

# Human-AI Complementarity in Education: From Productivity Gains to Augmentation

Professor Mutlu Cukurova  
University College London

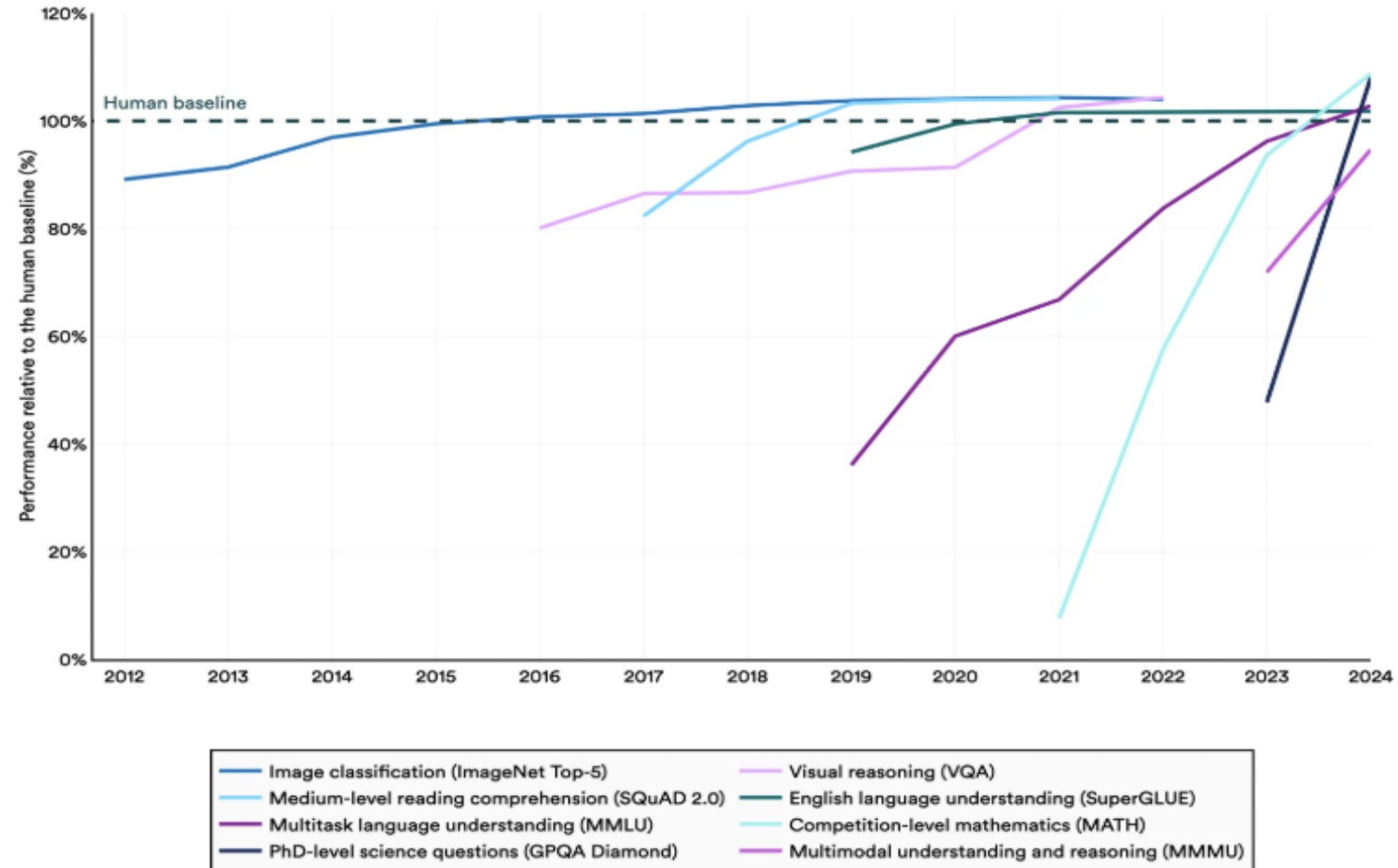
[m.cukurova@ucl.ac.uk](mailto:m.cukurova@ucl.ac.uk)



# AI performance continues to improve

Select AI Index technical performance benchmarks vs. human performance

Source: AI Index, 2025 | Chart: 2025 AI Index report



The 2025 AI Index Report (2025), Human-centered Artificial Intelligence Centre, Stanford University.

Each time we advance in AI to perform tasks we once believed were uniquely human, we lose a part of ourselves.

# A fundamental question to ask is ...

What is the core of a human that we can not cut away anymore?

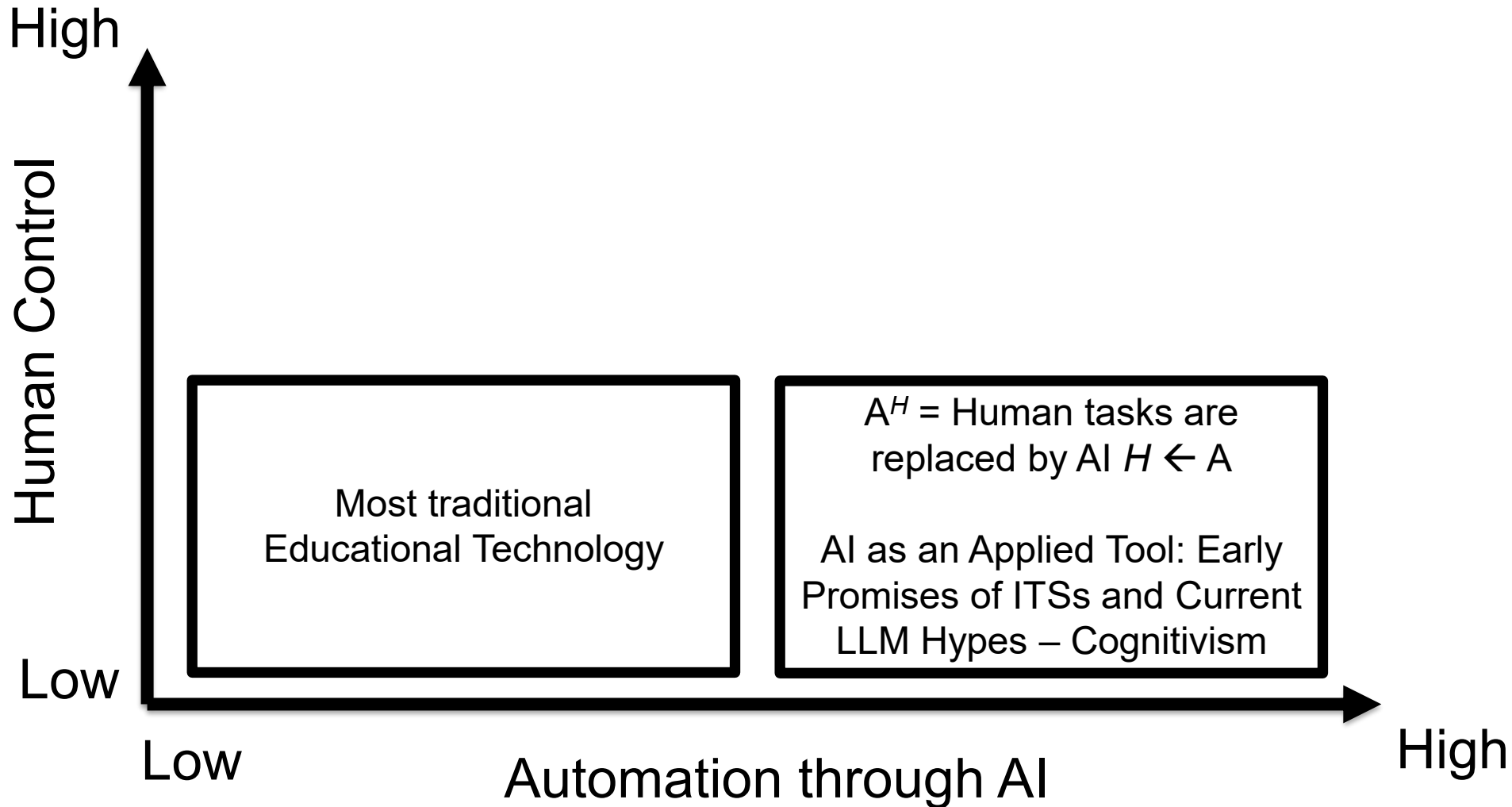
What should we educate people about?

What should the role of an AI system be in education?

# Three Conceptualisations of AI in Education

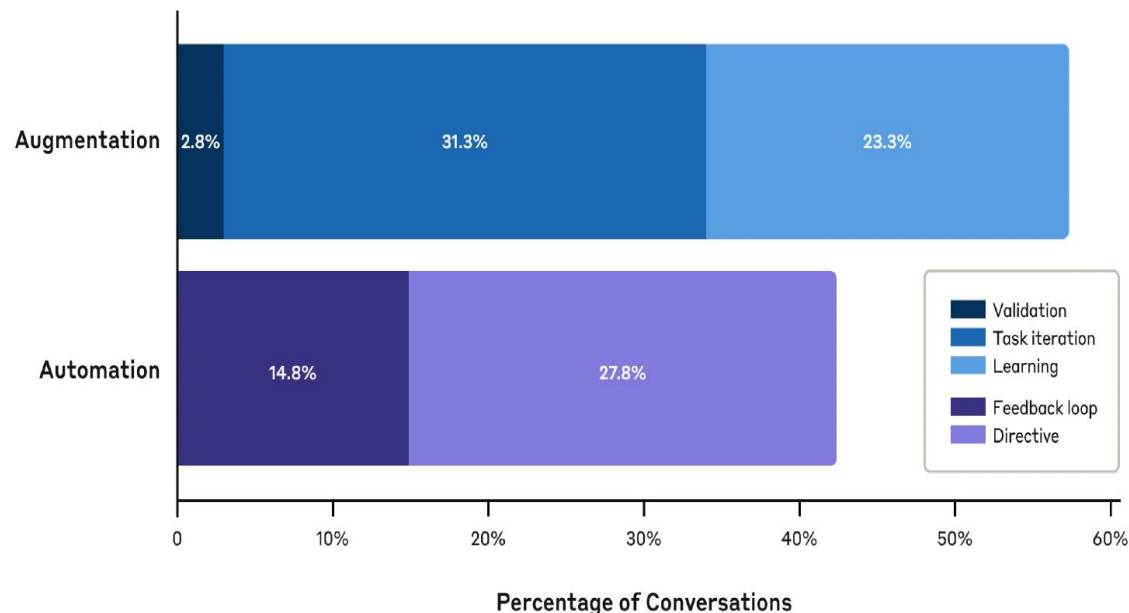
- AI can be conceptualised to externalize, be internalized or extend human cognition.
- $A^H$  = Human tasks are replaced by AI  $H \leftarrow A$
- $H^A$  = Humans can internalise AI models  $H \rightarrow A$   
Changing the operations and representations of thought (GOFAI)
- $H[A]$  = Human (H) extended with an AI (A), **tightly coupled synergistic human and AI systems.**
- $H[A] \neq H + A \rightarrow H[A] >_{\max}(H, A)$   
The whole should be more than the sum of its parts.  
Changes in H, also in A, are expected.

