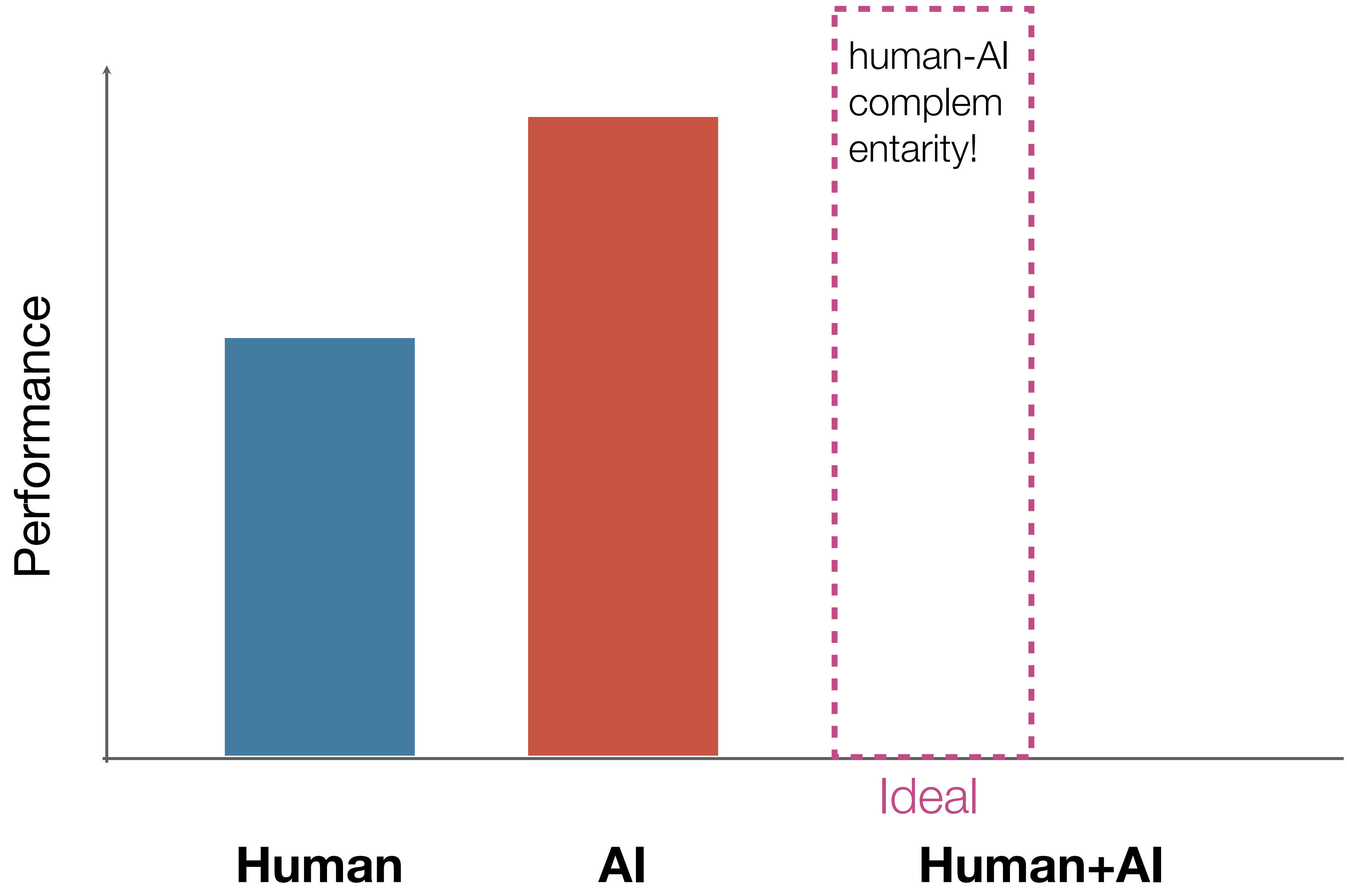
	<b>Not Rely</b>	<b>Rely</b>
<b>AI Correct</b>	Underreliance	Correct Reliance
<b>AI Incorrect</b>	Correct Non-Reliance	Overreliance

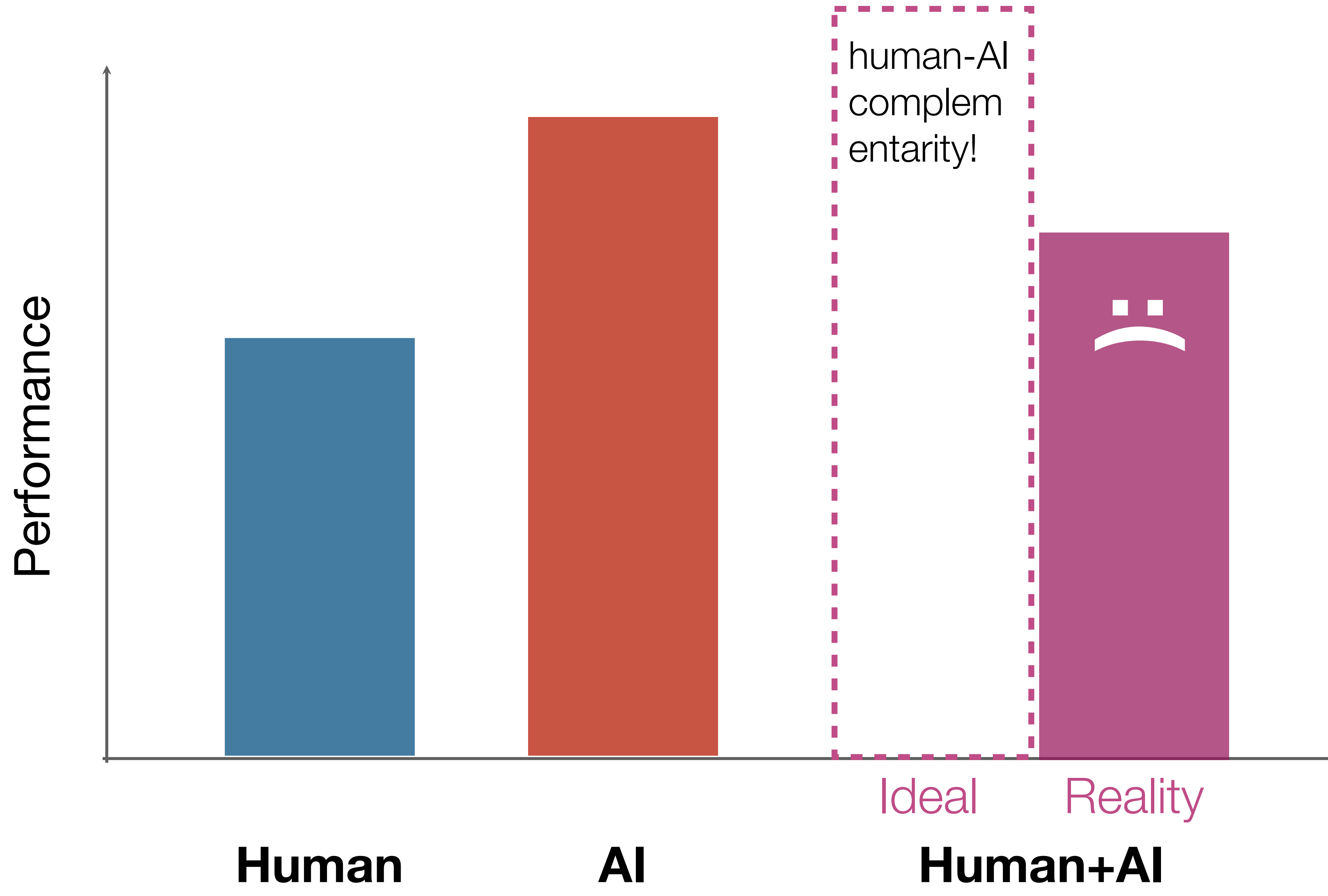
Complementarity requires **appropriate reliance (offloading)**:  
rely when and only when the AI output is correct (**green zones**)



Side note: The reliance action space may not be binary depending on the type of AI

The AI-assisted decision-making literature started with largely classification models





Performance

**Human**

**AI**

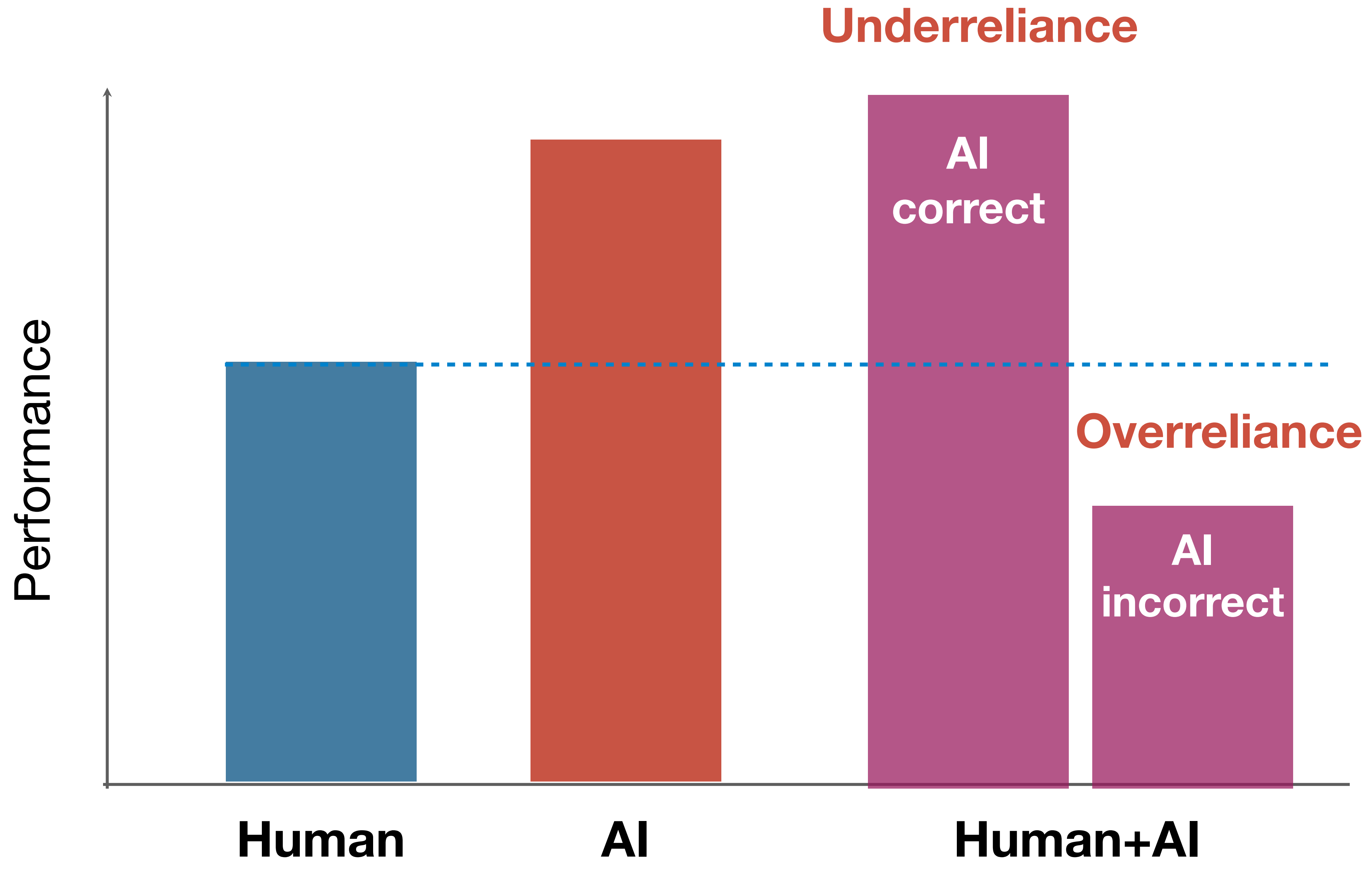
**Human+AI**

human-AI  
complem  
entarity!

Ideal

Reality





## Harms from overreliance:

- Poor human-AI performance
- Biased/homogenized errors (often a cause for inequality)

Table 1: Short and long-term impacts of overreliance on LLMs in two domains

Domain	Short-term	Long-term
Healthcare	<p><b>Individual:</b> Medical doctor follows incorrect AI diagnosis</p> <p><b>Institutional:</b> Hospital systems implement AI diagnostics without adequate verification protocols</p>	<p><b>Individual:</b> Healthcare professionals experience cognitive deskilling and atrophy of diagnostic abilities</p> <p><b>Institutional:</b> Medical departments restructure workflows around AI systems, reducing human oversight</p>
Personal advice	<p><b>Individual:</b> User accepts sycophantic or biased relationship or career advice without critical evaluation</p> <p><b>Societal:</b> Communities adopt similar AI-generated advice, creating homogenized decision patterns</p>	<p><b>Individual:</b> Self-worth becomes derived from AI companion approval [82]</p> <p><b>Societal:</b> Social norms shift toward algorithmic validation of personal choices</p>

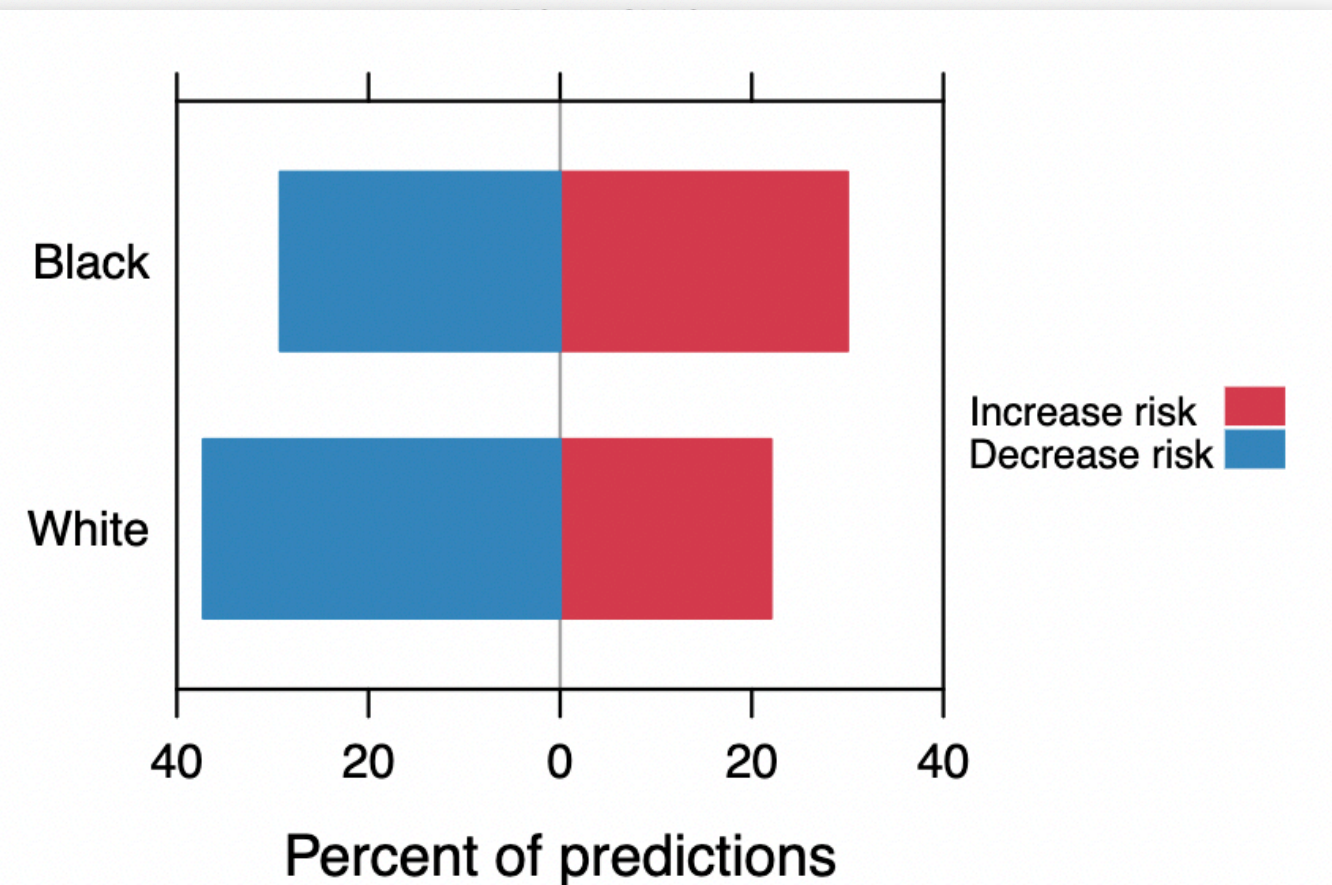
## Disparate Interactions: An Algorithm-in-the-Loop Analysis of Fairness in Risk Assessments

