



Pedagogy-driven Evaluation of Generative Al-powered Intelligent Tutoring Systems

Kaushal Kumar Maurya and Ekaterina Kochmar

Mohamed bin Zayed University of Artificial Intelligence (MBZUAI), Abu Dhabi, UAE

AIED 2025 (BlueSky Track)

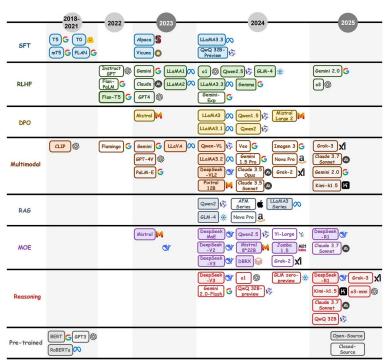


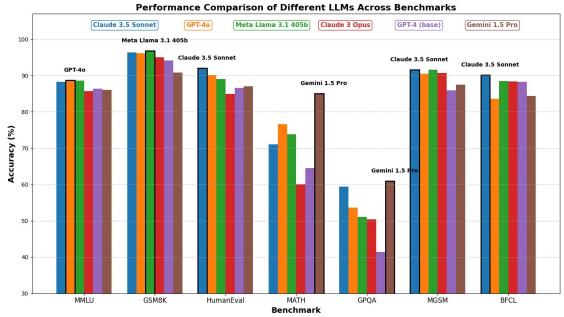


- ☐ Introduction
- Current State of Evaluation
- Key Challenges
- Path Forward
- Conclusion

- ☐ Introduction
- Current State of Evaluation
- □ Key Challenges
- Path Forward
- Conclusion

Generative AI models (aka. LLMs) have been impressive!

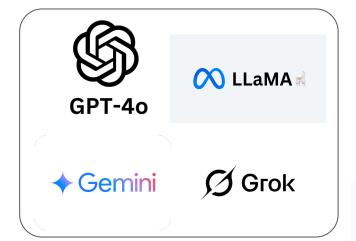




LLMs Performance on Diverse Benchmarks

Impact in Educational Domain: Raise of Al Tutors (aka. ITSs)











Implications in Educational Domain: Can we trust them?



George W. Bush

Problem (In 2002): American students do not stand high in education globally

NCLM Act: Attempt to fulfill the expectations

Pros: Standardized test (equal opportunity), higher test scores

Cons: Standardized test scores led to

shallow learning, reduced

engagement, lots of expectations from teachers

Generative Al Can Harm Learning (Süngü, A., et al. 2024)

Setup:

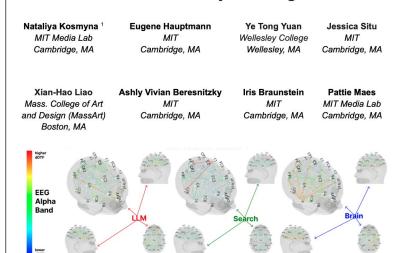
- Pre-registered, Randomized Controlled Trial
- 9-11th grades in Turkey, 2023-2024 Fall semester
- Fifty, 90-minute classes
- Nearly 1000 students with 15% of the math curriculum
- Models: ChatGPT Base, GPT Tutor (ChatGPT + Prompts), GPT4

Outcome:

- With GPT4, 127% improvement over GPT Tutor
- 17% reduction in performance without tutor access

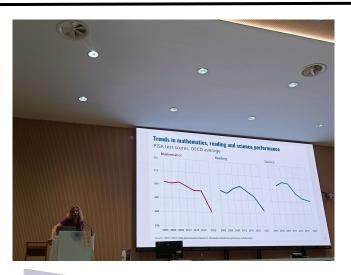
June 2025

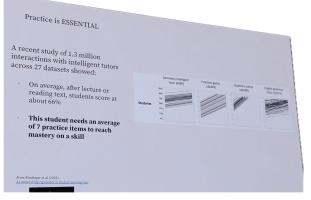
Your Brain on ChatGPT: Accumulation of Cognitive Debt when Using an Al Assistant for Essay Writing Task[△]



Implications in Educational Domain: Can we trust them?







