

# Assessment Design in the AI Era: Applying Psychometric Theory to Identify Items on Which Humans and Chatbots Diverge

Licol Zeinfeld, Alona Strugatski, Ziva Bar-Dov,  
Ron Blonder, Shelley Rap, Giora Alexandron



# Generative AI in Education Challenge

- GenAI matches/exceeds student performance on many assessments across domains
- Widespread use
- Learning and assessment may be compromised
- Need to rethink assessment

## RESEARCH ARTICLE

### ChatGPT performance on multiple choice question examinations in higher education

Philip Newton & Maira Xiromeriti

Assessment & Evaluation in Higher Education • 2024  
doi:10.1080/02602938.2023.2299059

**“ChatGPT can pass examinations based on multiple choice questions (MCQs), including those used to qualify doctors, lawyers, scientists etc.”**

Newton & Xiromeriti (2024) • Abstract • Pragmatic scoping review



The screenshot shows a news article from The Guardian. The header includes the Guardian logo and navigation links for 'Universities' and 'Students'. The article title is 'UK universities warned to 'stress-test' assessments as 92% of students use AI', with a sub-headline 'Survey of 1,000 students shows 'explosive increase' in use of generative AI in particular over past 12 months'. The author is Sally Weale, an education correspondent, and the article is dated Wednesday, 26 February 2025, 00:01 GMT. A 'Share' button is visible at the bottom right. The main image shows a white keyboard and a smartphone displaying the ChatGPT logo.

**The Guardian**

Universities Students

**Higher education**

**UK universities warned to 'stress-test' assessments as 92% of students use AI**

Survey of 1,000 students shows 'explosive increase' in use of generative AI in particular over past 12 months

**Sally Weale** Education correspondent

Wed 26 Feb 2025 00:01 GMT

Share

# Rethinking Assessment

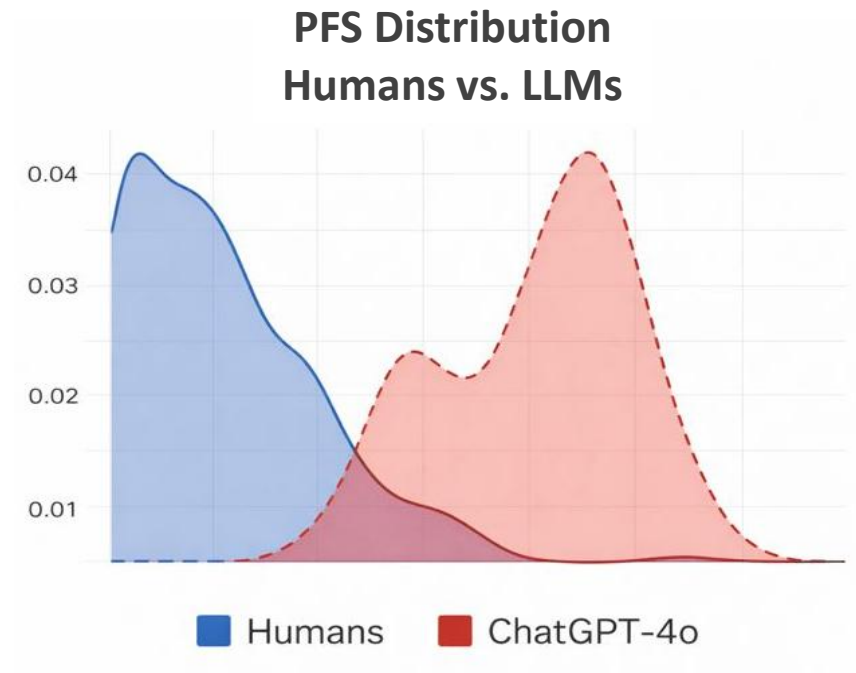
---

- Adapt assessment to growing GenAI presence
- Understand LLMs strengths and weaknesses
  - Current LLM evals rely on descriptive statistics
  - Psychometrics offer principled ways to study GenAI performance
- A lot of work on LLM responses to text-rich prompts
- Very little work on multiple-choice assessments