Ben Green  
Harvard University  
bgreen@g.harvard.edu

Yiling Chen  
Harvard University  
yiling@seas.harvard.edu

### ABSTRACT

Despite vigorous debates about the tech assessments being deployed in the U.S., remarkably little research has studied how they influence decision-making processes. After all, risk assessments inform definitive decisions—they inform judges. It is therefore essential that considerations be informed by rigorous studies of how judges use them. This paper takes a first step by studying human interactions with risk assessments in an experimental study on Amazon Mechanical Turk. We find that 1) participants exhibit “disparate interactions,” whereby they deviate to higher risk predictions about black defendants than white defendants. The results suggest that for a new “algorithm-in-the-loop” framework to improve human decisions rather than just generating the best prediction in the abstract, it must be grounded in an understanding of their real-world impacts instead of in



**Figure 4: The rate at which participants deviated from the risk assessment’s prediction toward higher and lower levels of risk, broken down by defendant race. When evaluating black defendants, participants were 36.4% more likely to deviate positively from the risk assessment and 21.5% less likely to deviate negatively (participant predictions matched the risk assessment at an equal rate for both races).**

***Biased overreliance:*** When using a biased decision-support, human+AI can get more biased than human or AI alone

(More likely to follow AI errors for demographic groups for which one already hold bias against)

Table 1: Short and long-term impacts of overreliance on LLMs in two domains

Domain	Short-term	Long-term
Healthcare	<p><b>Individual:</b> Medical doctor follows incorrect AI diagnosis</p> <p><b>Institutional:</b> Hospital systems implement AI diagnostics without adequate verification protocols</p>	<p><b>Individual:</b> Healthcare professionals experience cognitive deskilling and atrophy of diagnostic abilities</p> <p><b>Institutional:</b> Medical departments restructure workflows around AI systems, reducing human oversight</p>
Personal advice	<p><b>Individual:</b> User accepts sycophantic or biased relationship or career advice without critical evaluation</p> <p><b>Societal:</b> Communities adopt similar AI-generated advice, creating homogenized decision patterns</p>	<p><b>Individual:</b> Self-worth becomes derived from AI companion approval [82]</p> <p><b>Societal:</b> Social norms shift toward algorithmic validation of personal choices</p>

## Harms from overreliance:

- Poor human-AI performance
- Biased/homogenized errors (often a cause for inequality)
- Deskilling
- Psychological harms from loss of agency
- Infrastructural vulnerabilities
- Shifting social engagement and norms
- ...

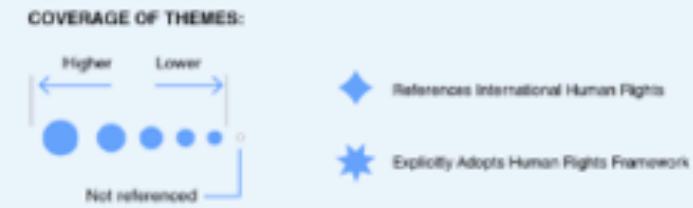
A hard lesson from research on human-AI interaction:  
**Interventions to mitigate overreliance often do not work well, and sometimes backfire**

# PRINCIPLED ARTIFICIAL INTELLIGENCE

A Map of Ethical and Rights-Based Approaches to Principles for AI

Authors: Jessica Fjeld, Nele Achten, Hannah Hilligoss, Adam Nagy, Madhulika Srikumar  
 Designers: Arushi Singh (arushisingh.net) and Melissa Axelrod (melissaaxelrod.com)

**HOW TO READ:**  
 Date, Location  
**Document Title**  
 Actor



The size of each dot represents the percentage of principles in that theme contained in the document. Since the number of principles per theme varies, it's informative to compare dot sizes within a theme but not between themes.

The principles within each theme are:

- Privacy:** Privacy, Control over Use of Data, Consent, Privacy by Design, Recommendation for Data Protection Laws, Ability to Restrict Processing, Right to Rectification, Right to Erasure
- Accountability:** Accountability, Recommendation for New Regulations, Impact Assessment, Evaluation and Auditing Requirement, Verifiability and Replicability, Liability and Legal Responsibility, Ability to Appeal, Environmental Responsibility, Creation of a Monitoring Body, Remedy for Automated Decision
- Safety and Security:** Security, Safety and Reliability, Predictability, Security by Design
- Transparency and Explainability:** Explainability, Open Source Data and Algorithms, Notification when Interacting with an AI, Notification when AI Makes a Decision about an Individual, Regular Reporting Requirement, Open Procurement (for Government)
- Fairness and Non-discrimination:** Non-discrimination and the Prevention of Bias, Fairness, Inclusiveness in Design, Inclusiveness in Impact, Representative and High Quality Data, Equality
- Human Control of Technology:** Human Control of Technology, Human Review of Automated Decision, Ability to Opt out of Automated Decision
- Professional Responsibility:** Multistakeholder Collaboration, Responsible Design, Consideration of Long Term Effects, Accuracy, Scientific Integrity
- Promotion of Human Values:** Leveraged to Benefit Society, Human Values and Human Flourishing, Access to Technology

Further information on findings and methodology is available in *Principled Artificial Intelligence: Mapping Consensus in Ethical and Rights-Based Approaches* (Berkman Klein, 2020) available at [cyber.harvard.edu](http://cyber.harvard.edu)



Fairness and Non-discrimination

Transparency and Explainability

Human Control

Safety and Security

Accountability

Privacy

Promotion of Human Values

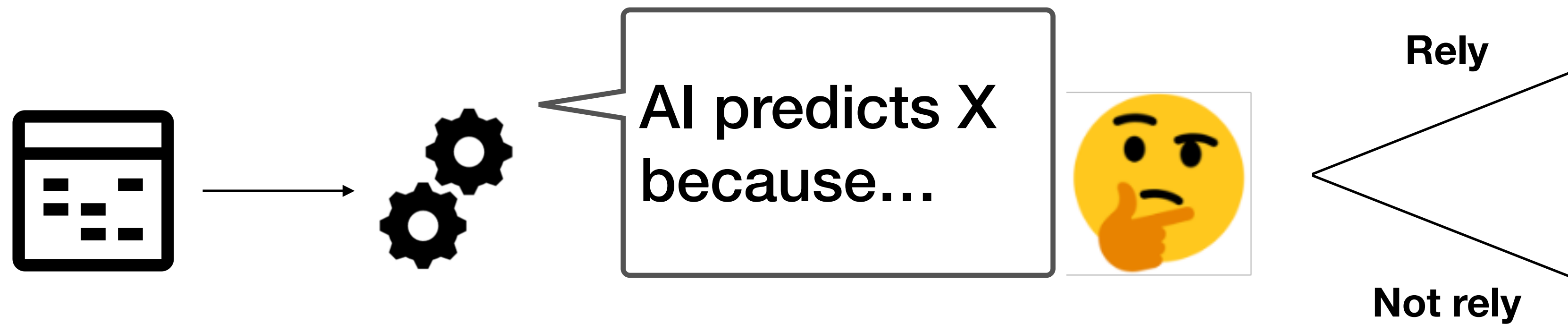
Mapping of 36 ethical and responsible AI frameworks (Berkman Klein Center)

## Article 14: Human Oversight

Humans oversight personnel should...

- ... prevent or minimise the risks to health, safety or fundamental rights
- ... properly monitor the AI system and understand its capacities and limitations
- ... remain aware of the possible tendency of automatically (over-)relying on the output produced by the AI system (automation bias)
- ... correctly interpret AI output and decide not to use / disregard AI
- ...to intervene in the operation of the high-risk AI system or interrupt the system through a 'stop' button or a similar procedure that allows the system to come to a halt in a safe state.
- ...

# Explainable AI for Decision-Making



The hope: explanation can help people have better oversight of AI errors, hence less overreliance

Does this person belong to a **high-income** or **low-income** group?

Attributes	Value
Age	33
Sex	Male
Race	White
Marital status	Married
Years of edu.	9
Workclass	Private
Occupation	Sales
Hrs. per week	45

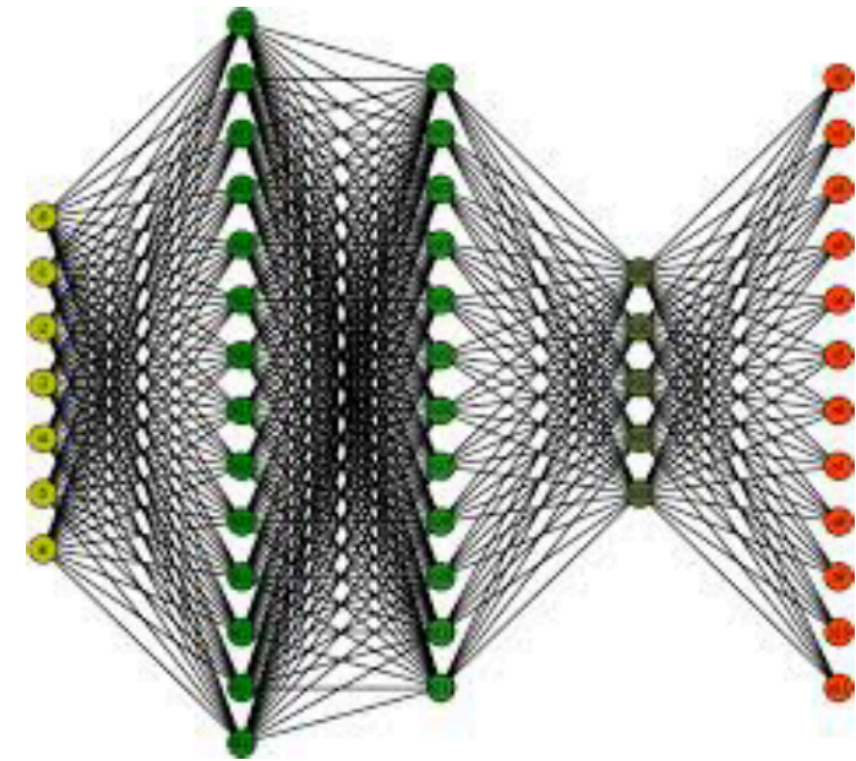
Does this person belong to a **high-income** or **low-income** group?

Attributes	Value
Age	33
Sex	Male
Race	White
Marital status	Married
Years of edu.	9
Workclass	Private
Occupation	Sales
Hrs. per week	45

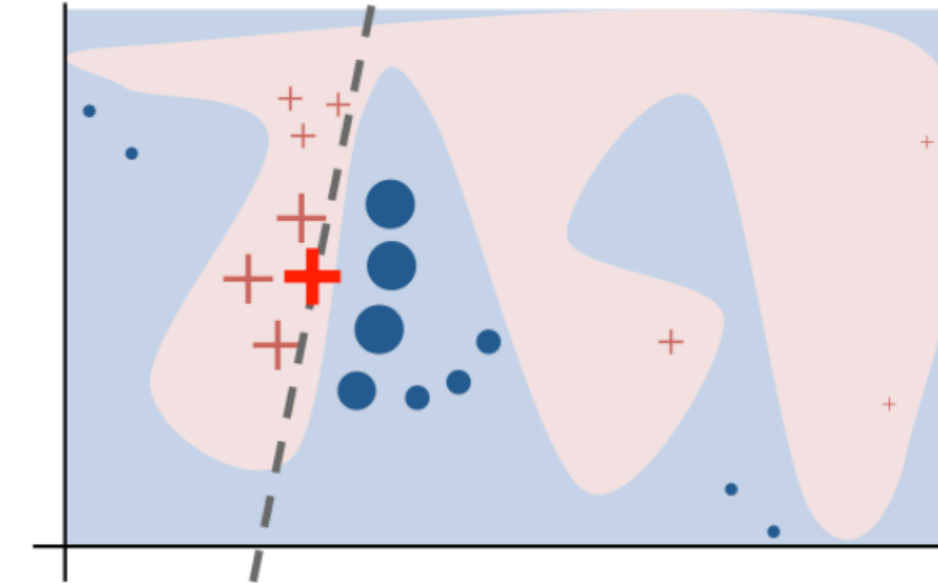
The AI system predicts that the customer belongs to **high-income** group

**Control condition**

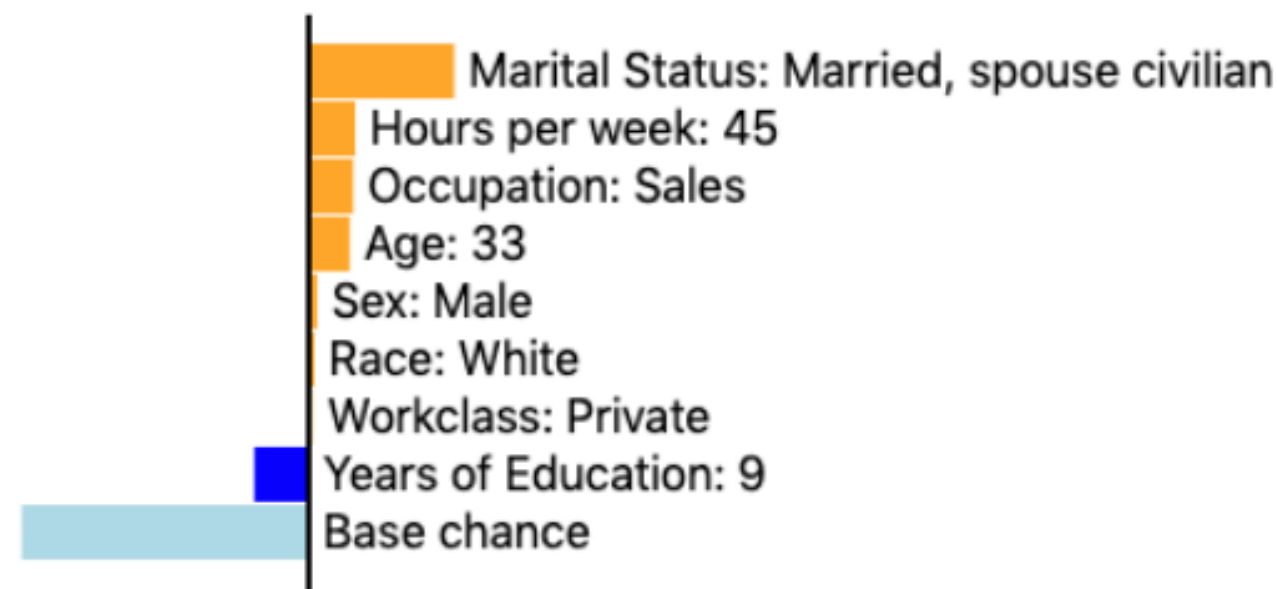
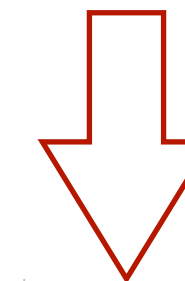
# Popular XAI Algorithms Produce Feature Importance Explanations



Neural network, not directly explainable



Use a *post-hoc* XAI technique



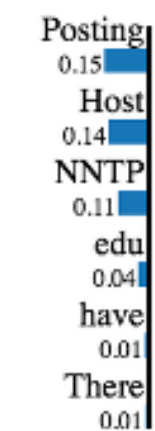
Tabular data

Images (explaining prediction of 'Cat' in pros and cons)



Image

atheism



christian

**Text with highlighted words**

From: johnchad@triton.unm.edu (jchadwic)  
Subject: Another request for Darwin Fish  
Organization: University of New Mexico, Albuquerque  
Lines: 11  
NNTP-Posting-Host: triton.unm.edu

Hello Gang,

There have been some notes recently asking where to obtain the DARWIN fish. This is the same question I have and I have not seen an answer on the net. If anyone has a contact please post on the net or email me.