Root Causes: Unreliable and Inconsistent Evaluation Practices

Large Exploration Space:

- Educational research is multidisciplinary
- Many learning theories, limited RCTs
- Many tutor moves, little clarity on which lead to learning gains
- Personalized protocols and benchmarks with homogeneous populations
- Evaluation and guidelines are subjective or inconsistent
- Generic metrics used to evaluate pedagogical capabilities
- Models and research are centered in fewer (or WEIRD) countries — less inclusive and less adaptive
- Many more

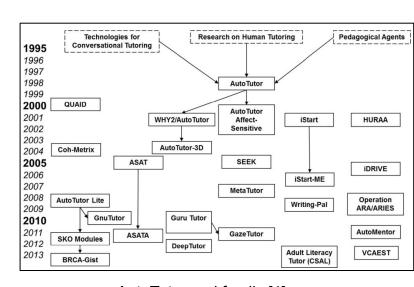
- Introduction
- Current State of Evaluation
- □ Key Challenges
- Path Forward
- Conclusion

Current State of Evaluation: Traditional Approaches

- Assessing teachers' practices [2]
 - Analysis of classroom artifacts (teacher' assignments and student work)
 - Teaching portfolios, self-reports, logs, interviews, etc.
- Role-simulation reports & extrinsic studies [2]
- Self-assessment with self-regulatory learning theories [3]

Limitations [4]:

- Not applicable to AI responses due to static nature of evaluation
- Not-scalable
- Hard to handle ethical and biased Al responses



AutoTutor and family [1]

Ref: [1] Nye, Benjamin D., Arthur C. Graesser, and Xiangen Hu. "AutoTutor and family: A review of 17 years of natural language tutoring." International Journal of Artificial Intelligence in Education 24.4 (2014): 427-469.

^[2] Goe, L., Bell, C., Little, O.: Approaches to Evaluating Teacher Effectiveness: A Research Synthesis. National Comprehensive Center for Teacher Quality (2008)

^[3] Crossley, S.A., Varner, L.K., Roscoe, R.D., McNamara, D.S.: Using Automated Indices of Cohesion to Evaluate an Intelligent Tutoring System and an Automated Writing Evaluation System. AIED 2013.

^[4] Macina, J., Daheim, N., Wang, L., Sinha, T., Kapur, M., Gurevych, I., Sachan, M.: Opportunities and Challenges in Neural Dialog Tutoring. In: Vlachos, A., Augenstein, I. (eds.) EACL 2024.

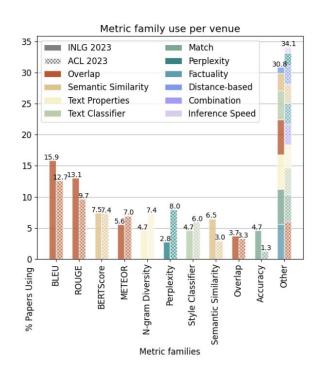
Current State of Evaluation: NLG-based Evaluation

NLG-based Metrics for Gen-Al ITS [1,2]

- Lexical and semantic match
- Validate coherence, fluency, human-likeness, etc.
- Ex: BLEU, BERTScores, etc.

Limitations [3]:

- Do not capture the pedagogical values
- Require a gold reference (often unavailable or non-unique)
- Can be hacked with superficial responses [4]



Ref: [1] Patricia Schmidtova, Saad Mahamood, Simone Balloccu, Ondrej Dusek, Albert Gatt, Dimitra Gkatzia, David M. Howcroft, Ondrej Platek, and Adarsa Sivaprasad. Automatic Metrics in Natural Language Generation: A survey of Current Evaluation Practices. INLG 2024

^[2] Tack, A., Kochmar, E., Yuan, Z., Bibauw, S., Piech, C.: The BEA 2023 Shared Task on Generating Al Teacher Responses in Educational Dialogues. In: Proceedings of BEA 2023. pp. 785–795 (2023)

^[3] Jurenka, I., Kunesch, M., McKee, K.R., Gillick, D., Zhu, S., Wiltberger, S., Phal, S.M., Hermann, K., et al.: Towards Responsible Development of Generative Al for Education: An Evaluation-Driven Approach. arXiv. 2024.