1	(A) (B) ● (D) (E)	11	(A) (B) (C) (D) (E)
2	● (B) (C) (D) (E)	12	(A) (B) (C) (D) (E)
3	(A) (B) (C) (D) (E)	13	(A) (B) (C) (D) (E)
4	● (B) (C) (D) (E)	14	(A) (B) (C) (D) (E)
5	(A) (B) (C) (D) (E)	15	(A) (B) (C) (D) (E)
6	(A) (B) ● (D) (E)	16	(A) (B) (C) (D) (E)
7	(A) (B) (C) (D) (E)	17	(A) (B) (C) (D) (E)
8	(A) ● (C) (D) (E)	18	(A) (B) (C) (D) (E)
9	(A) (B) (C) (D) (E)	19	(A) (B) (C) (D) (E)
10	(A) (B) (C) (D) (E)	20	(A) (B) (C) (D) (E)

# Applying psychometrics to evaluate GenAI capabilities\*

- Using person-fit statistics (PFS) to distinguish between humans and LLM 'test-takers' on MCQ assessments
  - **No information on the item and the task dimensions** that distinguish humans/LLMs
- => **We need item-level analysis**



# Item Level Analysis Using Differential Item Functioning (DIF)

---

- DIF are methods to identify items on which two groups differ

## Example:

- A math test where some items rely heavily on reading comprehension
- Second language learners (SLLs) may be **disadvantaged** compared to native speakers of the same math skill
- SLLs are termed the 'focal' group, and native speakers the 'reference' group
- **We apply DIF to LLMs:** LLMs are the 'focal' group and Humans are the 'reference'

# Goals and Research Questions (RQs)

---

**Goals:** Examine whether DIF methods can be useful to **identify** items on which humans and LLMs differ and **characterize** the task dimensions that make items easy/difficult for LLMs

## **Research Questions:**

**RQ1:** Can DIF techniques identify items that function differently for humans and LLMs?

**RQ2:** What are the key characteristics that subject-matter experts identify in items on which LLMs exhibit differential performance compared to humans?

## Two DIF Methods being used

---

### Mantel–Haenszel (MH-DIF):

- Simple method that is based on odds-ratio for each ability level
- Requires relatively less response data
- Identifies DIF that is uniform across all ability levels

	LLMs	Humans
Correct	a	b
Incorrect	c	d

$$\text{odds(LLM)} = a/c$$

$$\text{odds(human)} = b/d$$

$$\text{odds ratio} = \frac{\text{odds(LLM)}}{\text{odds(human)}}$$

### Logistic Regression-based (LR-DIF):

- Fits a logistic regression function to the responses
- Can identify more sophisticated DIF patterns (non-uniform, sign-change)
- Requires more data

# Instruments & Data

---

## Instruments & human data:

- Chemistry exam (931 high school students)
- Math sections from the Psychometric exam\* (~4800 examinees)

## GenAI models:

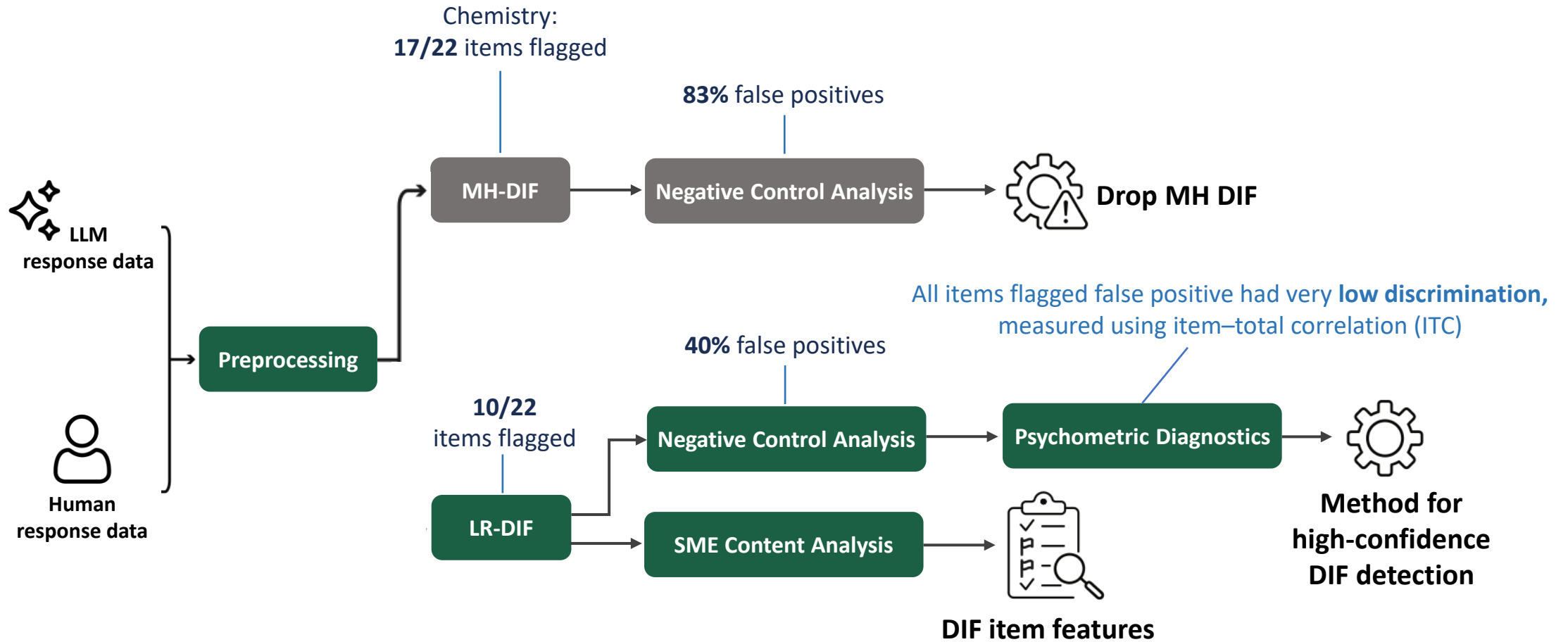
- ChatGPT-4o & 5.2, Gemini 1.5 & 3 Pro, Claude 3.5 & 4.5 Sonnet



 Claude

 Gemini

# Procedure



## Results – LR-DIF (partial results for RQ1 and RQ2)

### RQ1:

LR-DIF flagged items:		uniform	nonuniform
Inst.	POS	NEG	
Psych.	$\bar{11}$ , $\tilde{13}$ - $\tilde{15}$ , $\tilde{34}$ , $\bar{40}$	$\tilde{16}$ , $\bar{21}$ , $\tilde{25}$ , $\tilde{31}$	
Chem.	$\tilde{12}$ , $\bar{16}$ , $\bar{17}$ , $\bar{19}$	$\bar{5}$ - $\bar{7}$ , $\tilde{14}$ , $\bar{15}$ , $\tilde{22}$	

### POS:

- LLM **overperforms** when compared to humans at same ability level

### NEG:

- LLM **underperforms**

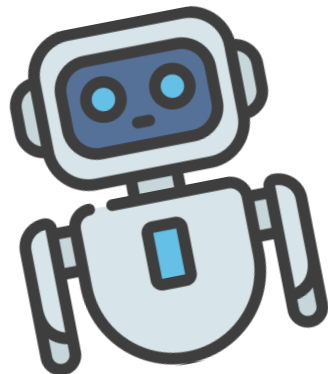
### RQ2:

- POS (LLM overperformance): “strength in tasks that prioritize rule-governed reasoning, conceptual clarity, canonical chemical principle”
- NEG (LLM underperformance): “These items place greater demands on visual interpretation, linguistic nuance, and the execution of complex, multi-step procedures”

To conclude...

---

- LR-DIF on items with reasonable ( $ITC \geq 0.2$ ) discrimination properties is a **robust method** to identify items on which humans and LLMs differ
- Using DIF-identified items can inform item analysis around GenAI and identify **task dimensions** that make items easy/difficult for LLMs
- This can inform assessment redesign that is **GenAI-resilient**



**Thank You!**