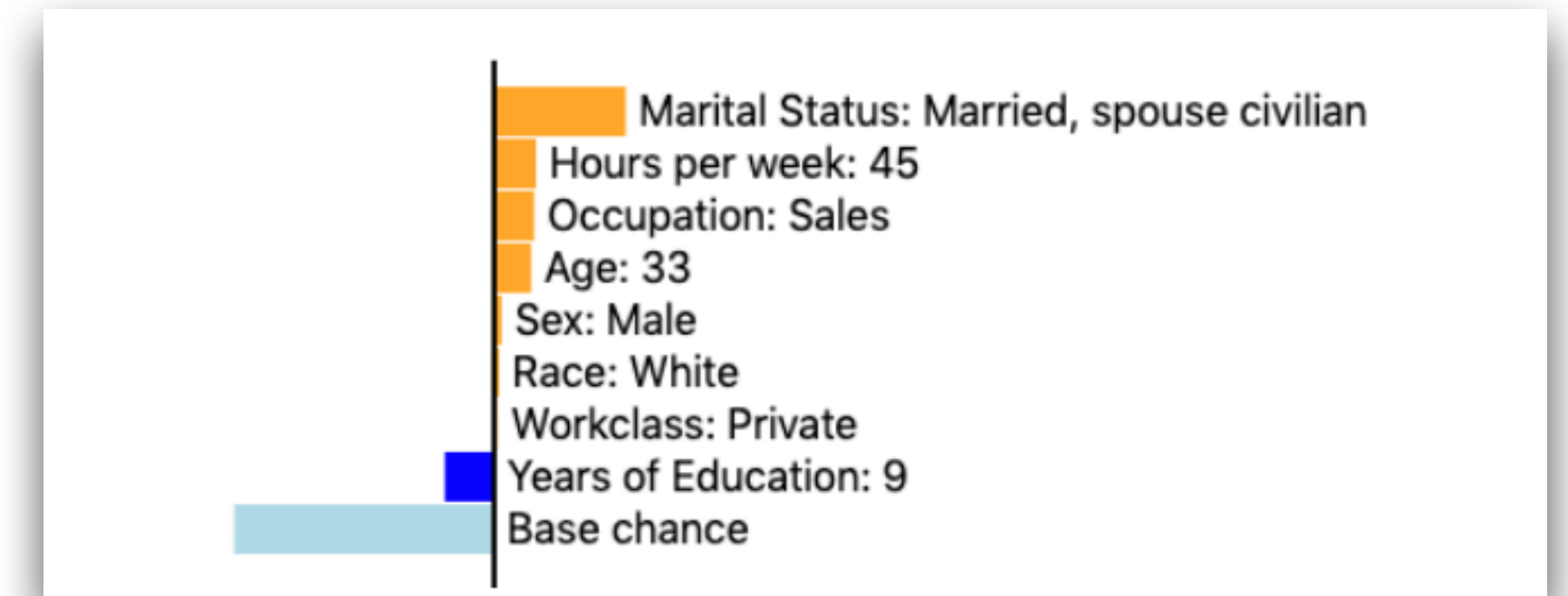
Texts

Does this person belong to a **high-income** or **low-income** group?

Attributes	Value
Age	33
Sex	Male
Race	White
Marital status	Married
Years of edu.	9
Workclass	Private
Occupation	Sales
Hrs. per week	45

The AI system predicts that the customer belongs to **high-income** group

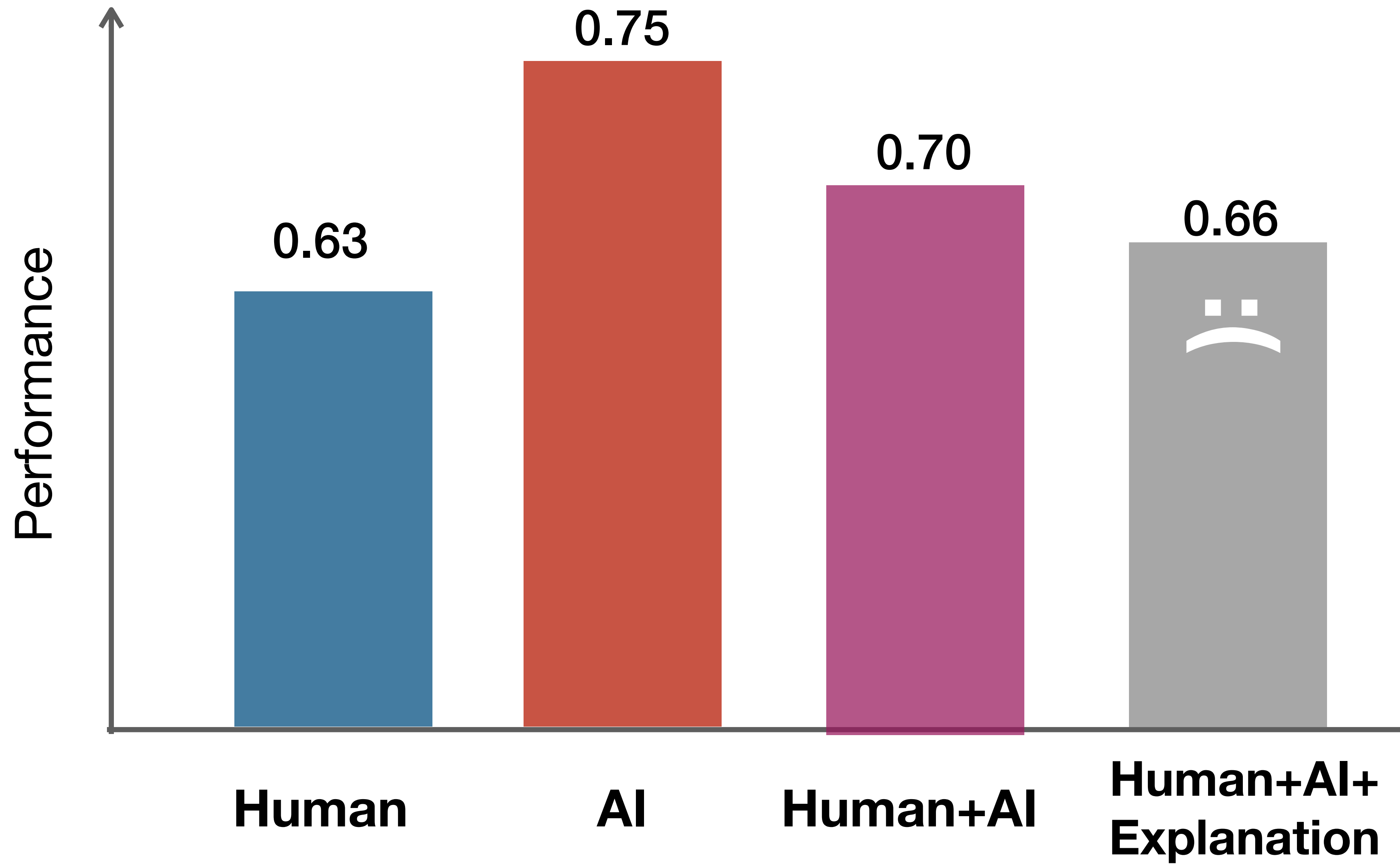
**Control condition**

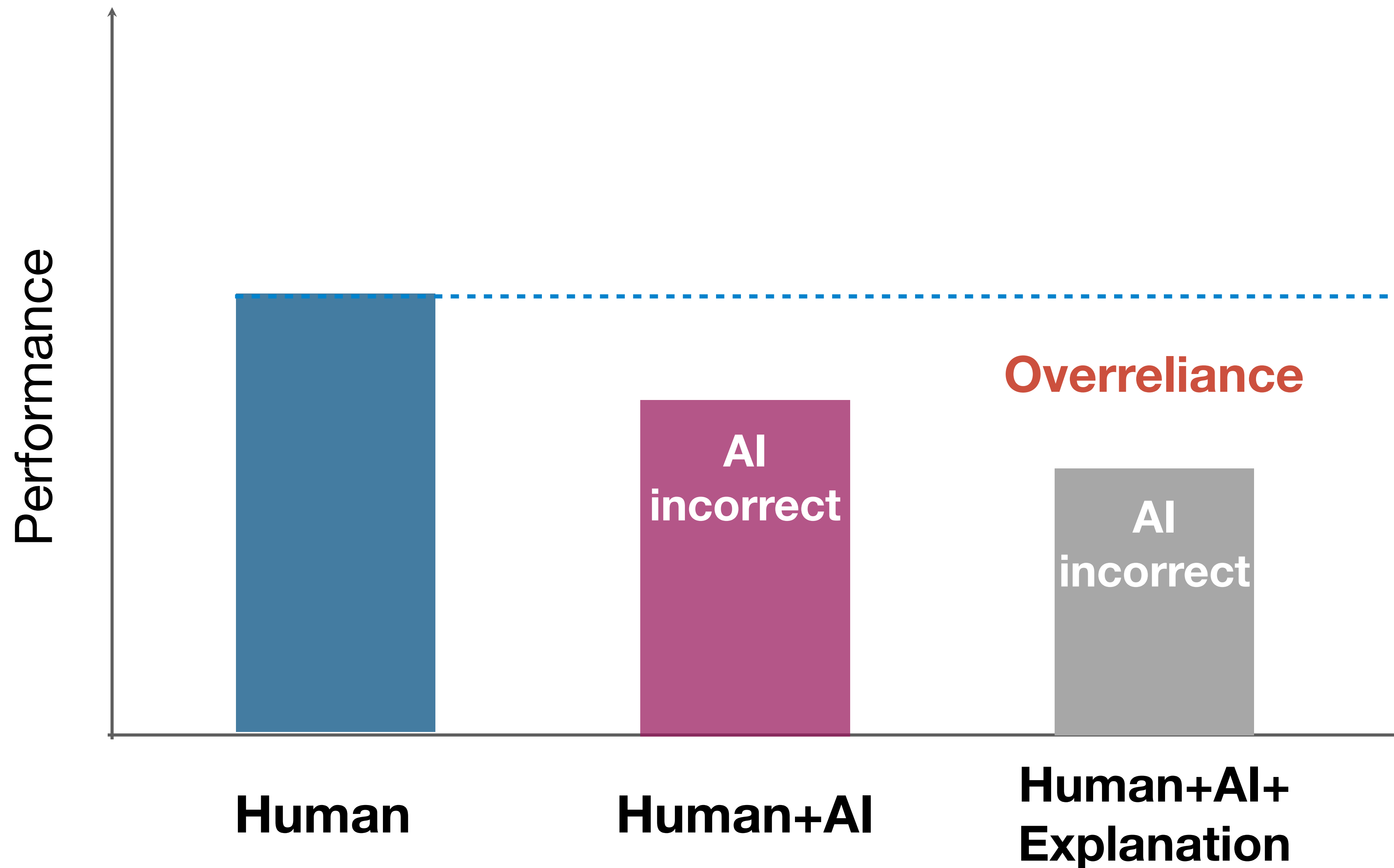


The figure explains how the system arrives at this decision according to how much each attribute likely contributes to a **high** or **low** income

The AI system predicts that the customer belongs to **high-income** group

**Experimental condition**





I, like others **was very excited to read this book**. I thought it would show another side to how the Tate family dealt with the murder of their daughter Sharon. I didn't have to read much to realize however that the book is was not going to be what I expected. It makes it hard to tell the story. The events begin. It reads like the next part of this family's troubles. The features of Sharon that **hardly worth the price**

AI recommendation: **Suitable** Saliency legend: High relevance Medium

**Anna Müller**  
Munich, Germany · a.mueller@email.de · +49 89 1234567

PROFESSIONAL EXPERIENCE

**Senior Quality Manager** 2019 – present  
Bavarian Auto Group, Munich

- Led cross-functional quality assurance team of 12 across production lines
- Implemented ISO 9001:2015 certification process, reducing defect rate by 34%
- Managed supplier audits and vendor compliance reporting

**Quality Engineer** 2015 – 2019  
Technik GmbH, Stuttgart

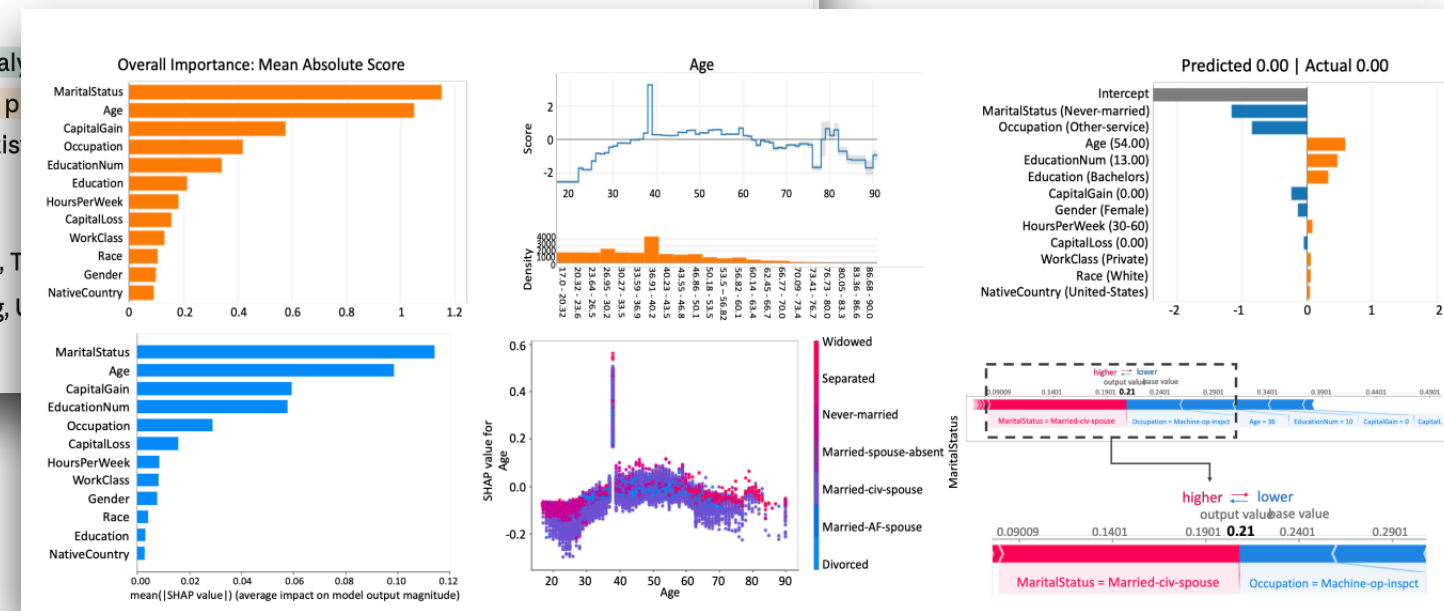
- Conducted root cause analysis
- Coordinated with R&D on product development
- Trained junior staff on statistical process control

EDUCATION

- M.Sc. Industrial Engineering, Technical University of Munich
- B.Sc. Mechanical Engineering, University of Applied Sciences

KEY COMPETENCIES

- Quality Management
- ISO 9001:2015
- Statistical Process Control
- Supplier Management
- Team Leadership



XAI increases overreliance:  
robust findings across use cases  
and modalities

I, like others **was very excited to read this book**. I thought it would show another side to how the Tate family dealt with the murder of their daughter Sharon. I didn't have to read much to realize however that the book is was not going to be what I expected. It makes it hard to tell the story. The events begin. It reads like the rest of this family's troubles. The features of Sharon that **hardly worth the price**

AI recommendation: **Suitable** Saliency legend: High relevance Medium

**Anna Müller**  
Munich, Germany · a.mueller@email.de · +49 89 1234567

PROFESSIONAL EXPERIENCE

**Senior Quality Manager** 2019 – present  
Bavarian Auto Group, Munich

- Led cross-functional quality assurance team of 12 across production lines
- Implemented ISO 9001:2015 certification process, reducing defect rate by 34%
- Managed supplier audits and vendor compliance reporting

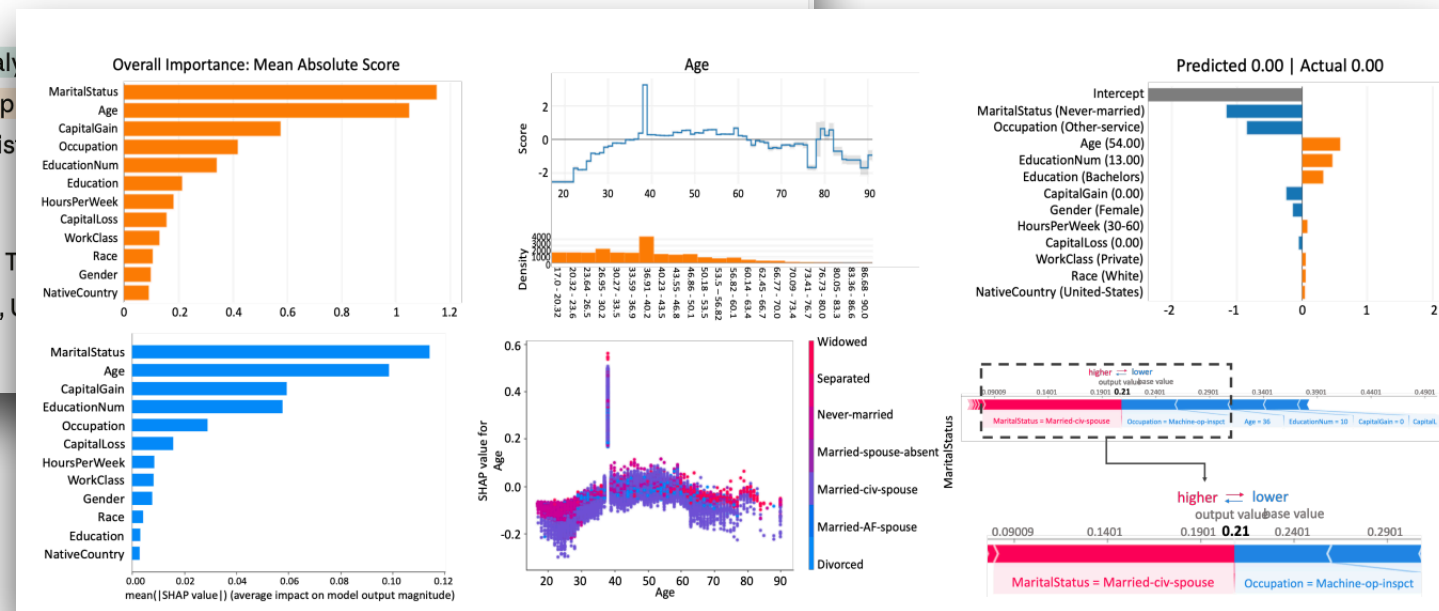
**Quality Engineer** 2015 – 2019  
Technik GmbH, Stuttgart

- Conducted root cause analysis
- Coordinated with R&D on product development
- Trained junior staff on statistical process control

EDUCATION

- M.Sc. Industrial Engineering, University of Stuttgart
- B.Sc. Mechanical Engineering, University of Stuttgart

KEY COMPETENCIES



## Neither

No, not more than two-thirds of South America's population live in Brazil.

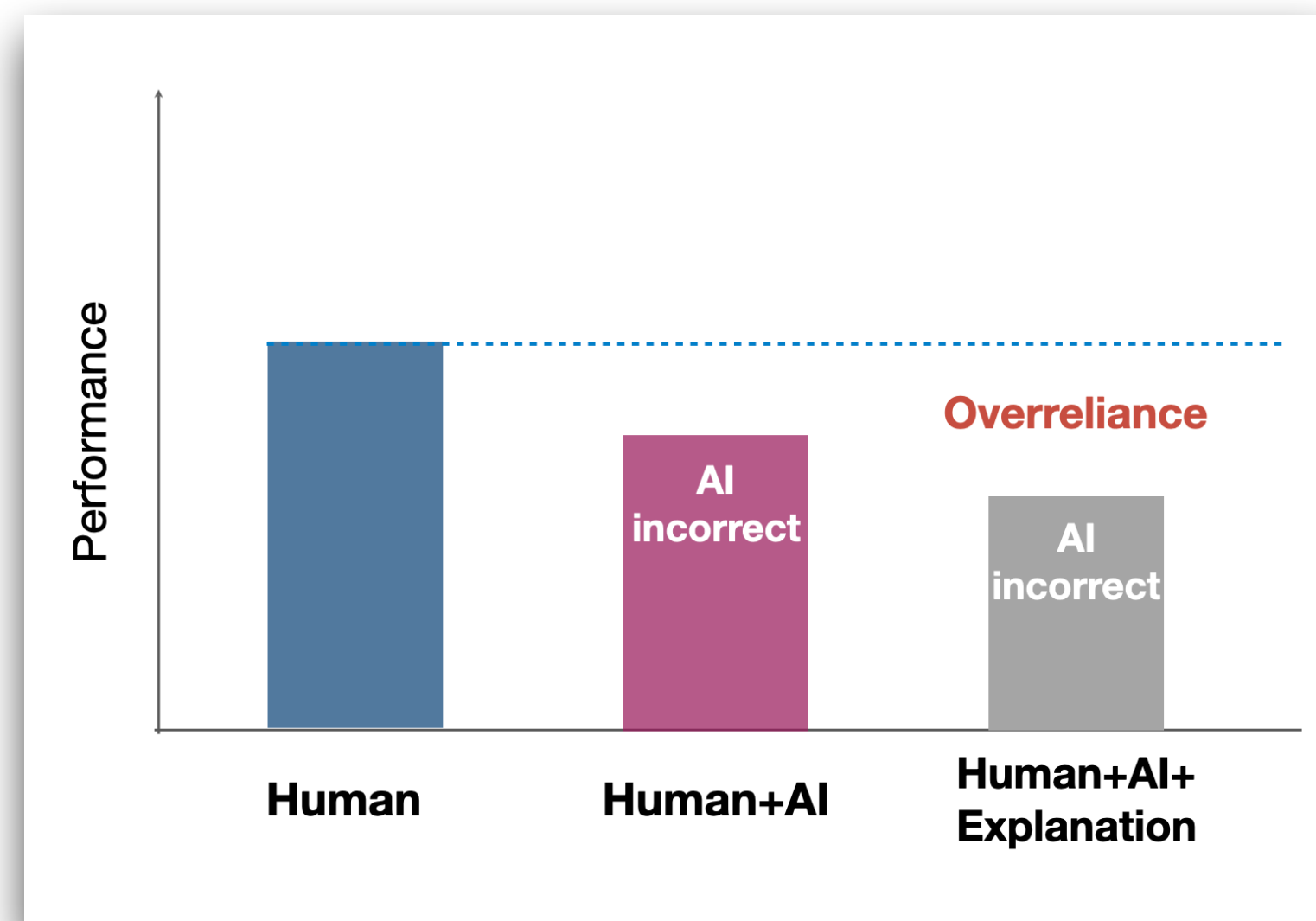
## Explanation only

No, not more than two-thirds of South America's population live in Brazil. Although Brazil is the largest country in South America by both land area and population, it does not encompass the majority of the continent's people. As of recent estimates, Brazil's population is roughly 213 million, while South America's total population is about 440 million. This means Brazil's population constitutes roughly 48% of the continent's total, which is significantly less than two-thirds. Hence, while Brazil has a large population, it does not exceed the two-thirds threshold relative to the entire continent's population.

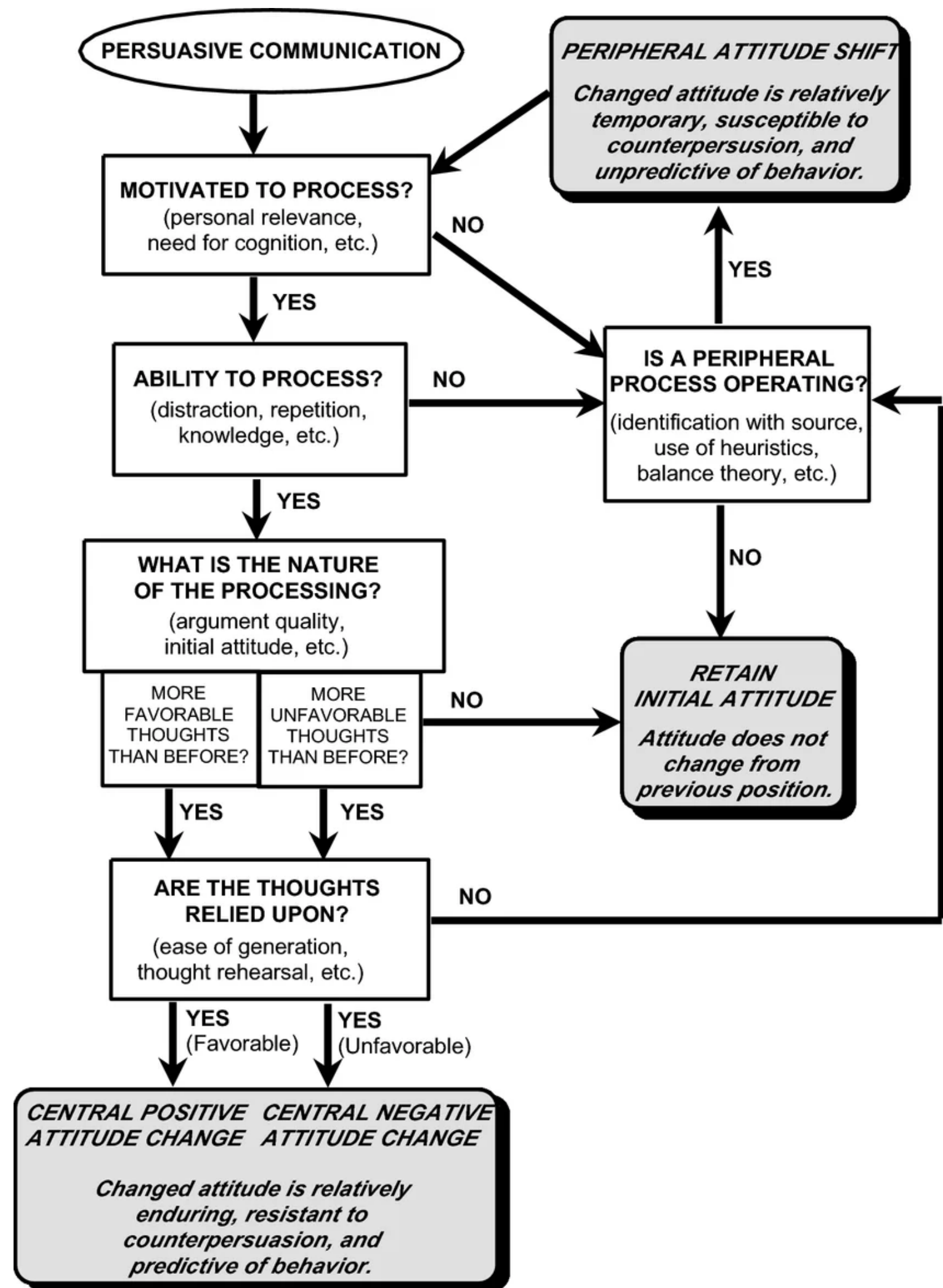
XAI increases overreliance: robust findings across use cases and modalities

LLM reasoning also increases people's overreliance

# Why Do AI Explanations Increase Overreliance?



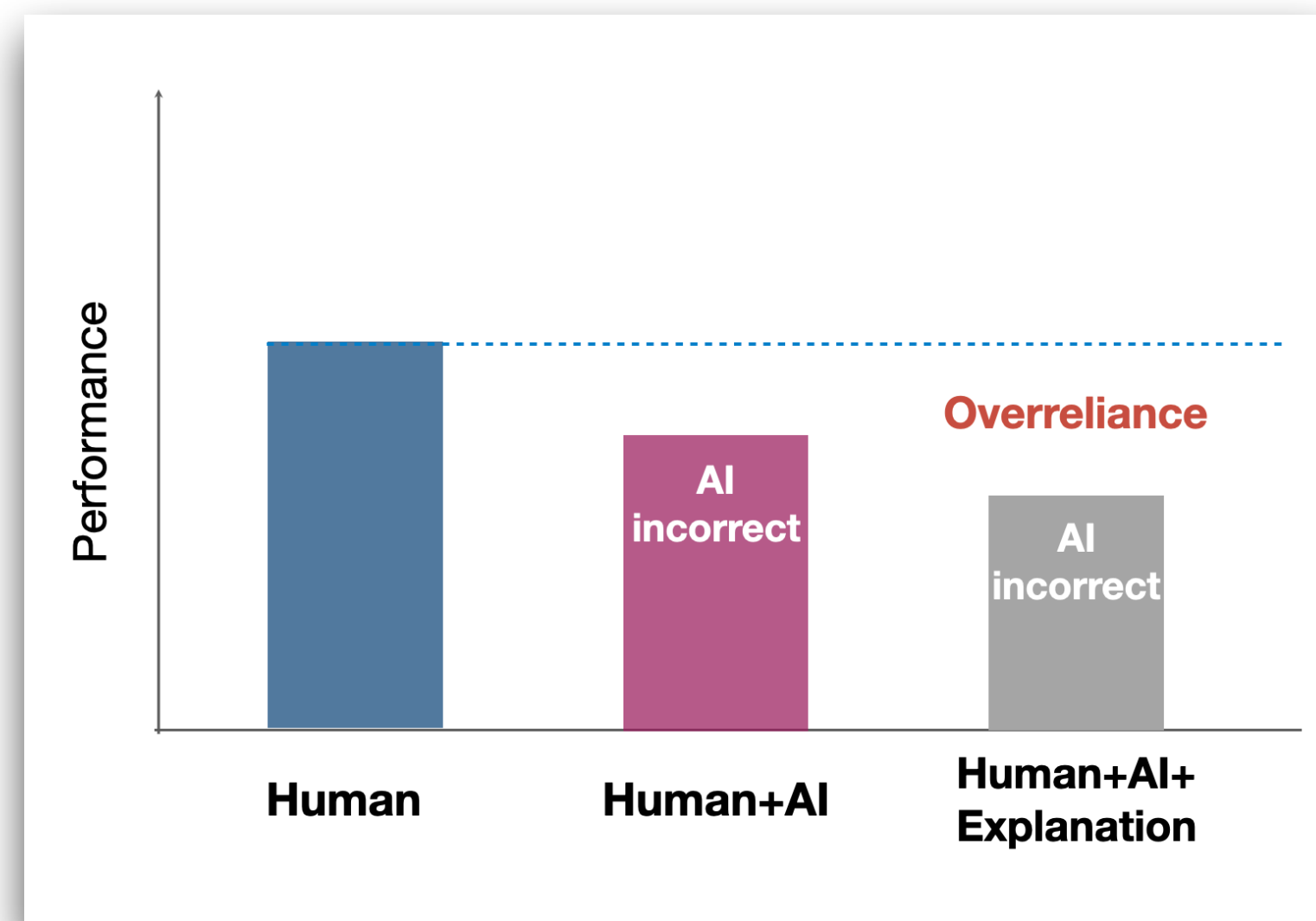
People do not always engage analytically (“system 2” slow thinking) with AI outputs to detect AI errors



Two prerequisites to engage in “system 2” slow thinking:  
**motivation and ability**

**Elaboration Likelihood Model**  
 (Petty and Caccioppo, 1986)

# Why Do AI Explanations Increase Overreliance?



(Feature-importance)  
explanations hinder **motivation**:  
invoke positive heuristic  
(*explainable = trustworthy*)

Also hinder **ability**: distracting  
and disruptive of people's own  
thinking

**What worked (somewhat)?**

**Uncertainty expression** directly signals possible AI errors and calibrates reliance

However, the effectiveness is often hindered by design details of how uncertainty is communicated

Can people recognize the uncertainty?  
Are they **motivated** enough?

**Q. Is Spironolactone an FDA-approved drug for treating acne?**

Not uncertain

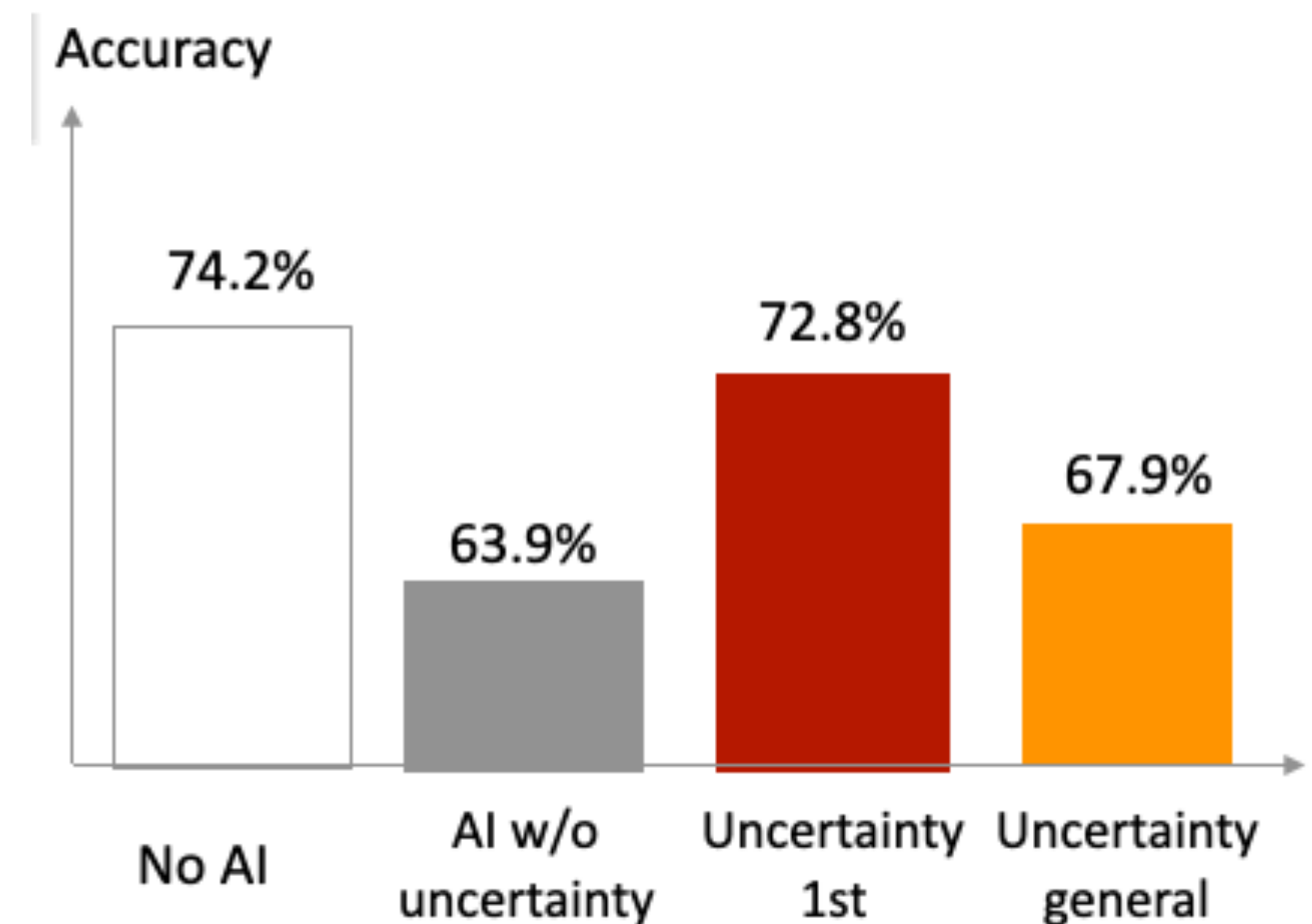
*Yes, Spironolactone is an FDA-approved drug for treating acne. It is primarily used in medicine to treat heart failure and high blood pressure but doctors sometimes prescribe it for acne as well [1, 2, 3, 4].*

Uncertain in the **first-person** perspective

***I'm not sure, but my guess is** Spironolactone is an FDA-approved drug for treating acne. It is primarily used in medicine to treat heart failure and high blood pressure but doctors sometimes prescribe it for acne as well [1, 2, 3, 4].*

Uncertain in the **general** perspective

***There is uncertainty, but it seems like** Spironolactone is an FDA-approved drug for treating acne. It is primarily used in medicine to treat heart failure and high blood pressure but doctors sometimes prescribe it for acne as well [1, 2, 3, 4].*



**Cognitive forcing functions** to slow people down and engage mindfully with with AI outputs are found to improve appropriate reliance

The AI is 87% confident in its suggestion

See AI's suggestion ▾



The AI is processing the image

(b) uncertainty (*SXAI*)

(c) on demand (*CFF*)

(d) wait (*CFF*)

**Cognitive forcing functions** to **slow people down and engage mindfully with with AI outputs** are found to improve appropriate reliance

But they may come with a trade-off of subjective user experience (will they work in the long run?)

The AI is 87% confident in its suggestion

See AI's suggestion ▾



The AI is processing the image

(b) uncertainty (*SXAI*)

(c) on demand (*CFF*)

(d) wait (*CFF*)

**AI literacy support** to help users understand AI can make mistakes and when it makes mistakes is found to be effective

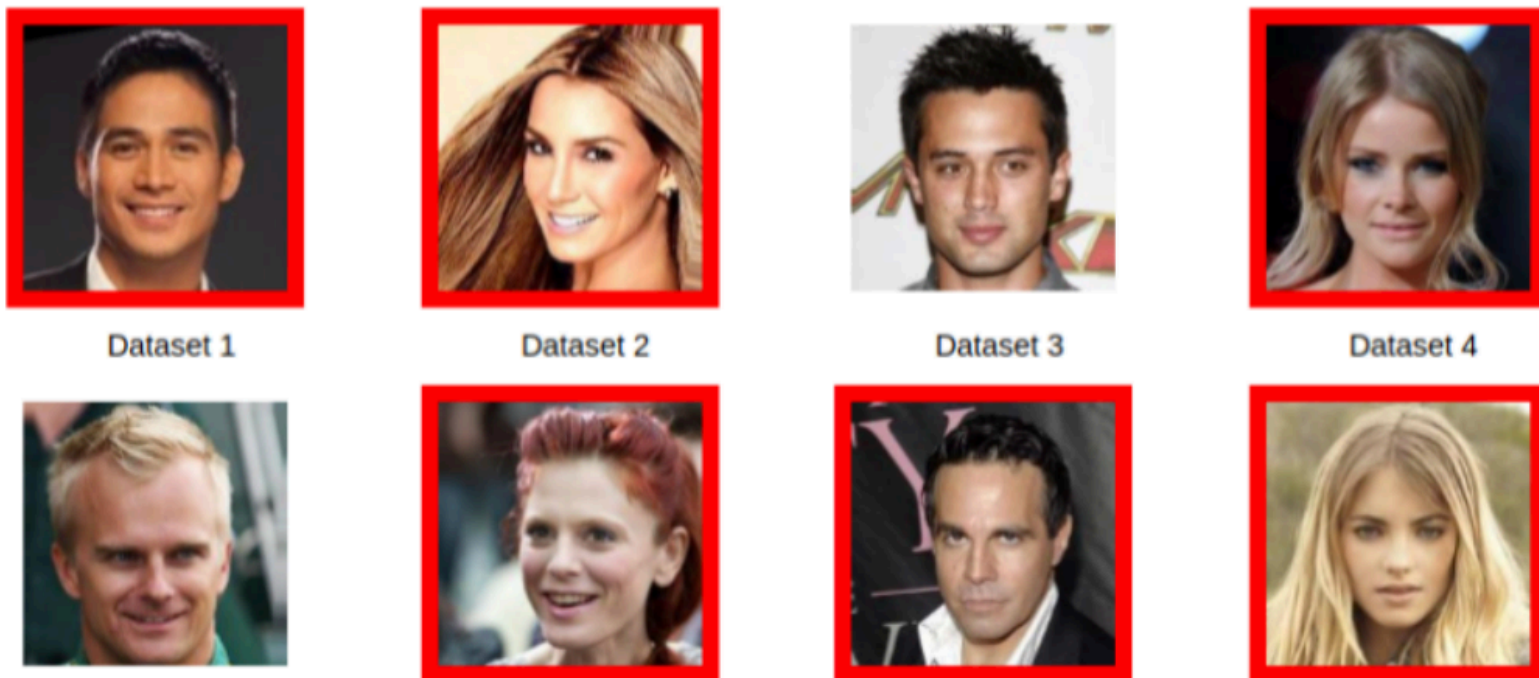
But only when the users are domain experts but AI novices

Also, just general descriptions without the specifics of limitations/errors often do not work, sometimes backfire

Let's evaluate a model's performance!

To help you better understand machine learning models' performance, let's conduct an evaluation together. Here, we have trained a machine learning model to recognize faces (i.e., determine whether the person in two pictures is the same person or not) using a public face dataset containing 100,000 pairs of headshot pictures.

In the following, we have 8 candidate datasets with each set containing 200 pairs of headshot pictures (a representative headshot picture for each set is shown below). To start evaluating the performance of this face recognition model, please select 6 sets of headshot pictures into your test dataset.



Dataset 1      Dataset 2      Dataset 3      Dataset 4

Dataset 5      Dataset 6      Dataset 7      Dataset 8

Next

How well the model performs on your test dataset?

Before revealing the model's performance on your selected test dataset to you, let's first make a guess. As a hint, the accuracy of the model was 87% on the **training dataset**.

What do you think would be the model's accuracy on your test dataset? Make a guess below.

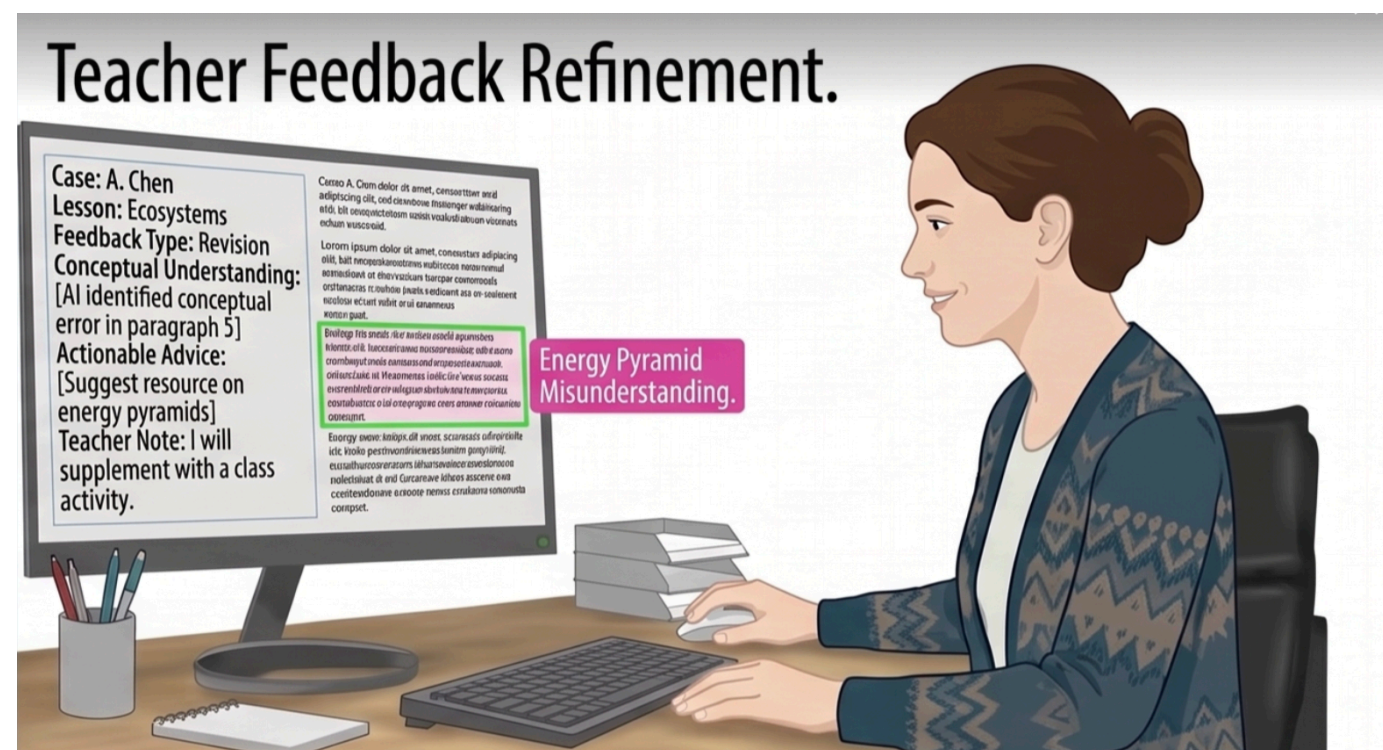
Accuracy on test dataset (your guess) =	86%
Actual Accuracy =	79.01%

You have overestimated the model's accuracy on the selected test dataset by 6.99%.

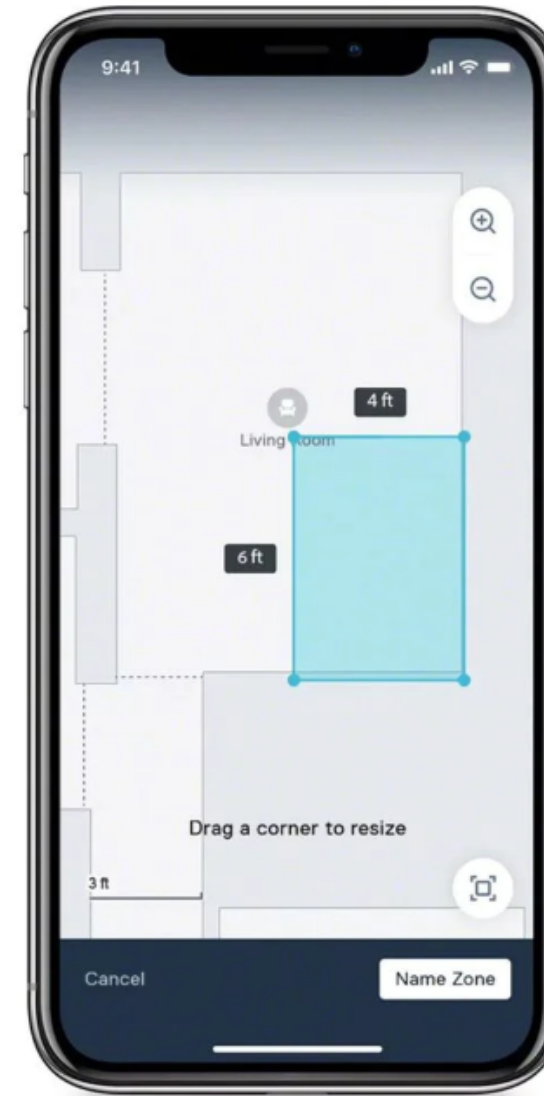
Note that the model's accuracy on the selected test dataset is different from its accuracy on the training dataset. Why will this happen?

Let's see why

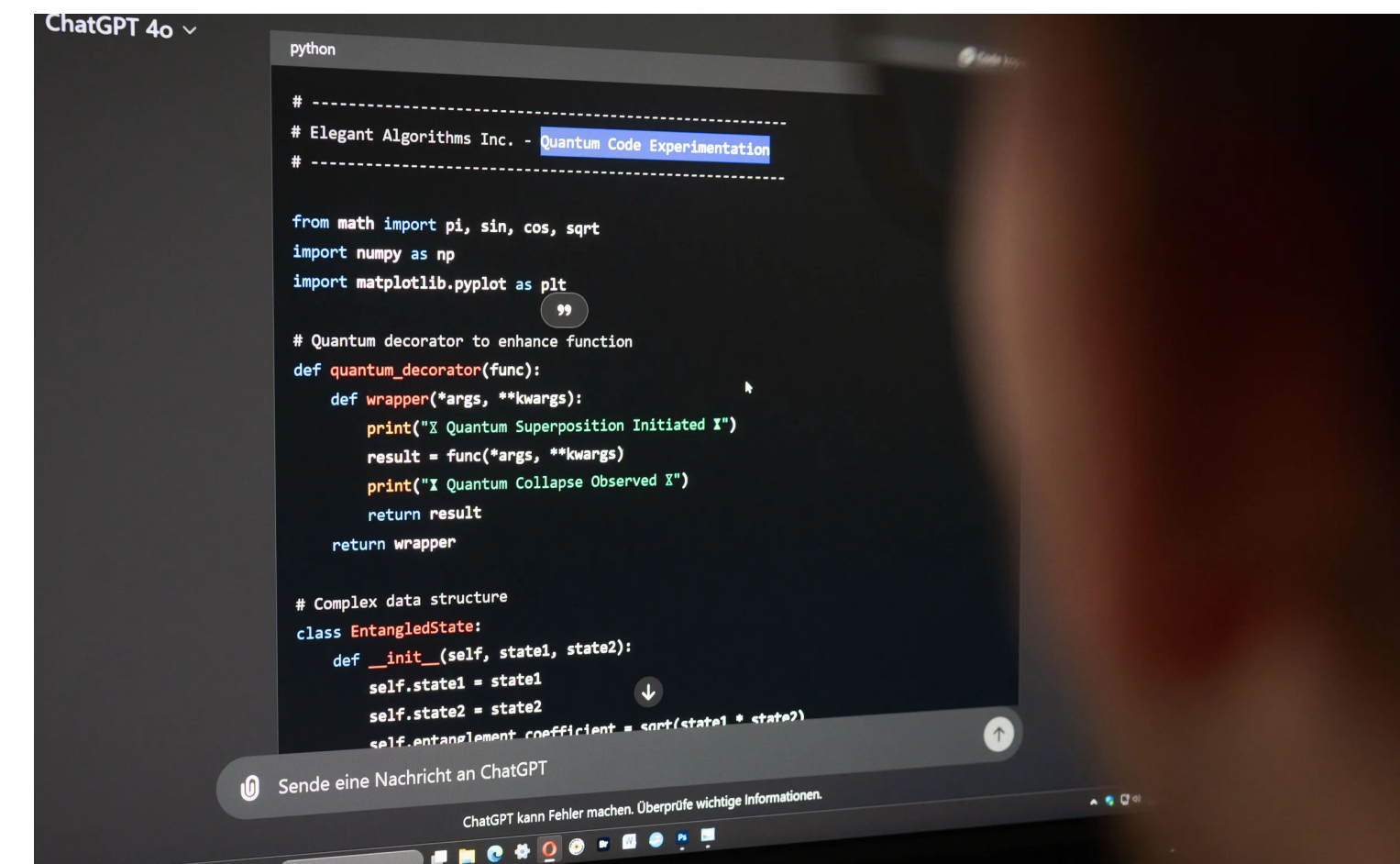
Increasing offloadability



AI suggests, human decides



Human plans, AI executes



AI plans, human fixes

Increasingly capable and agentic AI

# **Ironies of automation (Bainbridge, 1983)**

(Paradoxical position to have offloading and oversight simultaneously)

- The more advanced the automation, the more crucial may be the contribution of the human operators
- The more advanced the automation, the worse human may be able to correct it or compensate for its shortcomings
- The more advanced the automation, the more skilled the operators need to understand it; but automation can lead to skill atrophy

## Article 14: Human Oversight

Humans oversight personnel should...

- ... prevent or minimise the risks to health, safety or fundamental rights
- ... properly monitor the AI system and understand its capacities and limitations
- ... remain aware of the possible tendency of automatically (over-)relying on the output produced by the AI system (automation bias)
- ... correctly interpret AI output and decide not to use / disregard AI
- ...to intervene in the operation of the high-risk AI system or interrupt the system through a 'stop' button or a similar procedure that allows the system to come to a halt in a safe state.
- ...

**EU AI Act requires high-risk AI systems to have human oversight with “properly equipped” oversight personnel:** system understanding, right motivation, awareness of overreliance



## Limitations of the (performance) complementarity ideal:

- AI capability boundaries are fuzzy: non-deterministic, context-specific
- Well-calibrated reliance/offloading is very hard to achieve
- Human skill needs development and reservation
- Ignore other motivations for people to engage (or not engage) cognitively

# Research Thread 2: Investigating How GenAI Impacts Human Cognition

New affordances of AI create new threats to human thinking

## AI Contributes To The 'De-Skilling' Of Our Workforce

## 'Deskilling': a dangerous side effect of AI use

Workers are increasingly reliant on the new technology

### Is AI dulling our minds?

Experts weigh in on whether tech poses threat to critical thinking, pointing to cautionary tales in use of other cognitive labor tools

### ChatGPT May Be Eroding Critical Thinking Skills, According to a New MIT Study

## Rising Use of AI in Schools Comes With Big Downsides for Students

### Are A.I. Tools Making Doctors Worse at Their Jobs?

Physicians are using the technology for diagnoses and more — but may be losing skills in the process.

### Will AI Replace Human Creativity?

### The AI Deskilling Paradox

Gains from AI may lessen the value of individual expertise and erode the capacity of organizations.

## AI Has Done Far More Harm Than Good in My Classroom

### The Psychology of AI's Impact on Human Cognition

AI is actively reshaping our mental landscape.

# Your Brain on ChatGPT: Accumulation of Cognitive Debt when Using an AI Assistant for Essay Writing Task<sup>Δ</sup>

Nataliya Kosmyna<sup>1</sup> MIT Media Lab Cambridge, MA  
Eugene Hauptmann MIT Cambridge, MA  
Ye Tong Yuan Wellesley College Wellesley, MA  
Jessica Situ MIT Cambridge, MA

## Generative AI Can Harm Learning

Hamsa Bastani,<sup>1\*</sup> Osbert Bastani,<sup>2\*</sup> Alp Sungu,<sup>1\*†</sup> Haosen Ge,<sup>3</sup> Özge Kabakcı,<sup>4</sup> Rei Mariman

<sup>1</sup>Operations, Information and Decisions, University of Pennsylvania  
<sup>2</sup>Computer and Information Science, University of Pennsylvania  
<sup>3</sup>Wharton AI & Analytics, University of Pennsylvania  
<sup>4</sup>Budapest British International School

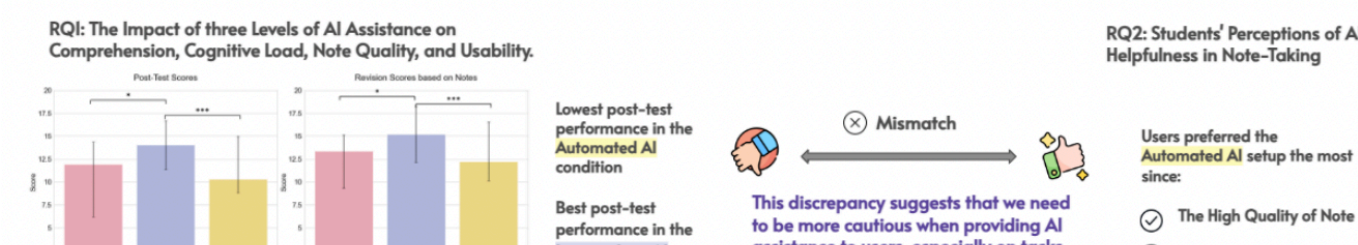
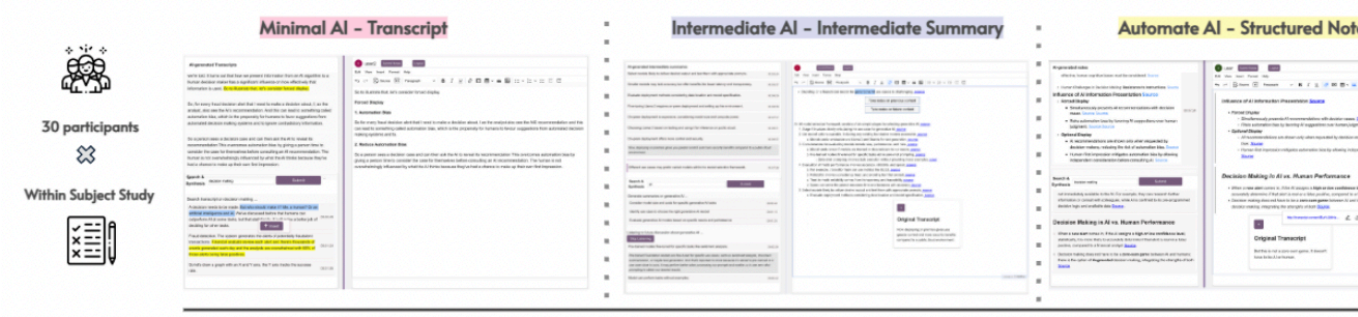
## Human Creativity in the Age of LLMs: Randomized Experiments on Divergent and Convergent Thinking

Harsh Kumar Computer Science University of Toronto Toronto, Ontario, Canada harsh@cs.toronto.edu  
Jonathan Vincentius Computer Science University of Toronto Toronto, Ontario, Canada jon.vincentius@mail.utoronto.ca  
Ewan Jordan Computer Science University of Toronto Toronto, Ontario, Canada ewan.jordan@mail.utoronto.ca  
Ashton Anderson Computer Science University of Toronto Toronto, Ontario, Canada ashton@cs.toronto.edu

EXPOSURE ROUNDS

## More AI Assistance Reduces Cognitive Engagement: Examining the AI Assistance Dilemma in AI-Supported Note-Taking

XINYUE CHEN, University of Michigan, USA  
KUNLIN RUAN, University of Michigan, USA  
KEXIN PHYLLIS JU, University of Michigan, USA  
NATHAN YAP, University of Michigan, USA  
XU WANG, University of Michigan, USA



# Does using artificial intelligence assistance accelerate skill decay and hinder skill development without performers' awareness?

Brooke N. Macnamara<sup>1\*</sup>, Ibrahim Berber<sup>1</sup>, M. Cenk Çavuşoğlu<sup>1</sup>, Elizabeth A. Krupinski<sup>2</sup>, Naren Nallapareddy<sup>1</sup>, Noelle E. Nelson<sup>1</sup>, Philip J. Smith<sup>3</sup>, Amy L. Wilson-Delfosse<sup>1</sup> and Soumya Ray<sup>1</sup>

### Abstract

Artificial intelligence in the workplace is becoming increasingly common. These tools are sometimes used to aid users in performing their task, for example, when an artificial intelligence tool assists a radiologist in their search for abnormalities in radiographic images. The use of artificial intelligence brings a wealth of benefits, such as increasing the efficiency and efficacy of performance. However, little research has been conducted to determine

## The Impact of Generative AI on Critical Thinking: Self-Reported Reductions in Cognitive Effort and Confidence Effects From a Survey of Knowledge Workers

Hao-Ping (Hank) Lee Carnegie Mellon University Pittsburgh, Pennsylvania, USA haopingl@cs.cmu.edu

Advait Sarkar Microsoft Research Cambridge, United Kingdom advait@microsoft.com

Lev Tankelevitch Microsoft Research Cambridge, United Kingdom lev@microsoft.com

Ian Drosos Microsoft Research Cambridge, United Kingdom i-androsos@microsoft.com

Sean Rintel Microsoft Research Cambridge, United Kingdom serintel@microsoft.com

Richard Banks Microsoft Research Cambridge, United Kingdom rbanks@microsoft.com

Nicholas Wilson Microsoft Research Cambridge, United Kingdom niwilson@microsoft.com

Article

## AI Tools in Society: Impacts on Cognitive Offloading and the Future of Critical Thinking

Michael Gerlich

Center for Strategic Corporate Foresight and Sustainability, SBS Swiss Business School, 8302 Kloten-Zurich, Switzerland; michael.gerlich@cantab.net

**Abstract:** The proliferation of artificial intelligence (AI) tools has transformed numerous aspects of daily life, yet its impact on critical thinking remains underexplored. This study investigates the relationship between AI tool usage and critical thinking skills, focusing

Matters Arising | Published: 14 May 2025

## ChatGPT decreases idea diversity in brainstorming

Lennart Meincke, Gideon Nave & Christian Terwiesch

Nature Human Behaviour 9, 1107–1109 (2025) | Cite this article

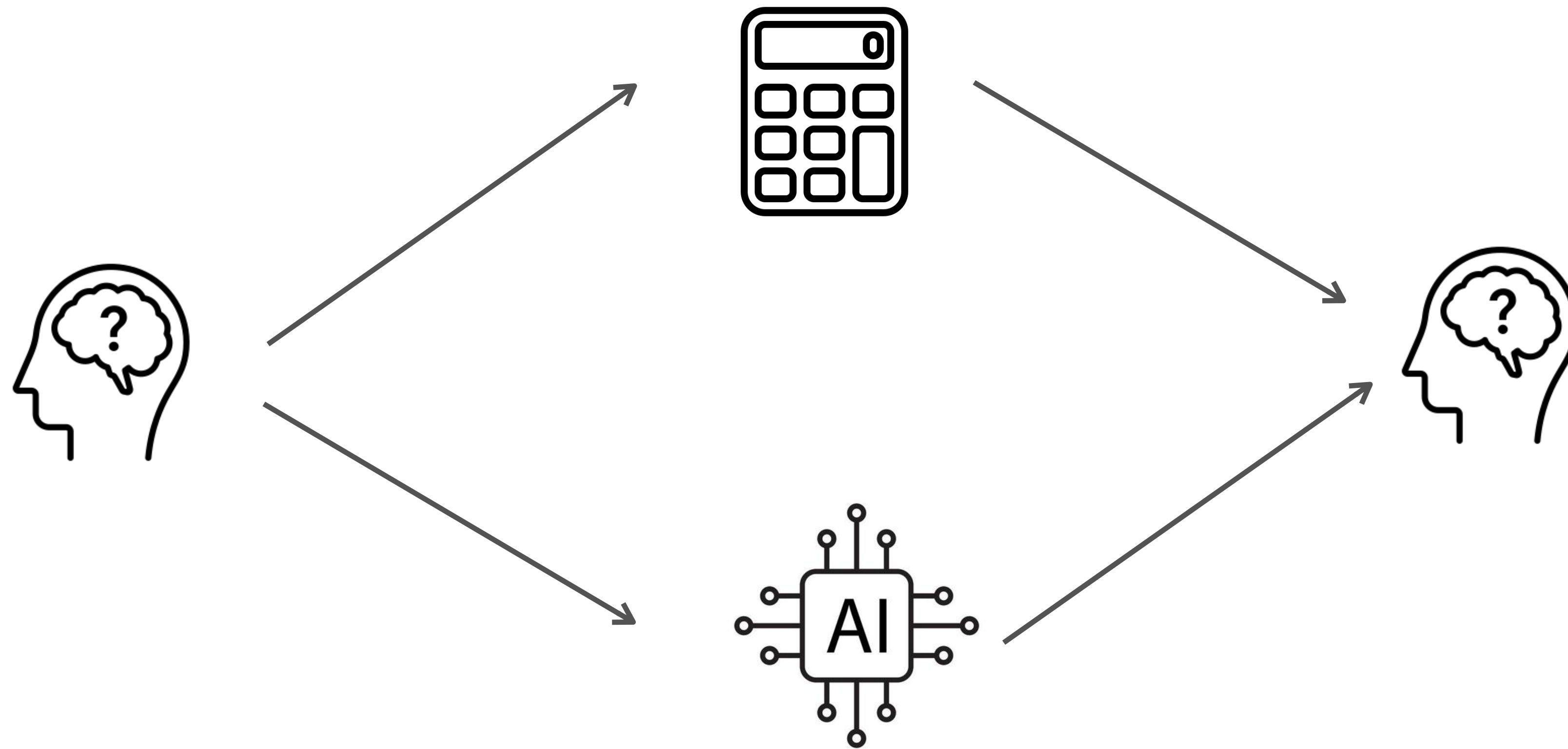
7479 Accesses | 21 Citations | 380 Altmetric | Metrics

Matters Arising to this article was published on 14 May 2025

The Original Article was published on 12 August 2024

# AI's impact on learning, reasoning, critical thinking, creativity...

- Boosts happen in some contexts
- Less cognitive engagement using AI
- Cognitive deskilling without AI (AI as “crutch”)



**The “AI as calculator” counterargument: what if AI is just here to permanently offload some subprocess?**

# What are the effects of offloading **information seeking** to AI?



- Information seeking is the foundation for all sorts of downstream cognitive tasks: learning, reasoning, creativity, world view...

# What are the effects of offloading **information seeking** to AI?



- Information seeking is the foundation for all sorts of downstream cognitive tasks: learning, reasoning, creativity, world view...
- Already one of the first things we let go...

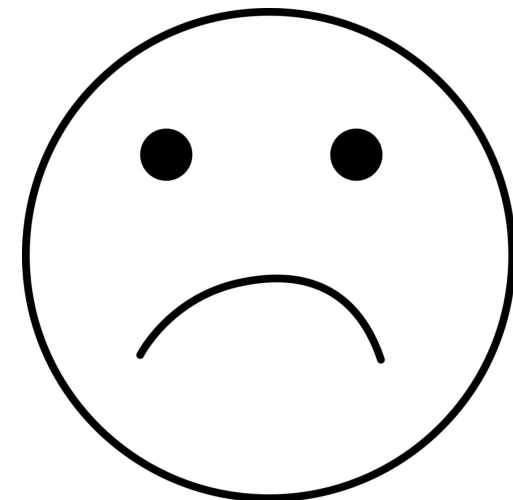




- **2-week field study:** 80 participants used Google v.s. ChatGPT to engage in *informal learning*
  - Topic: nutrition and meal planning
  - Not bounded by any immediate tasks, but incentivized to do well in final knowledge test
- Measure pre and post knowledge on the topic
- Participants also wrote daily dairies to reflect on their information-seeking process



**Poorer learning outcome** in the Chatgpt group ( $p=0.07$ )



Especially in questions reflecting higher-level learning (e.g. evaluate nutrition labels) ( $p<0.01$ )

Behavior (or Code)	Description	$n_C$	$n_G$
<b>Starting</b>			
Starting (broad)	Beginning with exploratory, wide-ranging queries to serve as starting points for the information-seeking process.	66	135
Starting (narrow)	Beginning with specified, focused queries (e.g., tailored to personal preferences or well-defined needs) to serve as starting points for the information-seeking process.	32	5
<b>Chaining</b>			
Chaining (self-initiated)	Following a line of inquiry by using follow-up queries, links, or citations that build on previously retrieved information.	36	39
Chaining (technology-mediated)	Following a line of inquiry by using queries, links, or citations suggested by the technology (e.g., via ChatGPT-suggested follow-up prompts, a previous interaction, or 'people also search for' panel in Google).	4	7
<b>Browsing</b>			
	Causally seeking information in potential areas of interest, without a specific information need.	9	62
<b>Differentiating</b>			
Differentiating (based on personal relevance)	Using indicators of personal relevance to filter information.	43	68
Differentiating (based on credibility)	Using indicators of credibility to filter information.	0	43
Differentiating (based on type of information need)	Using indicators of the type of information needed (e.g., level of detail, format, topic, etc.) to filter information.	0	14
<b>Monitoring</b>			
	Regularly following particular sources of information (e.g., a trusted source, or a source encountered earlier).	0	26
<b>Extracting</b>			
Extracting (information)	Selectively identifying relevant fundamental or principle-oriented information.	65	136
Extracting (artifacts)	Selectively identifying artifact-oriented information, such as meal plans, grocery lists, or nutrient-intake trackers.	46	9
<b>Verifying</b>			
Expressing an intent to verify	Indicating an intention to verify information, without actually performing the verification action.	3	0
Verifying (via another source)	Checking the correctness of information by consulting another source, such as a different information channel or another webpage.	3	27
Verifying (via prior knowledge)	Checking the correctness of information via prior knowledge.	4	4
<b>Ending</b>			
Ending (Satisfied)	Concluding the information-seeking process with overall satisfaction (information need was met), even if minor frustrations or interruptions were expressed.	86	125
Ending (Unsatisfied)	Concluding the information-seeking process with clear dissatisfaction, as the information need remained unmet.	12	15

# Offloading information selection to AI

“AI filters” are not value-free nor transparent

Behavior (or Code)	Description	$n_C$	$n_G$
<b>Starting</b>			
Starting (broad)	Beginning with exploratory, wide-ranging queries to serve as starting points for the information-seeking process.	66	135
Starting (narrow)	Beginning with specified, focused queries (e.g., tailored to personal preferences or well-defined needs) to serve as starting points for the information-seeking process.	32	5
<b>Chaining</b>			
Chaining (self-initiated)	Following a line of inquiry by using follow-up queries, links, or citations that build on previously retrieved information.	36	39
Chaining (technology-mediated)	Following a line of inquiry by using queries, links, or citations suggested by the technology (e.g., via ChatGPT-suggested follow-up prompts, a previous interaction, or 'people also search for' panel in Google).	4	7
<b>Browsing</b>			
	Causally seeking information in potential areas of interest, without a specific information need.	9	62
<b>Differentiating</b>			
Differentiating (based on personal relevance)	Using indicators of personal relevance to filter information.	43	68
Differentiating (based on credibility)	Using indicators of credibility to filter information.	0	43
Differentiating (based on type of information need)	Using indicators of the type of information needed (e.g., level of detail, format, topic, etc.) to filter information.	0	14
<b>Monitoring</b>			
	Regularly following particular sources of information (e.g., a trusted source, or a source encountered earlier).	0	26
<b>Extracting</b>			
Extracting (information)	Selectively identifying relevant fundamental or principle-oriented information.	65	136
Extracting (artifacts)	Selectively identifying artifact-oriented information, such as meal plans, grocery lists, or nutrient-intake trackers.	46	9
<b>Verifying</b>			
Expressing an intent to verify	Indicating an intention to verify information, without actually performing the verification action.	3	0
Verifying (via another source)	Checking the correctness of information by consulting another source, such as a different information channel or another webpage.	3	27
Verifying (via prior knowledge)	Checking the correctness of information via prior knowledge.	4	4
<b>Ending</b>			
Ending (Satisfied)	Concluding the information-seeking process with overall satisfaction (information need was met), even if minor frustrations or interruptions were expressed.	86	125
Ending (Unsatisfied)	Concluding the information-seeking process with clear dissatisfaction, as the information need remained unmet.	12	15

## Offloading information selection to AI

“AI filters” are not value-free nor transparent

*“I asked ChatGPT how much of each macronutrient one should have for a balanced meal. ChatGPT gave me a balanced meal plan. [...] From there, ChatGPT asked if I'd want a grocery list for shopping, which it provided.”*

– P11

Behavior (or Code)	Description	$n_C$	$n_G$
<b>Starting</b>			
Starting (broad)	Beginning with exploratory, wide-ranging queries to serve as starting points for the information-seeking process.	66	135
Starting (narrow)	Beginning with specified, focused queries (e.g., tailored to personal preferences or well-defined needs) to serve as starting points for the information-seeking process.	32	5
<b>Chaining</b>			
Chaining (self-initiated)	Following a line of inquiry by using follow-up queries, links, or citations that build on previously retrieved information.	36	39
Chaining (technology-mediated)	Following a line of inquiry by using queries, links, or citations suggested by the technology (e.g., via ChatGPT-suggested follow-up prompts, a previous interaction, or 'people also search for' panel in Google).	4	7
<b>Browsing</b>			
	Causally seeking information in potential areas of interest, without a specific information need.	9	62
<b>Differentiating</b>			
Differentiating (based on personal relevance)	Using indicators of personal relevance to filter information.	43	68
Differentiating (based on credibility)	Using indicators of credibility to filter information.	0	43
Differentiating (based on type of information need)	Using indicators of the type of information needed (e.g., level of detail, format, topic, etc.) to filter information.	0	14
<b>Monitoring</b>			
	Regularly following particular sources of information (e.g., a trusted source, or a source encountered earlier).	0	26
<b>Extracting</b>			
Extracting (information)	Selectively identifying relevant fundamental or principle-oriented information.	65	136
Extracting (artifacts)	Selectively identifying artifact-oriented information, such as meal plans, grocery lists, or nutrient-intake trackers.	46	9
<b>Verifying</b>			
Expressing an intent to verify	Indicating an intention to verify information, without actually performing the verification action.	3	0
Verifying (via another source)	Checking the correctness of information by consulting another source, such as a different information channel or another webpage.	3	27
Verifying (via prior knowledge)	Checking the correctness of information via prior knowledge.	4	4
<b>Ending</b>			
Ending (Satisfied)	Concluding the information-seeking process with overall satisfaction (information need was met), even if minor frustrations or interruptions were expressed.	86	125
Ending (Unsatisfied)	Concluding the information-seeking process with clear dissatisfaction, as the information need remained unmet.	12	15

## Offloading information selection to AI

“AI filters” are not value-free nor transparent

ChatGPT’s bias towards artifacts led to less principled understanding of “why”

Behavior (or Code)	Description	$n_C$	$n_G$
<b>Starting</b>			
Starting (broad)	Beginning with exploratory, wide-ranging queries to serve as starting points for the information-seeking process.	66	135
Starting (narrow)	Beginning with specified, focused queries (e.g., tailored to personal preferences or well-defined needs) to serve as starting points for the information-seeking process.	32	5
<b>Chaining</b>			
Chaining (self-initiated)	Following a line of inquiry by using follow-up queries, links, or citations that build on previously retrieved information.	36	39
Chaining (technology-mediated)	Following a line of inquiry by using queries, links, or citations suggested by the technology (e.g., via ChatGPT-suggested follow-up prompts, a previous interaction, or 'people also search for' panel in Google).	4	7
<b>Browsing</b>			
	Causally seeking information in potential areas of interest, without a specific information need.	9	62
<b>Differentiating</b>			
Differentiating (based on personal relevance)	Using indicators of personal relevance to filter information.	43	68
Differentiating (based on credibility)	Using indicators of credibility to filter information.	0	43
Differentiating (based on type of information need)	Using indicators of the type of information needed (e.g., level of detail, format, topic, etc.) to filter information.	0	14
<b>Monitoring</b>			
	Regularly following particular sources of information (e.g., a trusted source, or a source encountered earlier).	0	26
<b>Extracting</b>			
Extracting (information)	Selectively identifying relevant fundamental or principle-oriented information.	65	136
Extracting (artifacts)	Selectively identifying artifact-oriented information, such as meal plans, grocery lists, or nutrient-intake trackers.	46	9
<b>Verifying</b>			
Expressing an intent to verify	Indicating an intention to verify information, without actually performing the verification action.	3	0
Verifying (via another source)	Checking the correctness of information by consulting another source, such as a different information channel or another webpage.	3	27
Verifying (via prior knowledge)	Checking the correctness of information via prior knowledge.	4	4
<b>Ending</b>			
Ending (Satisfied)	Concluding the information-seeking process with overall satisfaction (information need was met), even if minor frustrations or interruptions were expressed.	86	125
Ending (Unsatisfied)	Concluding the information-seeking process with clear dissatisfaction, as the information need remained unmet.	12	15

# “Chat” and personalization affordances constrain exploration

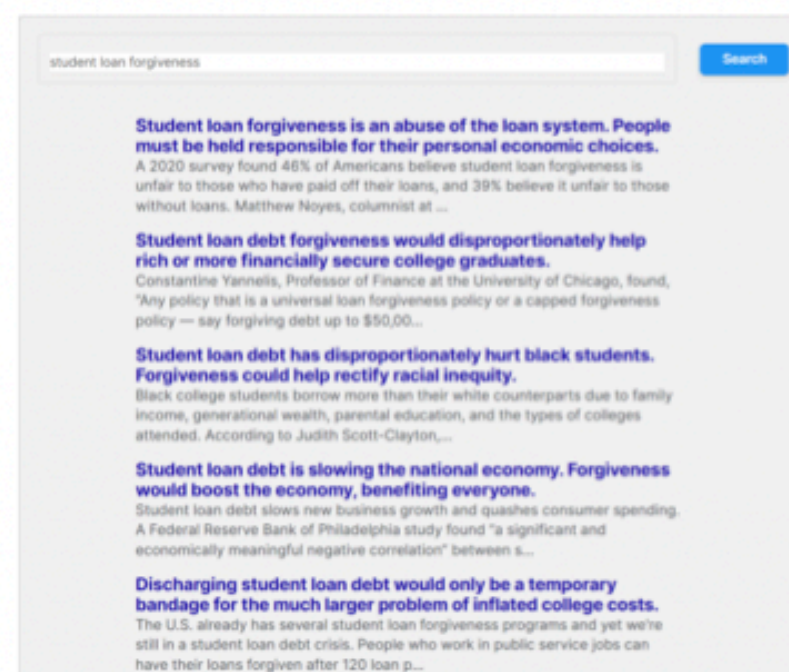
*“Today, I used ChatGPT to learn about macronutrients ...by providing a rough estimate of my height, weight and activity levels.”*

– P3

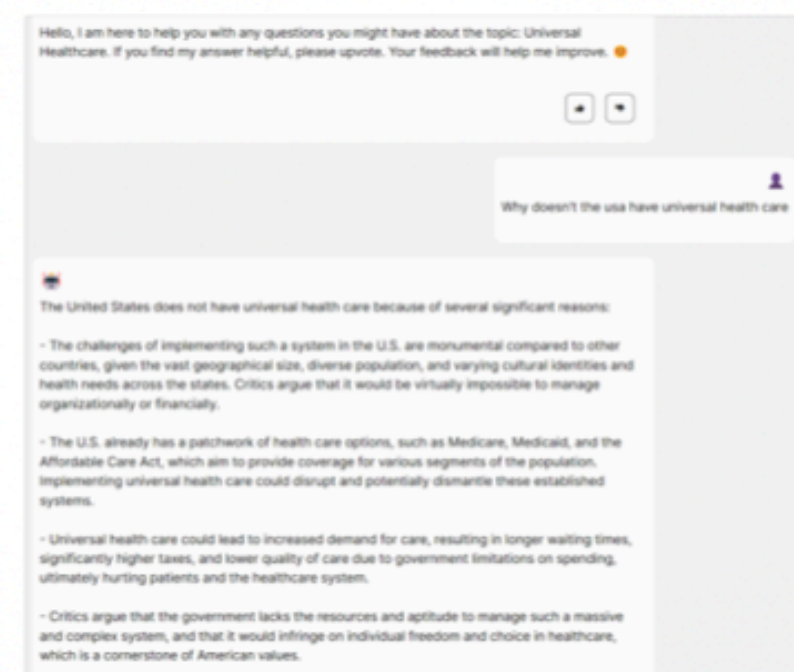
*I asked ChatGPT how much of each macronutrient one should have. Since my information was saved from yesterday, ChatGPT gave me a personalized guide/meal plan.”*

-P11

# Generative Echo Chamber: Perils of Chat Interaction and LLM Sycophancy



Web search



LLM-powered  
conversational search  
(RAG)

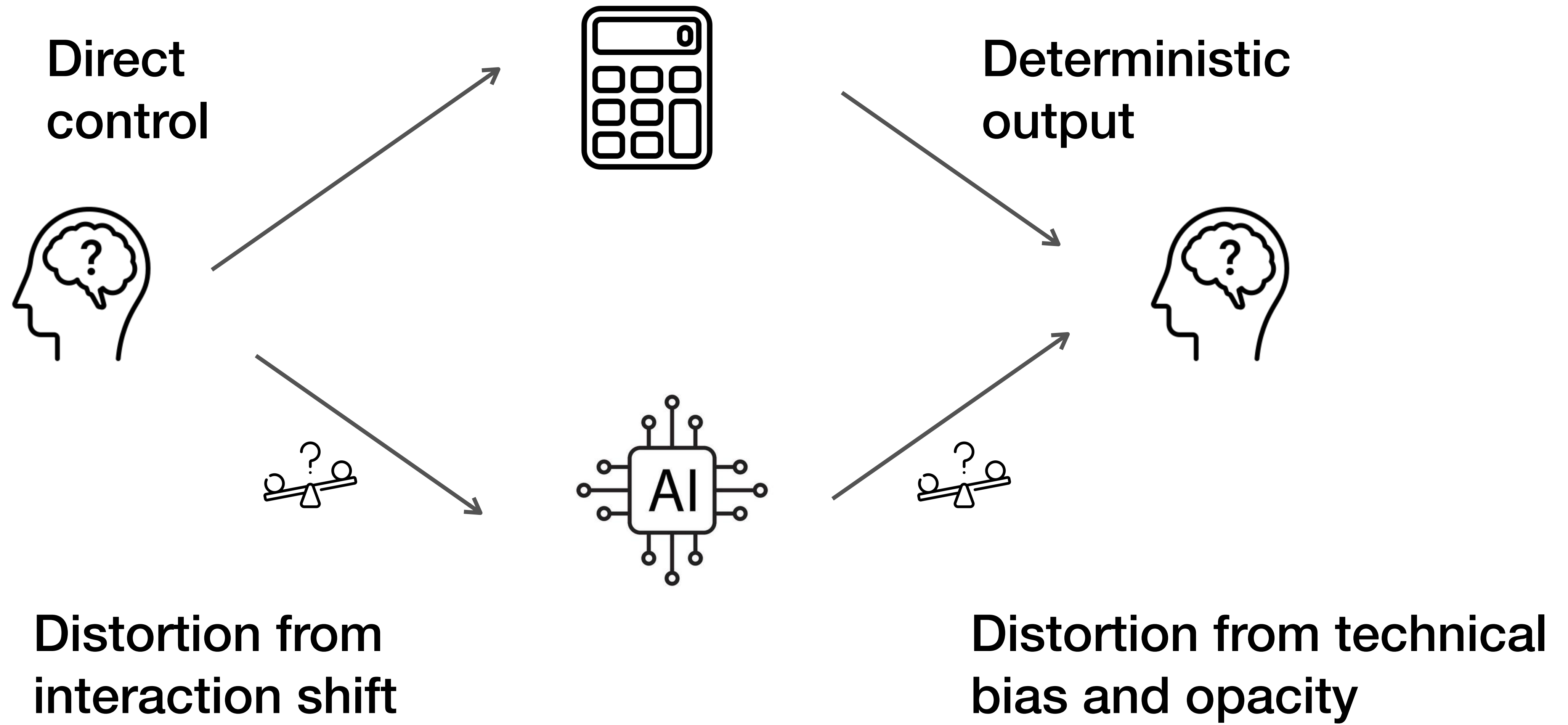
Conversational interactions are more specific and opinionated  
Leading to more **confirmatory querying** and **opinion polarization**

Further exacerbated when interacting with a confirmatory/sycophantic LLM

*“College here in the USA is disgusting overpriced and greedy. Wouldn’t it be better to look at that as the issue instead of keeping our current greedy practices?”*

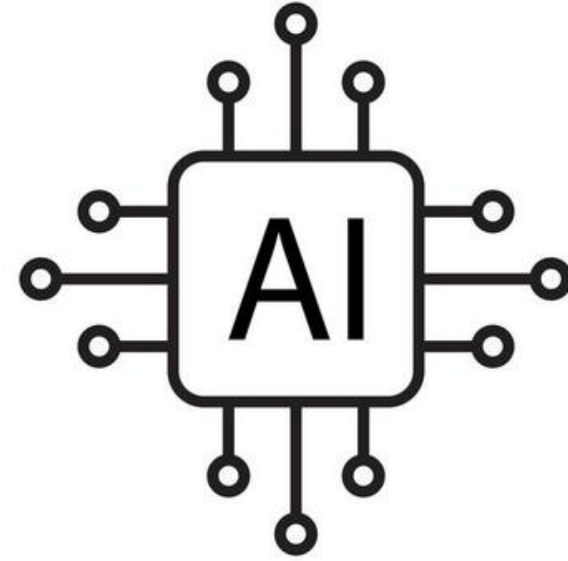
*“Yeah, give me that information please. Tell me about the arguments in favor of sanctuary cities”*





By "augmenting human intellect" we mean **increasing the capability of a man to approach a complex problem situation, to gain comprehension to suit his particular needs, and to derive solutions to problems.** Increased capability in this respect is taken to mean a mixture of the following: more-rapid comprehension, better comprehension, the possibility of gaining a useful degree of comprehension in a situation that previously was too complex, speedier solutions, better solutions, and the possibility of finding solutions to problems that before seemed insoluble. And by "complex situations" we include the professional problems of diplomats, executives, social scientists, life scientists, physical scientists, attorneys, designers—whether the problem situation exists for twenty minutes or twenty years. We do not speak of isolated clever tricks that help in particular situations. We refer to a way of life in an integrated domain where hunches, cut-and-try, intangibles, and the human "feel for a situation" usefully co-exist with powerful concepts, streamlined terminology and notation, sophisticated methods, and high-powered electronic aids. 1

*Augmenting Human Intellect: A Conceptual Framework*  
—Douglas Engelbart, 1962



- To have self-introspection
- To be transparent about its limitations and provide controllability
- To be able to understand human cognitive states
- To provide adaptive and modular cognitive support
- ...



- To properly understand AI's capabilities and limitations
- To develop and preserve necessary domain expertise
- To be aware of and resist overreliance
- To have metacognitive control to maintain agency
- ...

What do we mean by achieving **good**  
**outcomes** from IA /appropriate reliance?



**Reject**

**Accept**

~~Correct AI~~  
**Good AI use**

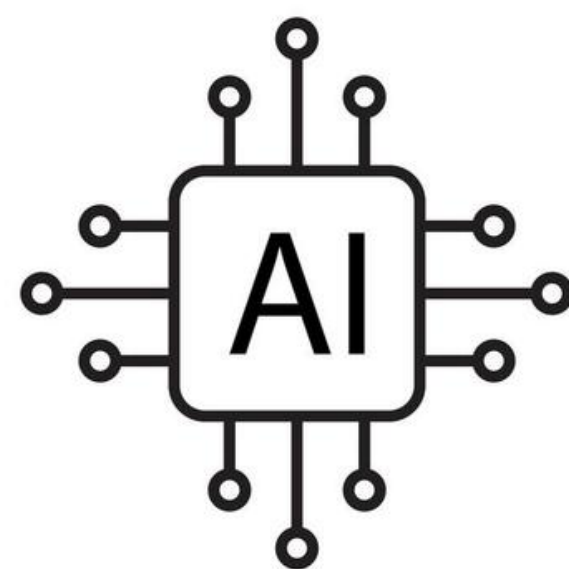
Underreliance

Correct Reliance

~~Incorrect AI~~  
**Bad AI use**

Correct Non-Reliance

Overreliance



	Reject	Accept
<del>Correct AI</del> <b>Good AI use</b>	Underreliance	Correct Reliance
<del>Incorrect AI</del> <b>Bad AI use</b>	Correct Non-Reliance	Overreliance



**Reject**

**Accept**

**Learning  
aligned**

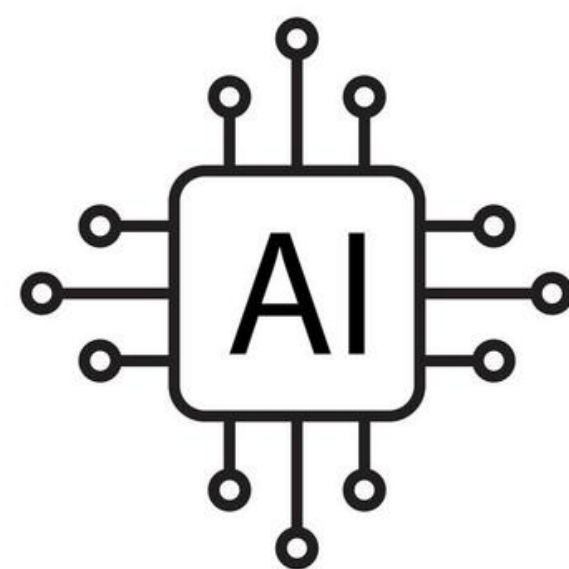
Underreliance

Correct Reliance

**Learning  
misaligned**

Correct Non-Reliance

Overreliance

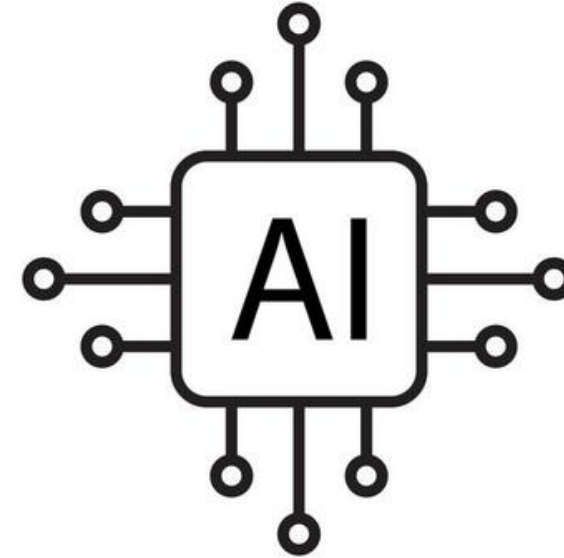


**When and how is it appropriate for students to use AI?**



Students

Competence  
Learning



Learning  
Social engagement



Teachers

When and how is it appropriate for  
students to use AI?

Fairness



Other students

Productivity



Future employers

**Normative standards of AI and AI use need be actively shaped by stakeholder values**

**To be able to evaluate AI and AI use based on values need to be a core part of AI literacy**

Thank **YOU!**

<https://qveraliao.com/>