^[4] asselli, J., Vasselli, C., Noheil, A., Watanabe, T.: NAISTeacher: A Prompt and Rerank Approach to Generating Teacher Utterances in Educational Dialogues. In: Proceedings of BEA 2023. pp. 772–784 (2023)

Current State of Evaluation: Pedagogically Oriented

Human experts to evaluate pedagogical performance [1,2]:

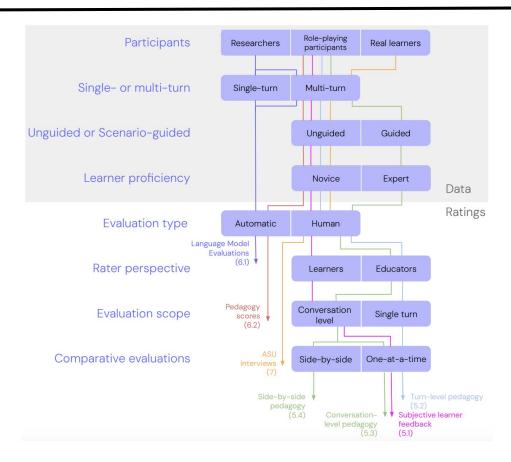
Better correlation with user satisfaction

Limitations [1,2,3]:

- Limited access to pedagogical experts
- Small number of annotators, leading to biases
- No unified protocol for evaluation
- Paid raters act as learners
- Evaluations with real students done in small number of participants - not scalable

- Introduction
- Current State of Evaluation
- □ Key Challenges
- Path Forward
- Conclusion

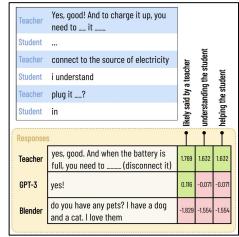
Key Challenges: Diverse Pragmatics of Evaluation

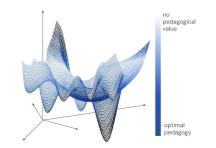


Key Challenges: Lack of Unified Pedagogical Practices

- Koedinger et al. (2013): Synthesized a list of 30 independent instructional principles
- Tack and Piech (2022): Proposed 3 strategies for good tutor
- Jurenka et al. (2024): Collected 28 strategies from various educational stakeholders







Key Challenges: Lack of Unified Evaluation Benchmarks

Dataset	Synthetic?	Domain	#Dialogues	#Moves	Grounding	Setting
CIMA	No	Language	391	5	Image, answer	1:1
Bridge	No	Math	700	10	Image, Confusion	1:1
MathDial	Yes	Math	2861	4	Confusion, answer	1:1
TSCC	No	Language	102	5	None	1:1
TalkMoves	No	Science	567	10	None	Classroom
NCTE	No	Math	1660	None	None	Classroom
MRBench	Mixed	Math	500	10	Confusion, answer	1:1

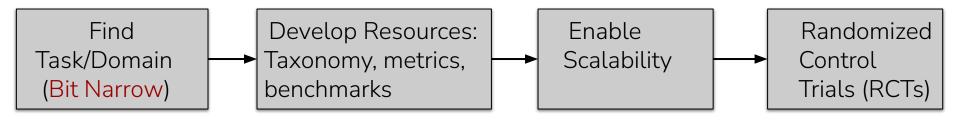
- Introduction
- Current State of Evaluation
- □ Key Challenges
- Path Forward
- Conclusion

Path Forward I: Evaluation Unification (Our Vision)

Research Hypothesis: Lack of task/domain-specific unified evaluation limits progress in ITSs.

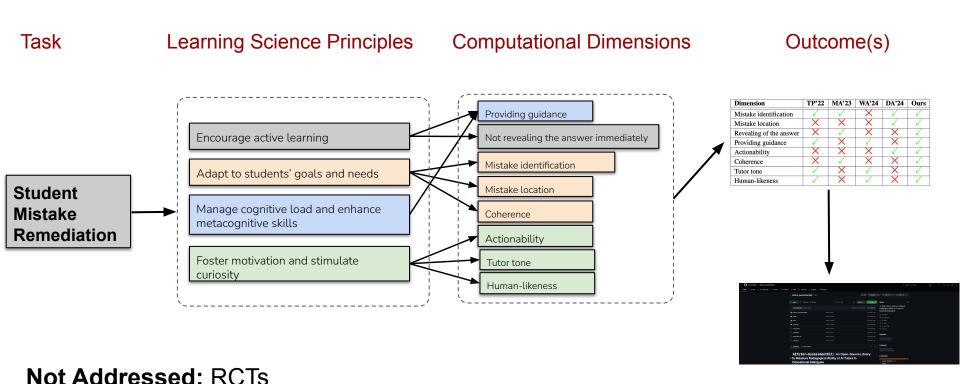
Research Statement: Develop unified evaluation taxonomies, metrics, and benchmarks tailored to specific tasks/domains to evaluate ITSs.