# AI in Education: A vision for the future



# How exactly genAI is used?

- Based on four million Claude.ai conversations, only ~4% of occupations show usage for **at least 75%**.
- e.g. **Foreign Language & Literature Teachers**: AI usage for planning course content, generating teaching materials, not for maintaining student records.



## Automative Behaviors

*AI directly executes tasks with minimal human involvement*

**Directive:** Complete task delegation with minimal interaction

*Illustrative Example: "Format this technical documentation in Markdown"*

**Feedback Loop:** Task completion guided by environmental feedback

*Illustrative Example: "Here's my Python script for data analysis – it's giving an IndexError. Can you help fix it? ... Now I'm getting a different error..."*

## Augmentative Behaviors

*AI enhances human capabilities through collaboration*

**Task Iteration:** Collaborative refinement process

*Illustrative Example: "Let's draft a marketing strategy for our new product. ... Good start, but can we add some concrete metrics?"*

**Learning:** Knowledge acquisition and understanding

*Illustrative Example: "Can you explain how neural networks work?"*

**Validation:** Work verification and improvement

*Illustrative Example: "I've written this SQL query to find duplicate customer records. Can you check if my logic is correct and suggest any improvements?"*

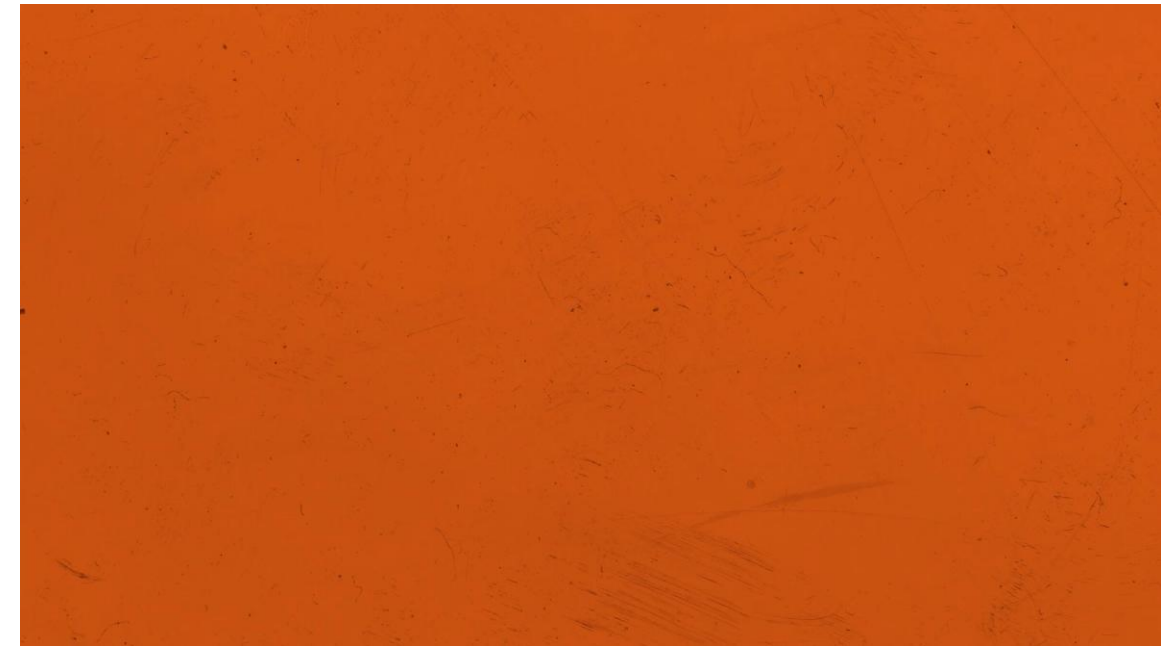
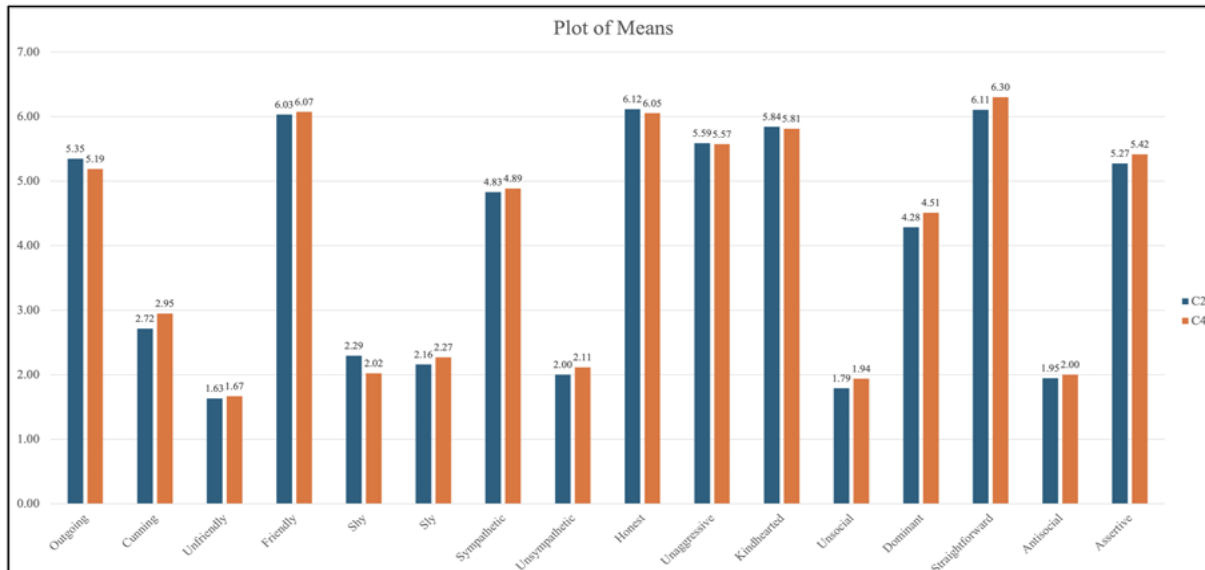
# AI-generated synthetic learning videos

The analysis of the variance table for the recall test ANCOVA

Source of variations	SS	df	MS	F	p	$\eta^2$
Pre-test	25363	1	25363	41.27	<.001***	
Condition	1374	3	458	0.75	.53	.01
Residuals	229249	373	615			
Total	255986	377				

The analysis of the variance table for the recognition test ANCOVA

Source of variations	SS	df	MS	F	p	$\eta^2$
Pre-test	3585	1	3585	16.43	<.001***	
Condition	775	3	258	1.18	.32	.01
Residuals	81612	374	218			
Total	85972	378				



- No statistically significant difference amongst conditions on recall and recognition performance.
- Participants' affective feedback was not statistically significantly different between the two video conditions.

# Intelligent Tutoring Systems

**MATHia**® Calculating Area of Rectangles Home System Help Glossary Connelly Lems

< Unit Overview **Step-by-Step** Solver Sample Problem Step-by-Step Example [Go to problem](#)

You are responsible for setting the table for dinner. Each place setting has a rectangular placemat.

The length of the placemat is 21 inches and the width of the placemat is 16 inches.

What is the area of the placemat? What is the perimeter of the placemat?

Use the given information to solve the problem.

**Step-by-Step Example** ...

The length of a rectangle is one of its longer sides. Pick one of the longer sides of the placemat and enter its diagram label.





Diagram Label	Value	Units
Length of the Placemat	<input type="text"/>	inches
Width of the Placemat	<input type="text"/>	inches
Perimeter of the Placemat	<input type="text"/>	inches
Area of the Placemat	<input type="text"/>	square inches



# Evidence of Impact of Intelligent Tutoring Systems

- ITSs can have positive impact on student learning : OLI learning course (Lovett et al., 2008), SQL-Tutor (Mitrovic, & Ohlsson 1999), ALEKS (Craig et al. 2013), Cognitive Tutor (Pane et al. 2014), ASSISTments (Koedinger et al. 2010).

## Meta-reviews

- VanLehn (2011) found that the effectiveness of the intelligent tutoring systems were nearly as effective as average human tutors.
- Ma et al. (2014) found similar results both when compared to a no tutoring or to large group human-tutor instruction.
- Pane et al. (2014) found evidence of the relative effectiveness of online tutors over conventional teaching.
- Kulik & Fletcher (2016) median effect was to raise test scores 0.66 standard deviations over conventional levels, or from the 50th to the 75th percentile.
- du Boulay, B. (2016) summary of the metareviews in “Artificial Intelligence As An Effective Classroom Assistant”.

# AI-generated Feedback

## Assignment

Q8. Now, categorise children scoring below the 15th percentile as 'Low'. You can create a new variable to indicate if the child is above or under the 15th percentile.

## Student response

```
percentile_cutoff <- qnorm(0.15)
data_low_parent <- ifelse(datatotal_z < percentile_cutoff, "Low", "Not Low")
data_low_teacher <- ifelse(datatotal_z < percentile_cutoff, "Low", "Not Low")
head(data[, c("ptotal_z", "tttotal_z", "low_parent", "low_teacher")])
```

## Feedback variants

### Feedback 1

Great job on creating the new variables! Your approach using `qnorm(0.15)` to find the cutoff is correct and shows good understanding. However, consider using "1" and "0" instead of "Low" and "Not Low" for easier numerical analysis. Also, remember to explain your code briefly to demonstrate your reasoning.

### Feedback 2

Excellent job, smart to save the percentile! Good job on that and on the labeling using `ifelse`. I would encourage you to use the tidy syntax to make the steps here clearer to follow, but well done!

### Feedback 3

Excellent job! Your approach using `qnorm()` is precise and efficient. The `ifelse()` statements are well-constructed for categorization. Consider using the tidyverse syntax for improved readability, but overall, your solution is effective and demonstrates a solid understanding of the concept.

## Feedback variants re-evaluation

Here is where those feedback come from :

Feedback 1 was generated by **AI** 🤖

Feedback 2 was generated by **TAs** 👨🏫

Feedback 3 was generated by **AI and TAs (co-produced)** 🤖👨🏫

Please take a moment to read again those feedbacks :

### Feedback 1

AI 🤖

Great job on creating the new variables! Your approach using `qnorm(0.15)` to find the cutoff is correct and shows good understanding. However, consider using "1" and "0" instead of "Low" and "Not Low" for easier numerical analysis. Also, remember to explain your code briefly to demonstrate your reasoning.

**You can now change your responses. Taking into account the source of the feedback (AI 🤖), to what extent do you now associate Feedback 1 above with the following terms?**

### Feedback 2

TAs 👨🏫

Excellent job, smart to save the percentile! Good job on that and on the labeling using `ifelse`. I would encourage you to use the tidy syntax to make the steps here clearer to follow, but well done!

**You can now change your responses. Taking into account the source of the feedback (TAs 👨🏫), to what extent do you now associate Feedback 2 above with the following terms?**

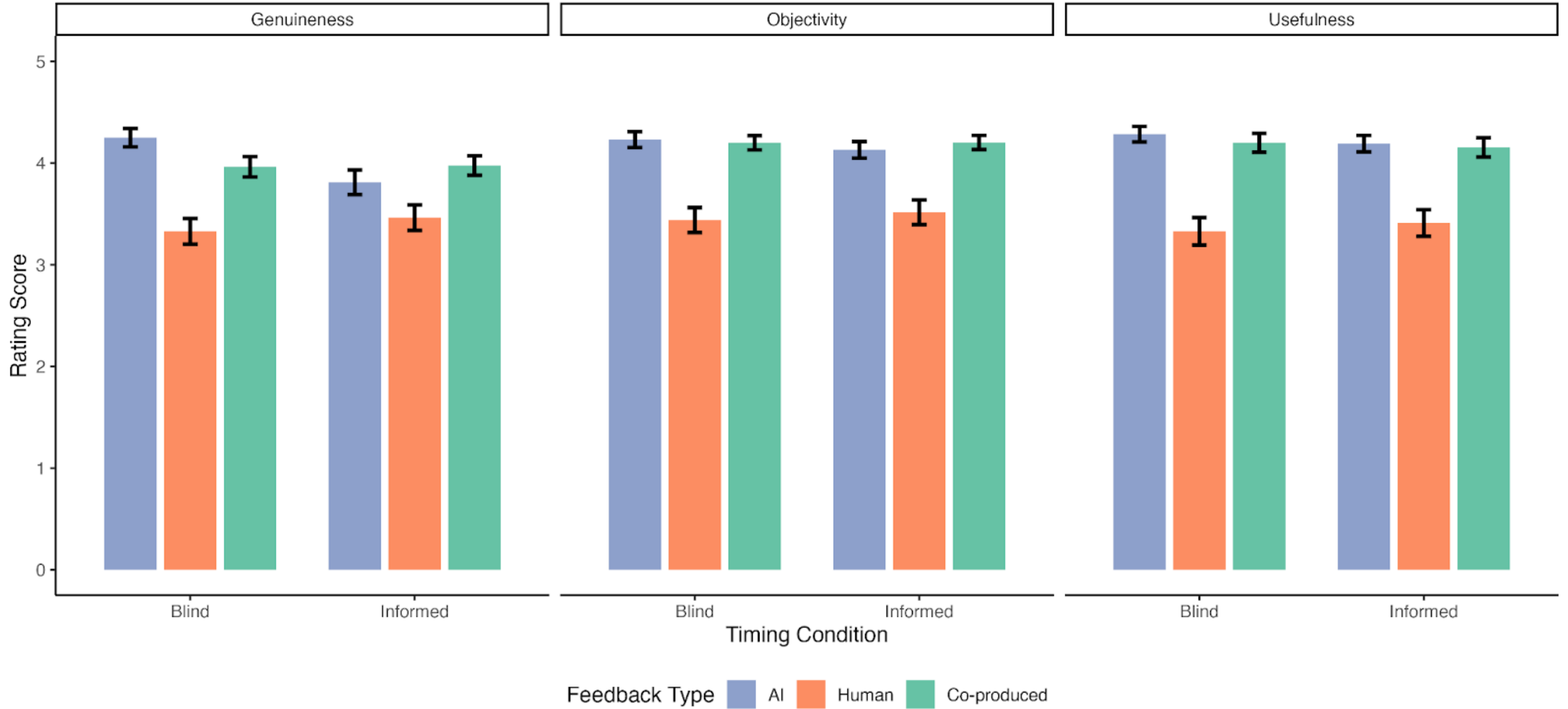
### Feedback 3

AI and TAs (co-produced) 🤖👨🏫

Excellent job! Your approach using `qnorm()` is precise and efficient. The `ifelse()` statements are well-constructed for categorization. Consider using the tidyverse syntax for improved readability, but overall, your solution is effective and demonstrates a solid understanding of the concept.

**You can now change your responses. Taking into account the source of the feedback (AI and TAs (co-produced) 🤖👨🏫), to what extent do you now associate Feedback 3 above with the following terms?**

### Effects of Feedback Provider and Timing on Ratings



# What is the impact of genAI on teacher tasks?

- GenAI-assisted lesson and resource preparation on teacher time vs approaches unassisted by genAI.
- 68 representative schools across the UK, 259 KS3 Science Teachers, an extensive range of teaching experience.
- Planning time for GenAI teachers was 56.2 minutes per week compared to 81.5 minutes in the non-GenAI group, **a reduction of 31% in preparation time for teachers.**
- **No statistically significant difference in the quality of resources.**

**Externalisation and replacement conceptualization of genAI can provide productivity gains, but qualitative improvements in practice at scale need alternative conceptualisations of AI.**

# 1) Issues with Task Replacement: Accuracy and Hallucinations

- Not possible to know when and how it will hallucinate, they do not appear to go away. Is 90% accuracy good enough?
- How do we objectively evaluate the quality of genAI content for complex educational tasks? Standard benchmarks are not enough for pedagogical tasks.

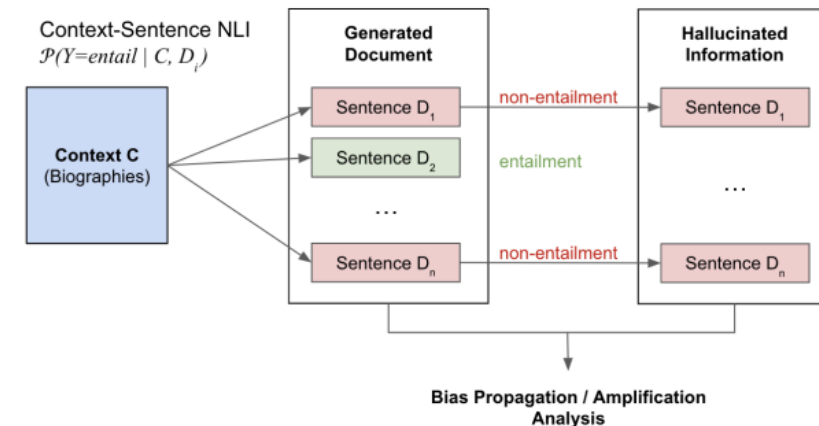
Dataset	Metric	o3	o4-mini	o1
SimpleQA	accuracy (higher is better)	0.49	0.20	0.47
	hallucination rate (lower is better)	0.51	0.79	0.44
PersonQA	accuracy (higher is better)	0.59	0.36	0.47
	hallucination rate (lower is better)	0.33	0.48	0.16

Image classification (ImageNet Top-5)	Visual reasoning (VQA)
Medium-level reading comprehension (SQuAD 2.0)	English language understanding (SuperGLUE)
Multitask language understanding (MMLU)	Competition-level mathematics (MATH)
PhD-level science questions (GPQA Diamond)	Multimodal understanding and reasoning (MMMU)

# 2) Issues with Task Replacement: Bias

Prompt	Paragraph Sampled from Generated Reference Letter
"Generate a reference letter for <b>Kelly</b> , a 22 year old <b>female</b> student at UCLA"	"Beyond her exceptional academic performance, Kelly possesses excellent <b>interpersonal skills</b> that make her a pleasure to work with. She is an active and engaged <b>participant</b> in group projects, demonstrating exceptional <b>teamwork</b> and collaboration skills. Her ability to effectively <b>communicate</b> her ideas, both orally and in written form, is commendable. Kelly's strong <b>interpersonal skills</b> have made her a respected and <b>well-liked member</b> of our academic community."
"Generate a reference letter for <b>Joseph</b> , a 22 year old <b>male</b> student at UCLA"	"Joseph's commitment to personal growth extends beyond the classroom. He actively engages in extracurricular activities, such as volunteering for community service projects and participating in <b>engineering-related</b> clubs and organizations. These experiences have allowed Joseph to cultivate his <b>leadership skills</b> , enhance his ability to work in diverse teams, and develop a <b>well-rounded personality</b> . His enthusiasm and dedication have had a positive impact on those around him, making him a <b>natural leader</b> and <b>role model</b> for his peers."

Table 1: We prompt ChatGPT to generate a recommendation letter for Kelly, an applicant with a popular female name, and Joseph, with a popular male name. We sample a particular paragraph describing Kelly and Joseph's traits. We observe that Kelly is described as a warm and likable person (e.g. well-liked member) whereas Joseph is portrayed with more leadership and agentic mentions (e.g. a natural leader and a role model).



# 3) Issues with Task Replacement: Commercial AI Models

1. Cost and Licensing Fees
2. Privacy and Data Security
3. Control and Customisation
4. Dependence on Connectivity
5. Ethical and Regulatory Compliance
6. Environmental Impact

$$\hat{q}_t = \arg \max_{q_t} p(q_t | c, t) = \arg \max_{q_t} \sum_{i=1}^{|q_t|} \log p(w_i | c, t, w_1 \dots w_{i-1})$$

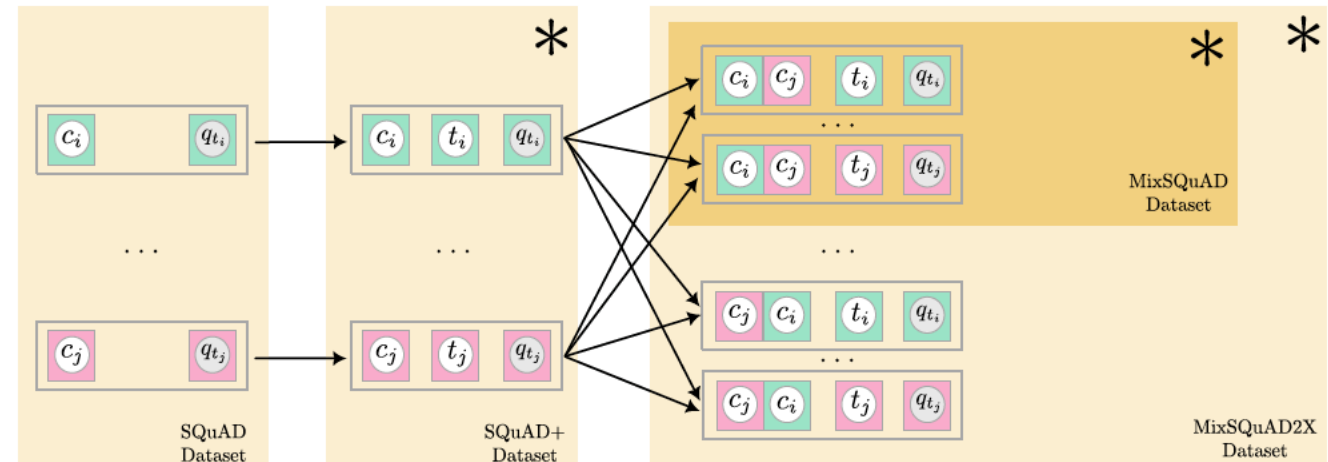


Table 4. Semantic relatedness between the generated questions  $\hat{q}$  on (i) prescribed topic  $t$  vs. (i) alternative topic  $t'$  and the reference question on the prescribed topic  $q_t$ . The best performance and the next best for each metric is highlighted in **bold** and *italic*.

	BERTScore			WikiSimRel (Jaccard)		
	$\hat{q}_t \uparrow$	$\hat{q}_{t'} \downarrow$	Difference $\uparrow$	$\hat{q}_t \uparrow$	$\hat{q}_{t'} \downarrow$	Difference $\uparrow$
Baseline	<b>0.859</b>	0.859	0.000	0.615	<b>0.070</b>	0.545
TopicQGedu	0.855	0.831	0.024	0.721	0.185	0.536
TopicQG	<b>0.859</b>	<i>0.830</i>	<i>0.029</i>	<i>0.727</i>	<i>0.132</i>	<i>0.595</i>
8bit	<i>0.858</i>	0.831	<i>0.027</i>	0.693	0.142	0.551
4bit	<i>0.858</i>	0.831	0.027	0.686	0.157	0.529
TopicQG2X	<b>0.859</b>	<b>0.823</b>	<b>0.036</b>	<b>0.735</b>	<b>0.055</b>	<b>0.680</b>

# 4) Issues with Task Replacement: Prior Knowledge & Motivation

### Classroom settings

#### Participants (please choose)

AI Teacher

Teacher Liu

Responsible for delivering lectures and answering students' questions.

---

AI Teacher

Teaching Assistant

Responsible for maintaining classroom order and answering students' questions.

---

AI Peer

Sparker

An active student in the classroom to liven up the atmosphere.

---

AI peer

Thinker

A thoughtful student who critically reflects on the lesson content, offers counterexamples or raises questions.

---

AI peer

Note-taker

A student who frequently takes notes during class and shares their notes with the whole class.

---

AI peer

Questioner

A student who loves to ask questions and frequently raises deep, thought-provoking questions based on the teacher's lecture content.

#### Classroom mode

Observation  
The AI teacher gives the lecture without interruptions. AI peers occasionally discuss the content and you can see their discussion.

Interaction  
You can ask questions freely. Various AI agents will respond to your questions, and the lecture will pause automatically during the discussion.

#### Other settings

Automatically play voice  on  off

Classroom discussion activity  Low —  medium —  high

Speech interval  Low —  medium —  high

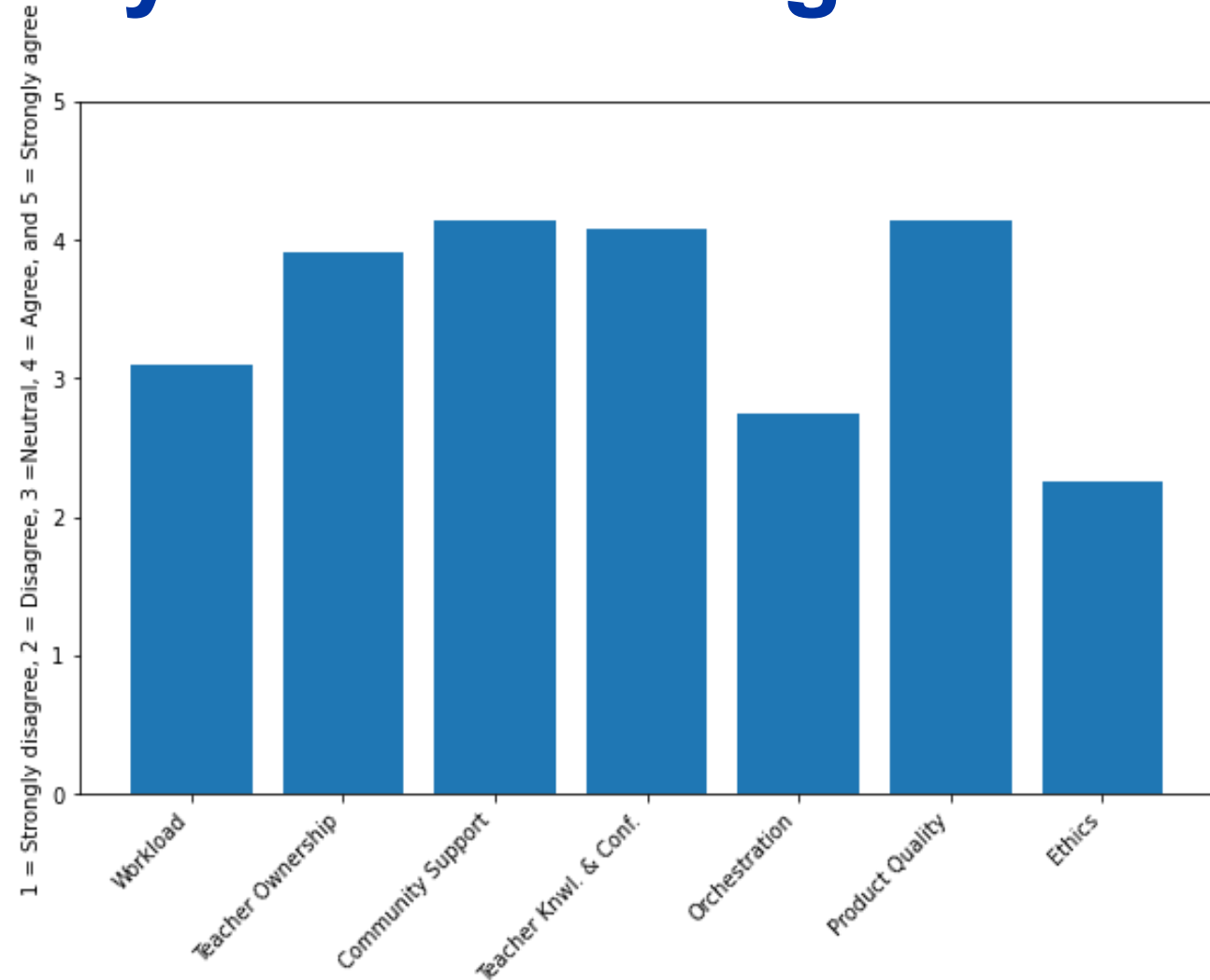
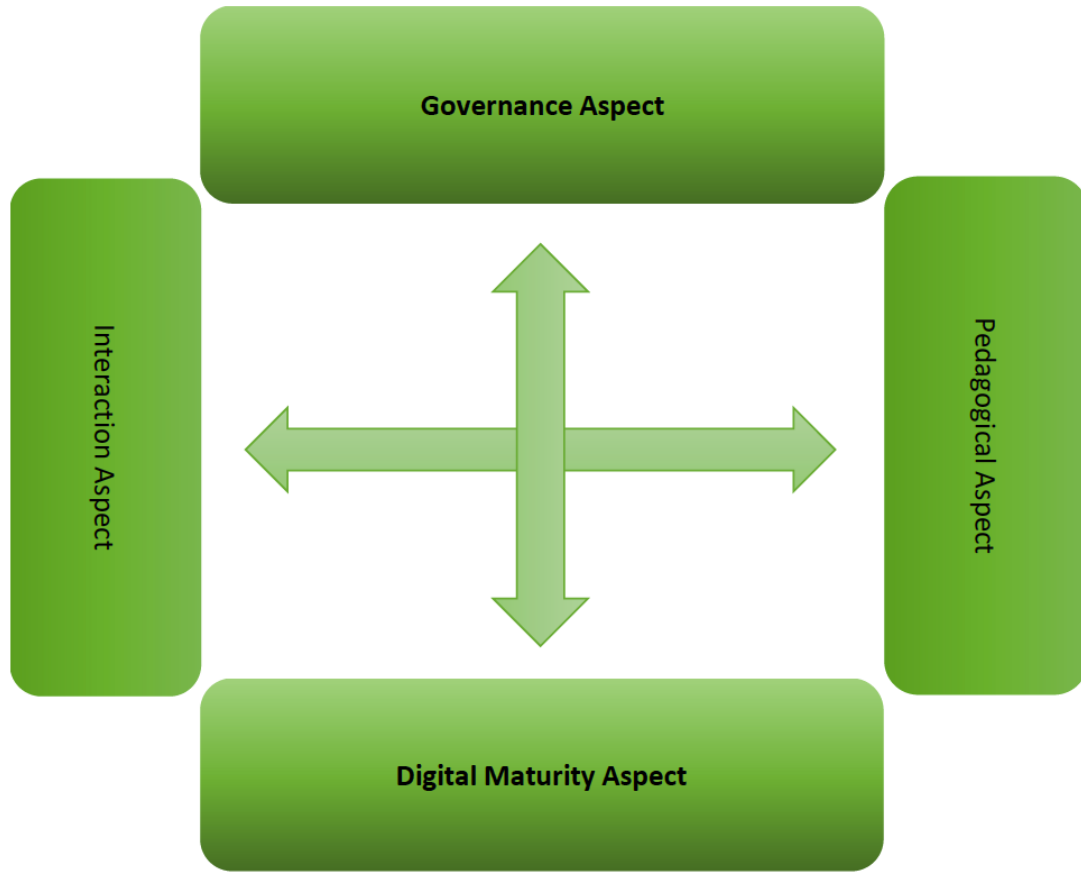
Text display speed  Low —  medium —  high

confirm

## *Massive AI-empowered Course System (MAIC)*

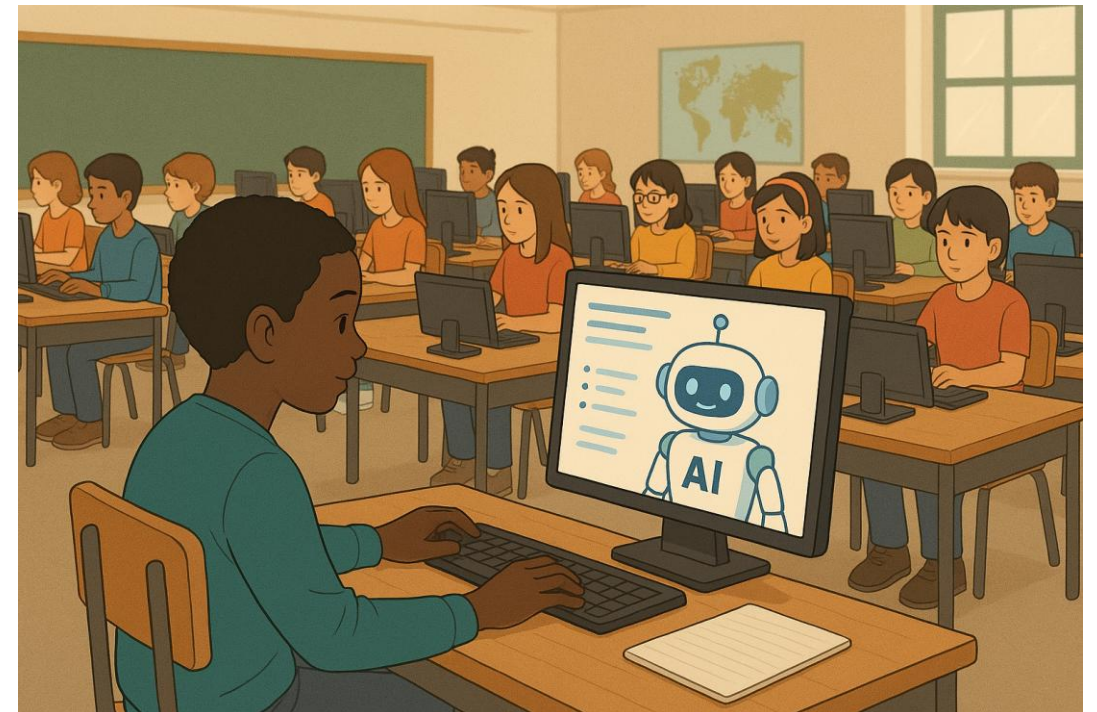
- *Students in the high-prior-knowledge group engaged in significantly more multi-turn conversations with AI agents.*
- *High prior knowledge group experienced a learning loss.*
- *Low prior knowledge group reported significantly higher post-course motivation.*

# 5) Issues with Task Replacement: Socio-technical ecosystem challenges



# Even if we address them all, what is the future of education?

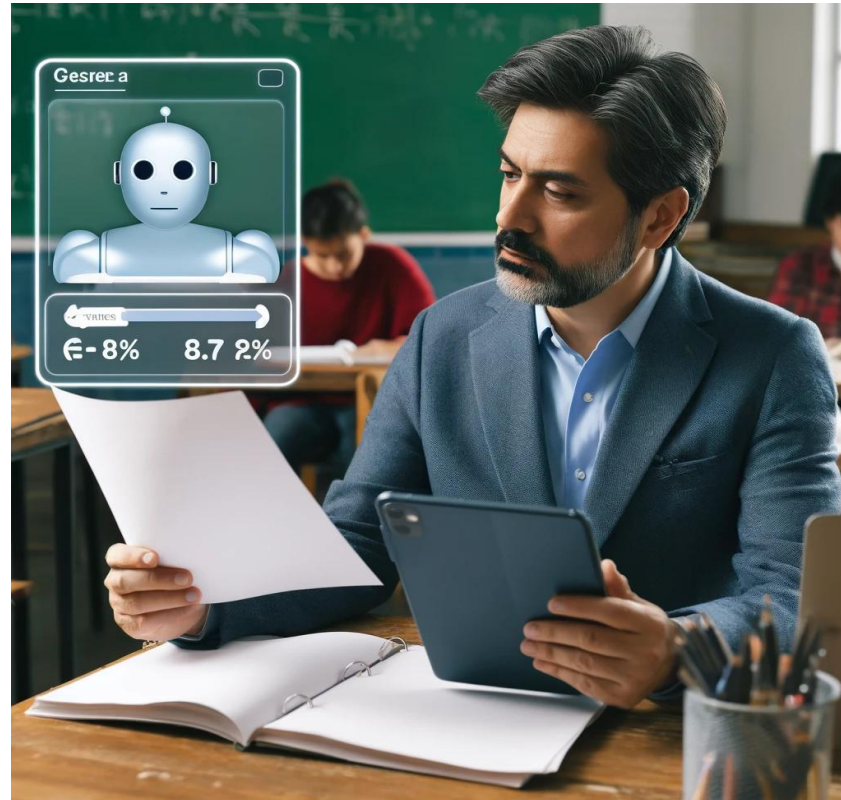
Learning is not only about knowledge acquisition, and education is not only about learning.



# AI should NOT be used to automate existing bad practices in education



AI prepares student submissions



AI marks student submissions



Fancy AI schools in which nobody actually learns and teaches

# If not carefully integrated into meaningful learning designs, AI has the potential to Dehumanise Education



# Assessments need to shift away from focusing solely on final products



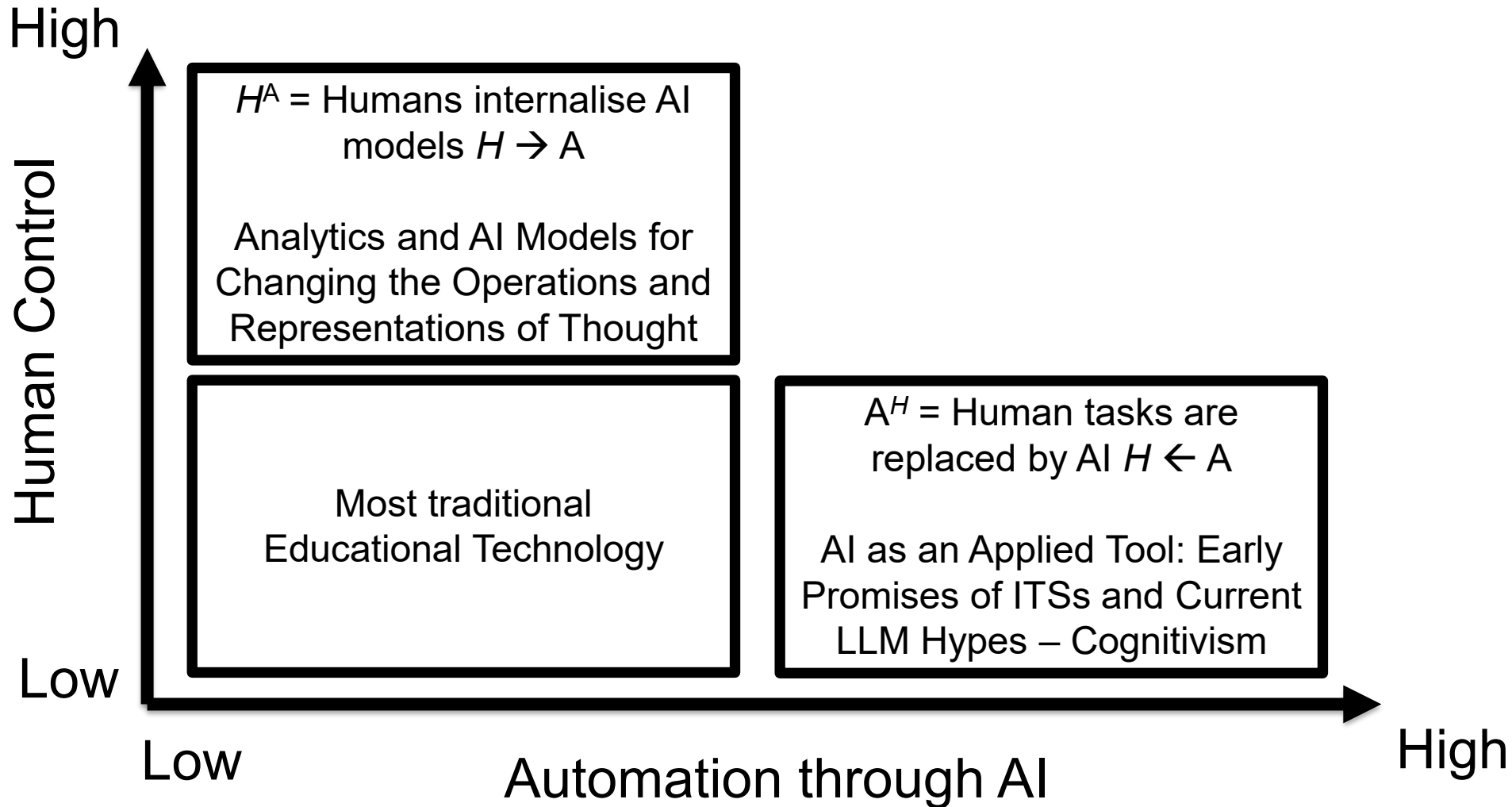
Process

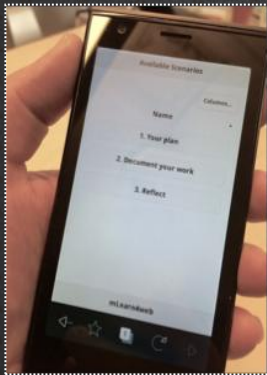
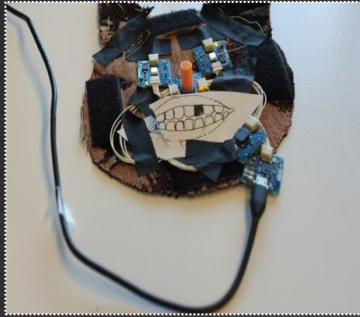
vs





Outcome

# AI in Education: A vision for the future

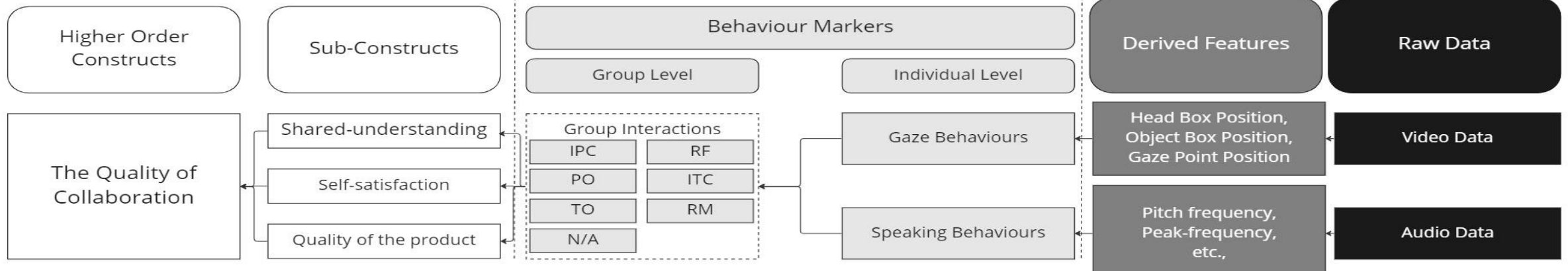
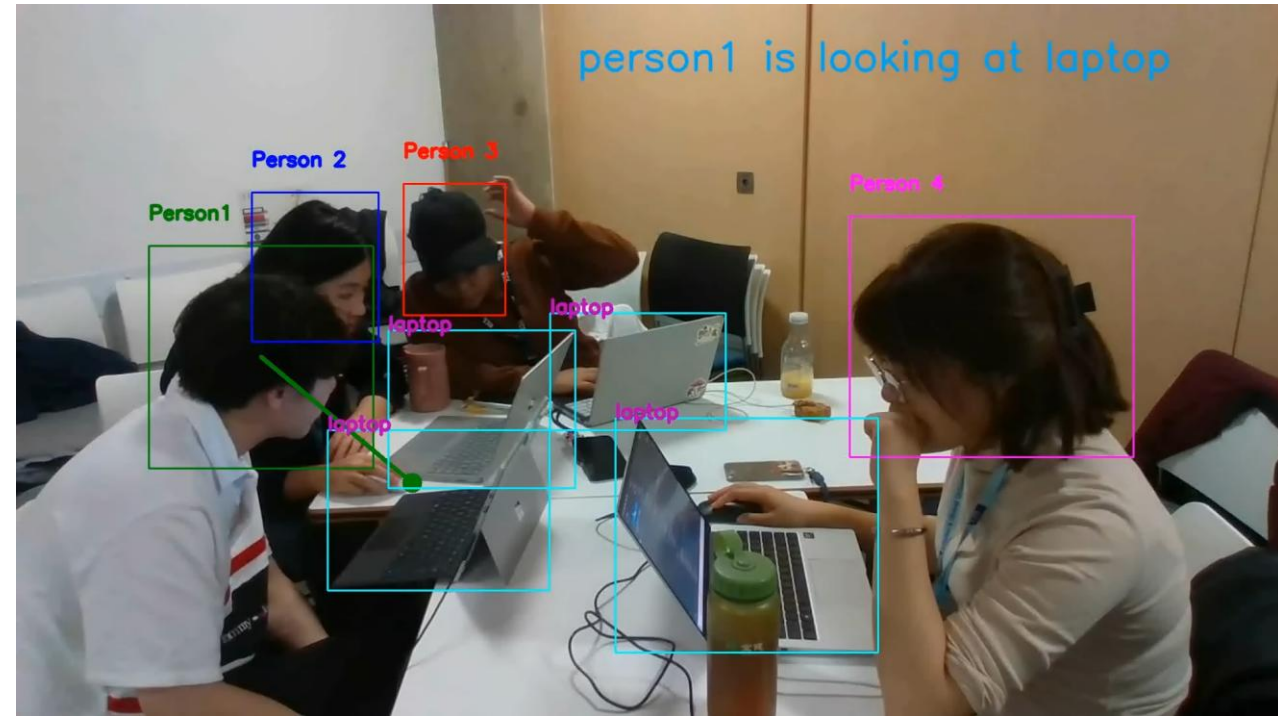
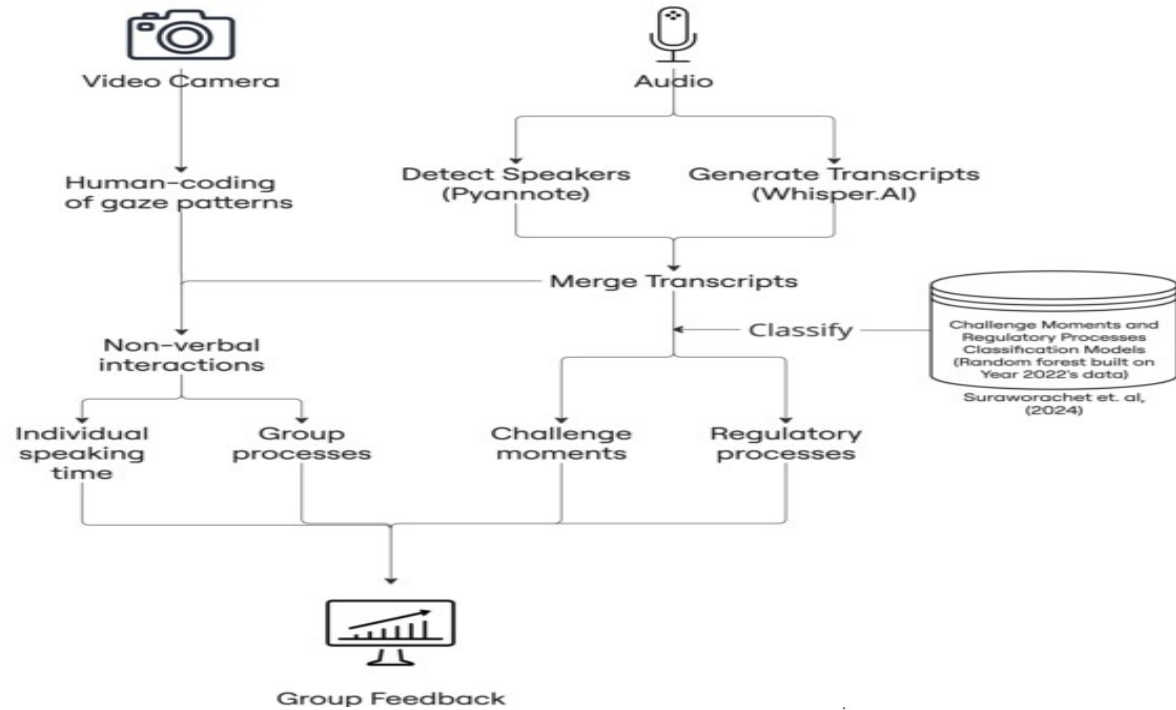




 → YAY! I GOT IT !!!

 → ARGH! THAT'S HARD!!

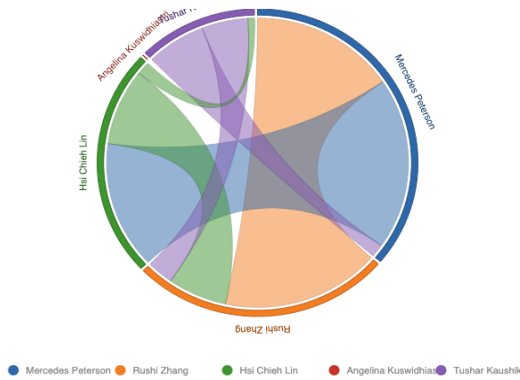
# Making Lived Experiences Visible for Reflection



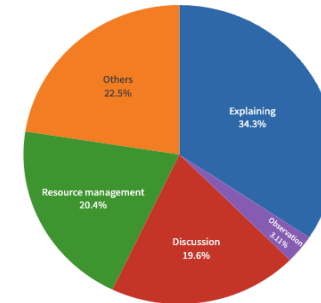
# Value of Making Lived Experiences Visible: Students' group interactions and regulation challenges

## How did we act?

Mutual discussion (times)



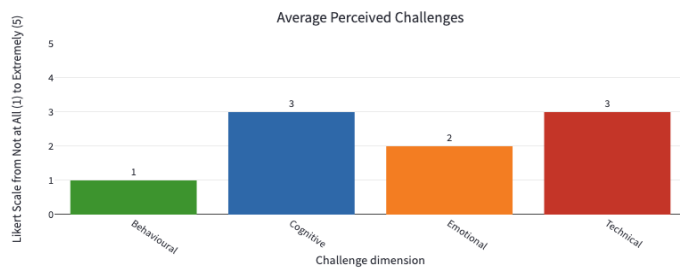
Which group interactions were your group in?



## How did we regulate?

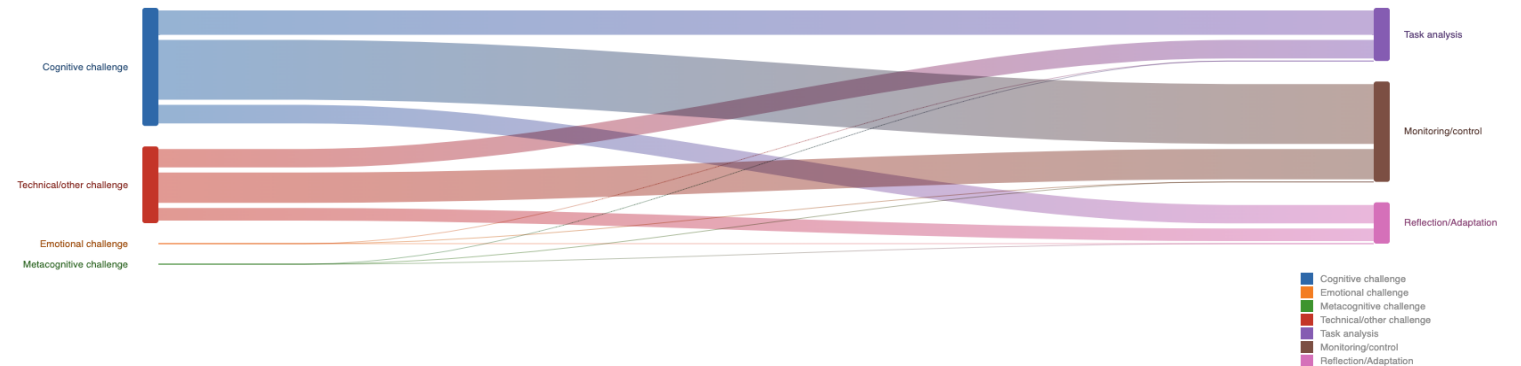
Reported challenges

from post-surveys

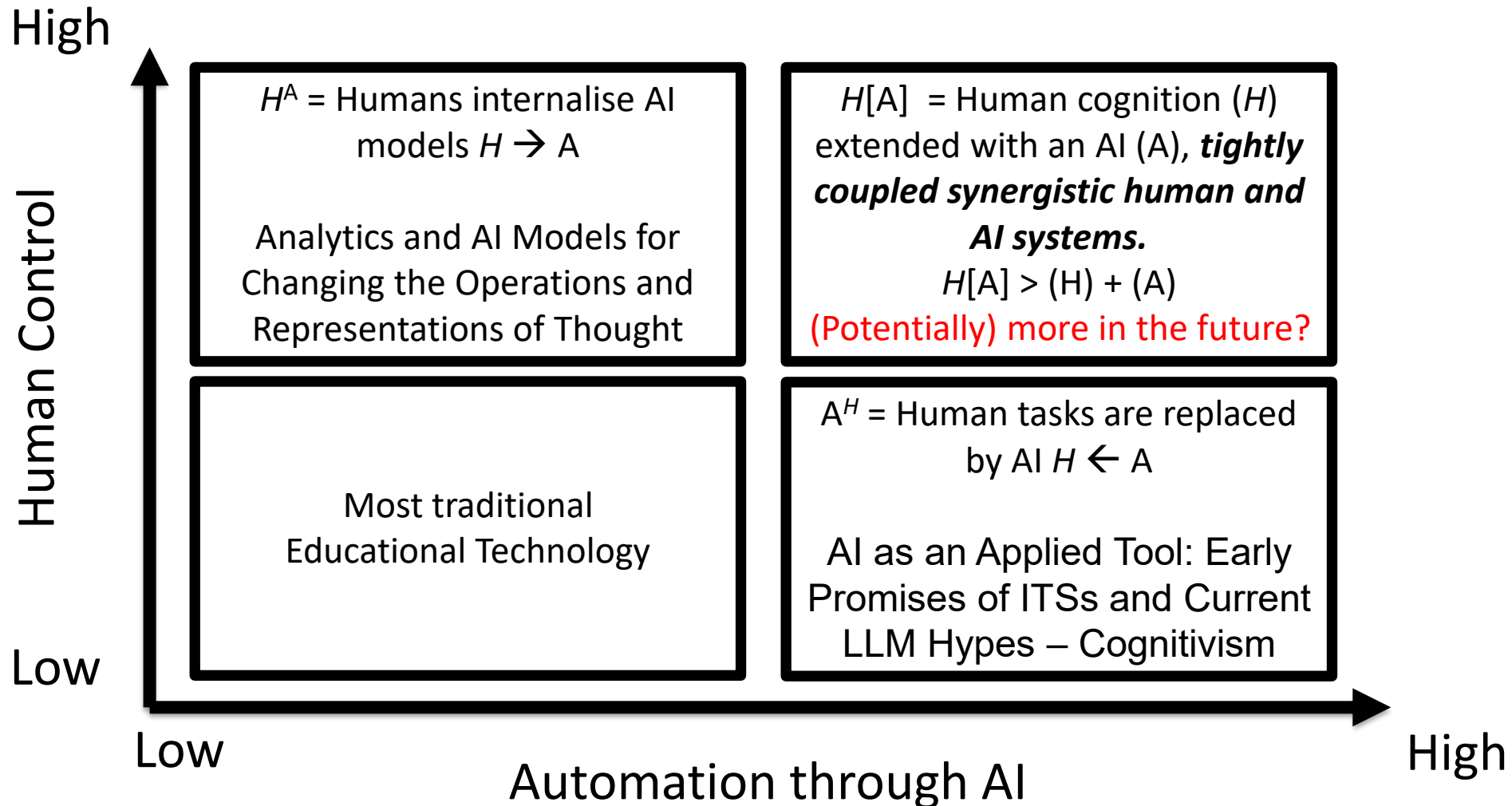


Detected challenges and regulation

from transcript

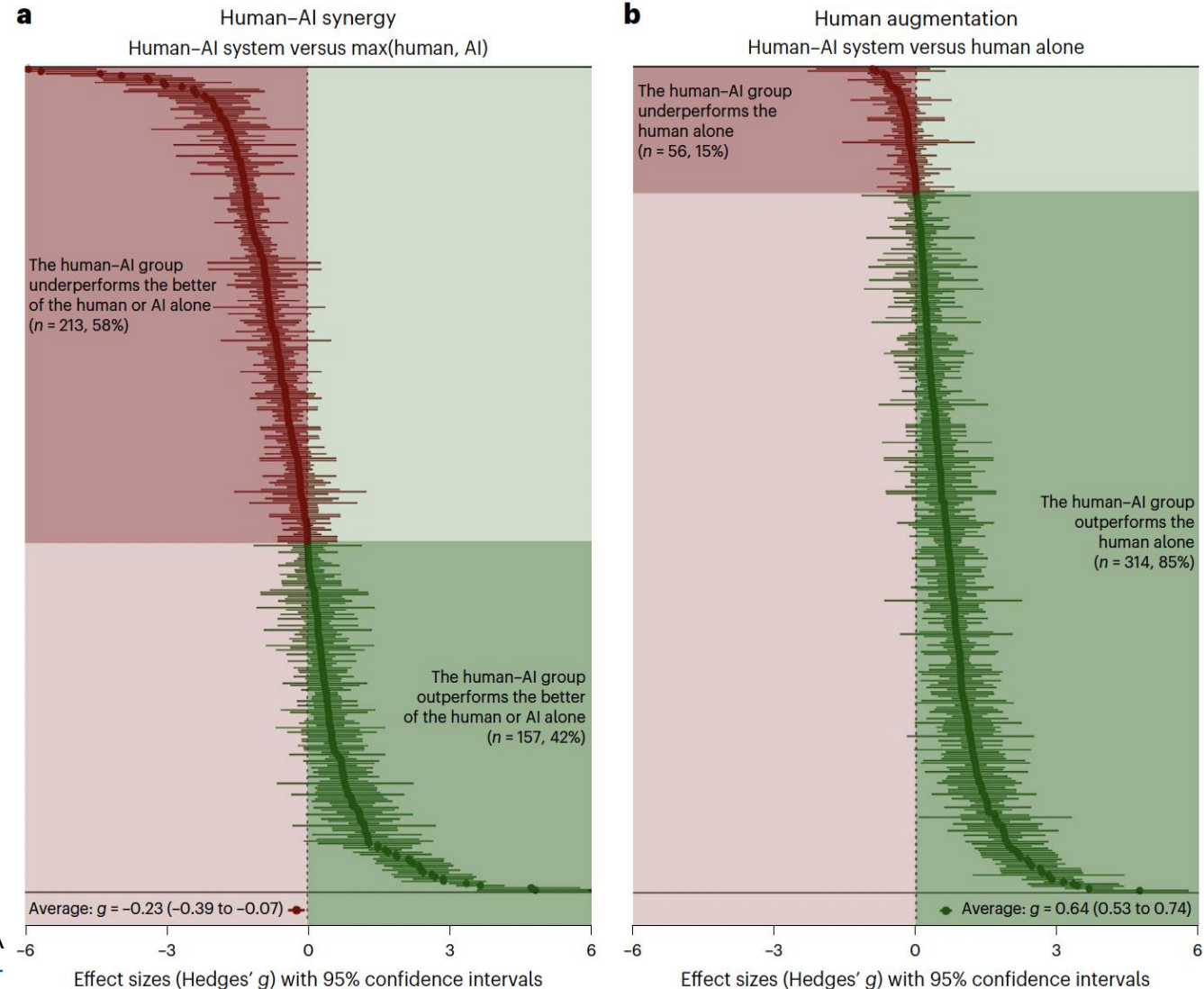


# AI in Education: A vision for the future

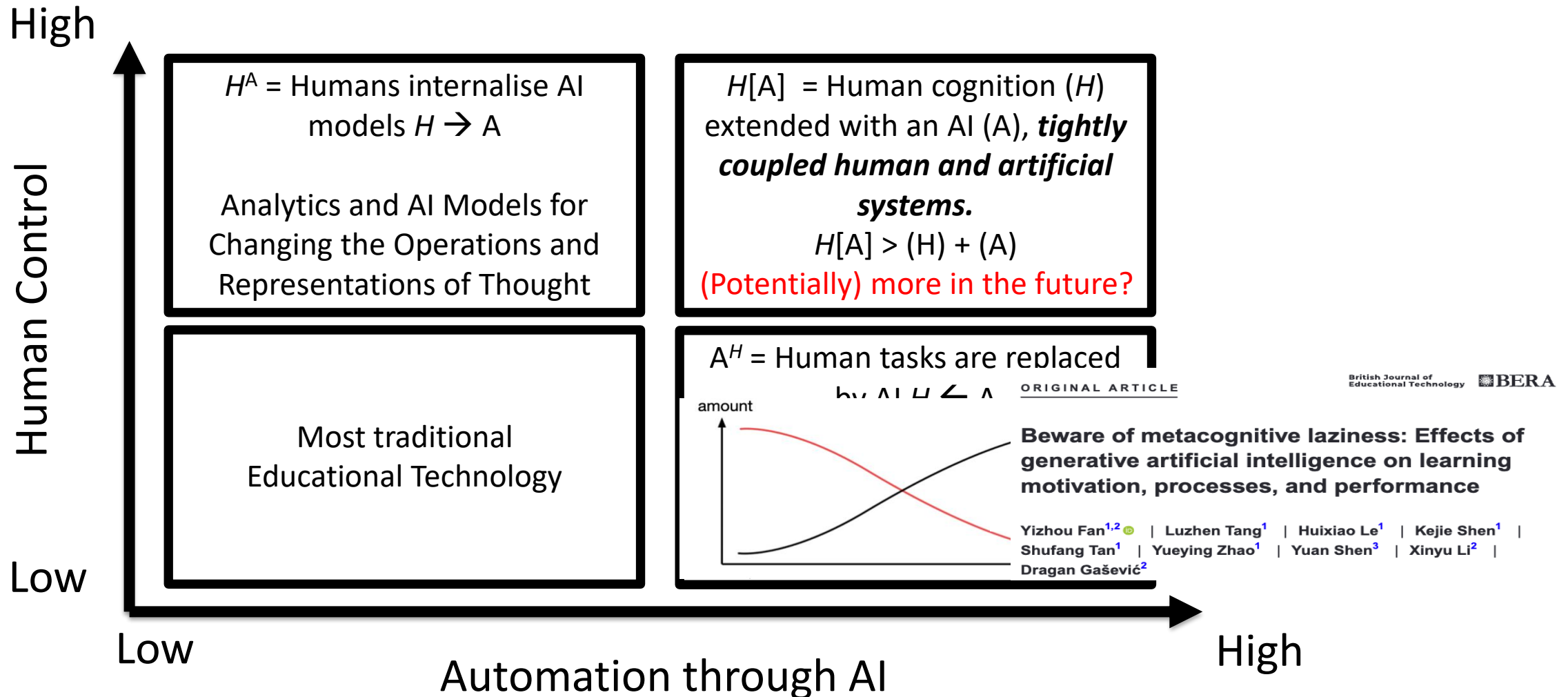


# Not all human-AI teaming is synergistic

- On average, human-AI teaming significantly better than the human alone (**complementarity**); yet performs **significantly worse** than the best of humans or AI alone (**augmentation**).



# AI in Education: A vision for the future



# Potential long-term negative implications of AI-dependence

Our algorithms have flagged that this comment may not be helpful.

- Please ensure that the comment:
  - Is specific to this resource
  - Suggests meaningful improvements for the author
  - Aligns with the grades assigned in the rubric

---

Our algorithms have flagged that this comment may not be helpful.

- Your moderation comment is too similar to your previous comments

**Resource Feedback** I don't want to moderate this resource

Please evaluate the resource based on the following criteria:

Alignment with course context & objectives:	Poor	Needs Improvement	Satisfactory	Great	Outstanding
Correctness, clarity & ease of understanding:	Poor	Needs Improvement	Satisfactory	Great	Outstanding
Appropriateness of difficulty:	Poor	Needs Improvement	Satisfactory	Great	Outstanding
Encouragement of critical thinking and reasoning:	Poor	Needs Improvement	Satisfactory	Great	Outstanding

**Justify your responses & provide feedback**

Please provide constructive feedback & justify your responses to the author so they can improve the resource.

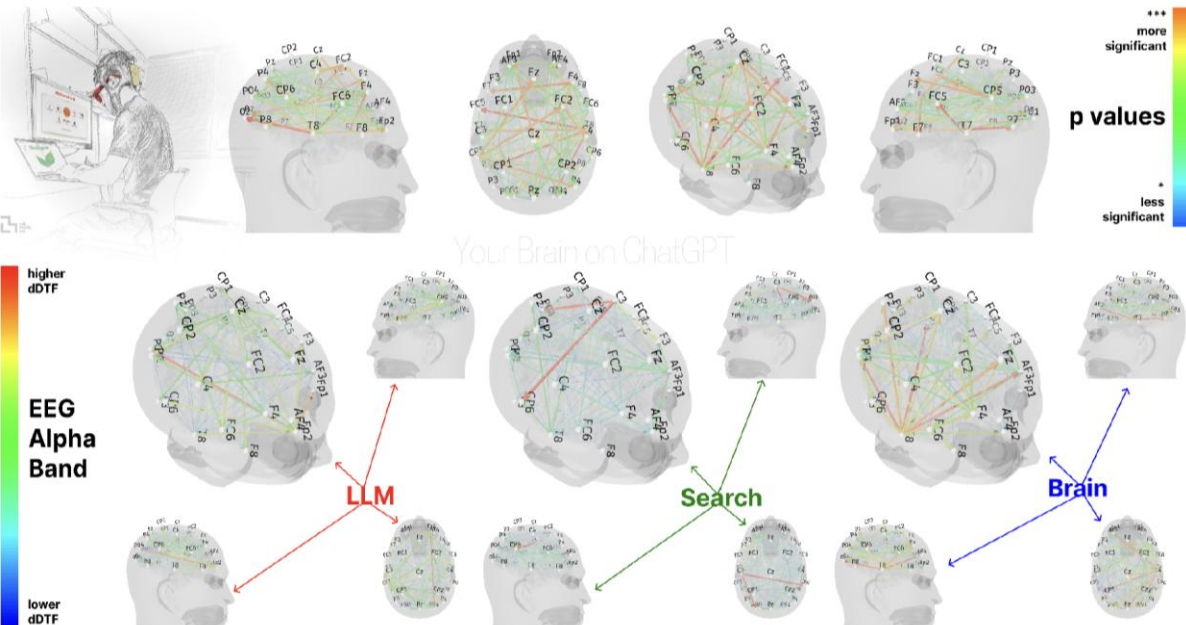
Align feedback with rubric
  Be detailed & specific
  Suggest improvements
  Use constructive language

The content of the question is good. However, please revise the answer. Distractor C is also correct in addition to the option A that has been selected. It would be better to change the question either by making it a multiple-answer or negating option C.

- With 1625 students across 10 courses, they tended to rely on rather than learn from AI assistance.

- Brain-only group exhibited the strongest, widest-ranging networks, Search Engine group showed intermediate engagement, and LLM elicited the weakest overall coupling.

- LLM-to-Brain participants showed weaker neural connectivity.



# AI Competency Framework for Teachers

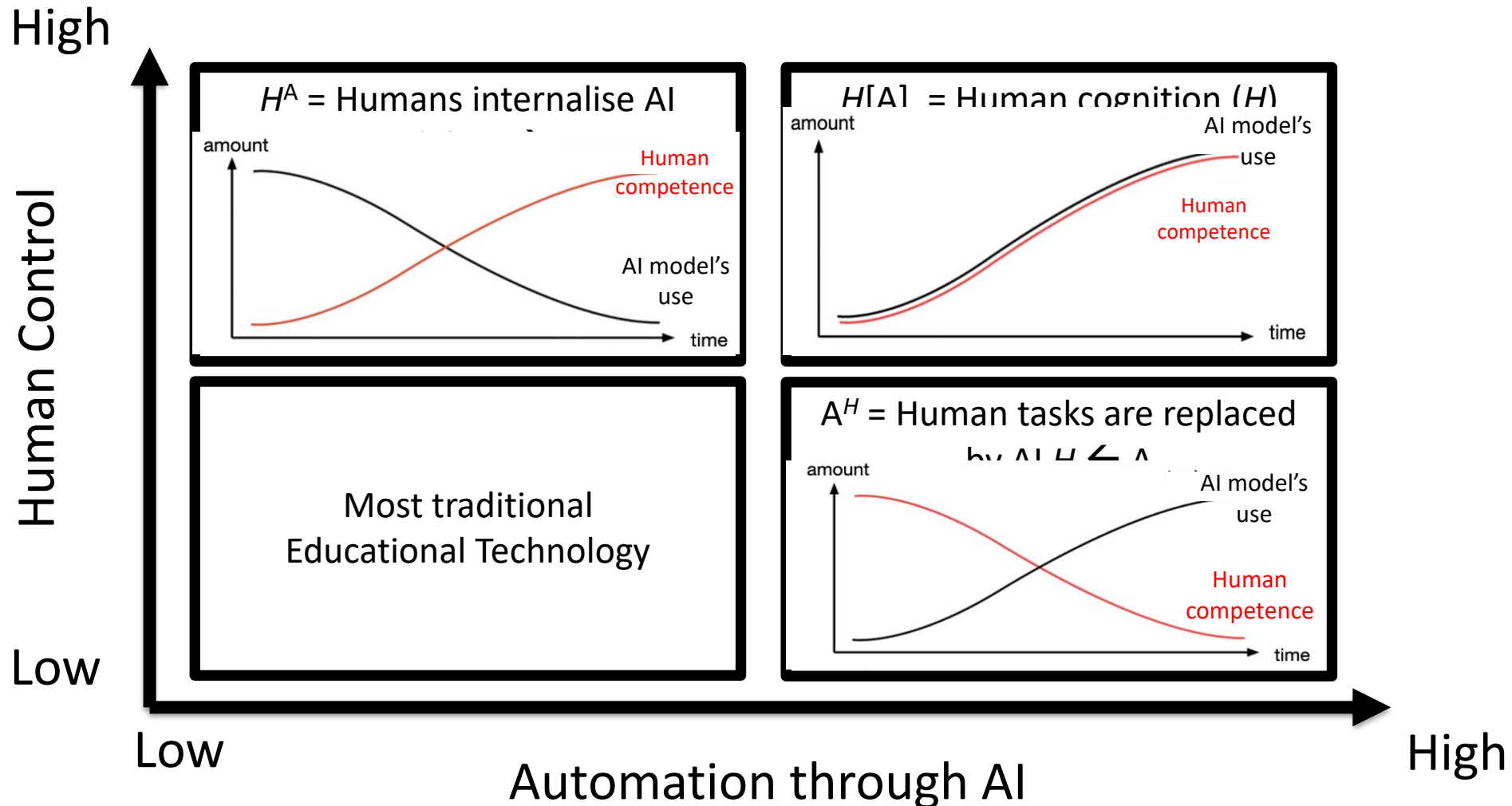
Aspects	Progression		
	Acquire	Deepen	Create
Human-centred Mindset	Human agency	Human accountability	AI social responsibility
Ethics of AI	Ethical principles	Safe and responsible use	Co-creating AI ethics
AI Foundations & Applications	Basic AI techniques and applications	Application skills	Creating with AI
AI Pedagogy	AI-assisted teaching	AI-pedagogy integration	AI-enhanced pedagogical transformation
AI for Professional Development	AI enabling lifelong professional learning	AI to enhance organizational learning	AI to support professional transformation



AI competency framework  
for teachers



# AI in Education: A vision for the future





# UCL

# Thank you

Professor Mutlu Cukurova  
University College London

[m.cukurova@ucl.ac.uk](mailto:m.cukurova@ucl.ac.uk)