Potential Approach:



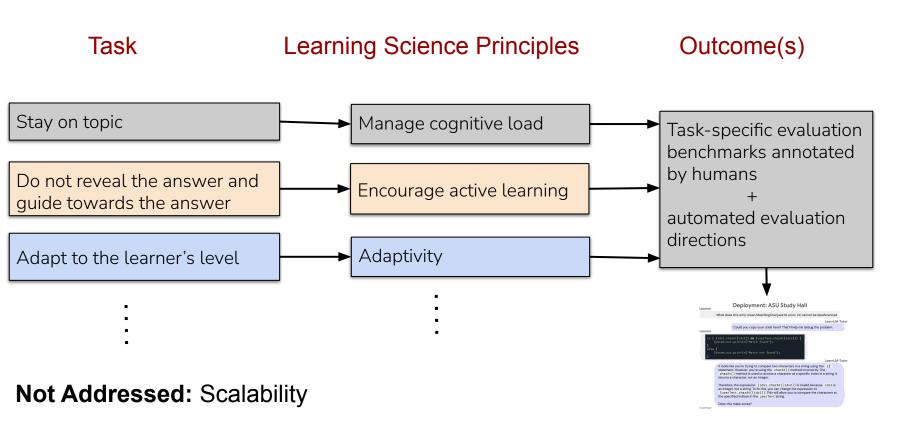
Extensions: Multi-Agent Systems

Evaluation Unification: Case Study I



Ref: Maurya, Kaushal Kumar, K. V. Srivatsa, Kseniia Petukhova, and Ekaterina Kochmar. "Unifying Al tutor evaluation: An evaluation taxonomy for pedagogical ability assessment of LLM-powered Al tutors." NAACI. 2025

Evaluation Unification: Case Study II



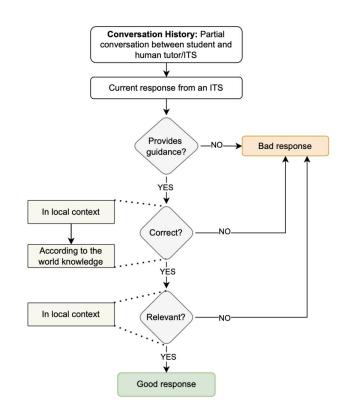
Path Forward II: Measuring Pedagogical Guidance

Research Hypothesis: By offering timely and context-aware appropriate guidance (e.g., hints, explanations, etc.), tutors can help students navigate their learning journey.

Research Statement: Develop a robust quantitative framework to assess the appropriateness and richness of the pedagogical guidance provided by GenAl-powered ITSs.

Grounded on: Metacognition

Potential Approach:



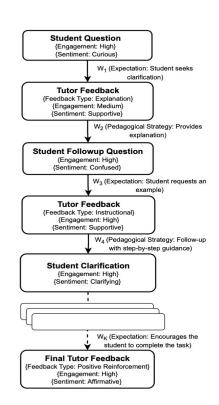
Path Forward III: Measuring Active Learning

Potential Approach:

Research Hypothesis: Active learning enables students to engage more deeply in the learning process through critical thinking and reflection.

Research Statement: Develop a robust quantitative framework to measure the active learning capabilities of GenAl-powered ITSs.

Grounded on: Constructivism and inquiry-based learning



- Introduction
- Current State of Evaluation
- □ Key Challenges
- Path Forward
- Conclusion

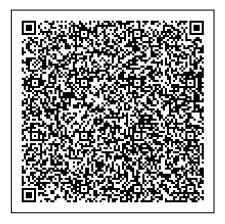
Conclusion

- ITS evaluation is challenging due to the diverse pragmatics of evaluation.
- Current state-of-the-art evaluation approaches are unreliable, inconsistent, or subjective.
- The community should work towards a pedagogy-oriented evaluation framework that is unified and grounded in learning science principles.









Paper

Let us know your thoughts!!!

Correspondence: