

AIME-Con 2025

**Artificial Intelligence in Measurement and Education
Conference (AIME-Con)**

Volume 2: Works in Progress

October 27-29, 2025

The AIME-Con organizers gratefully acknowledge the support from the following sponsors.

Platinum



Gold



*ets research institute

Silver



Gates Foundation



Supporters





duolingo english test

The future of language assessment is here

The Duolingo English Test is a computer adaptive test powered human-in-the-loop AI and supported by rigorous validity research. The test measures speaking, writing, reading, and listening skills, providing a deeper insight into English proficiency.



Built on the latest language assessment science

- ✓ Accessible by design, supporting test takers wherever they are for just \$70
- ✓ Built on rigorous research and industry- leading security
- ✓ Integrates the latest assessment science and AI for accurate results
- ✓ Accepted by over 5,800 programs worldwide



englishtest.duolingo.com



Evidence-based approach to AI in Measurement & Learning

At the intersection of artificial intelligence and educational measurement, Pearson stands as your trusted partner—delivering clarity, confidence, and innovation in every assessment moment.

Why Pearson?

- **AI-Enhanced Accuracy:** Using automated scoring and predictive analytics to provide insights that are accurate, fair, and timely.
- **Future-Ready Solutions:** Platforms that evolve with policy, pedagogy, and technology.
- **Personalized Learning Journeys:** Multi-lingual access and adaptive item generation to support each student's unique growth trajectory.
- **Ethical AI Practices:** Commitment to data security, transparency, explainability, and bias mitigation.
- **Collaborative Innovation:** Partnering with educators, researchers, and technologists to shape the future of assessment.

Human-Centric AI	Pearson believes AI's highest purpose is to elevate and empower human capabilities.
Assessment as a Learning Continuum	We reimagine assessments not as endpoints, but as integral parts of the learning journey.
AI as an Environment	Pearson is exploring how this shift impacts our approach to assessment—ensuring our tools are adaptive and future-ready.
Balancing Vision and Capabilities	We deliver reliable solutions today while building toward the future of AI in education.

The future of *i-Ready* Assessment is invisible.

- Voice technology is coming to *i-Ready* Literacy Tasks
- Built to hear students' voices of all accents and dialects
- Creating the best possible solution by collaboratively learning with teachers in the classroom



AI Labs
Curriculum Associates



Learn more about our vision for the future

*ets research institute

Shaping the Future of AI in Assessment

ETS advances responsible AI research to promote fairness, trust, and innovation.

As AI transforms education, ETS brings decades of expertise to ensure that new solutions are not only powerful, but also valid, equitable, and transparent. Our work is driving the next-generation of measurement science, standing at the intersection of AI, learning, and assessment.

Highlights from ETS research at NCME AIME 2025:

- Investigating racial and ethnic subgroup representation in automated essay scoring
- Using generative AI teaching simulations to support teacher training
- Designing fairness-promoting, automated fraud detection systems
- Validating AI generated scoring rationales

REVIEW OUR
GUIDELINES FOR
RESPONSIBLE AI →



Advancing Assessment with AI

Grounded in science and responsible best practices, we use AI to enhance how we measure what students know and can do.

● **19 states**
we serve use hybrid scoring

24M essays & short answers
auto-scored by our AI engines

2M verbal responses
auto-scored by our AI engines

More AI-Powered Features - Coming Soon!

- WriteOn with Cambi
- Item Parameter Estimation
- Cheating Analysis
- Teacher Authoring with AI passage generation
- Hotline for student-at-risk work detection



Data reflects the 2024-2025 academic year



College Board Is a Proud Sponsor of AIME-Con

Join our engaging sessions to learn how we're advancing innovative and responsible use of AI in educational measurement.



edCount is pleased to sponsor 2025 NCME AIME-Con

*Over 20 years of service to
students and educators!*



Our Belief Statement

Every individual brings unique experiences, skill sets, and perspectives that work to advance our purpose: continuously improving the quality, fairness, and accessibility of education for all students.

Our Services

- Assessment Design, Development, and Evaluation
- Instructional Systems and Capacity Building
- Policy Analysis and Technical Assistance



www.edCount.com

(202) 895-1502 | info@edCount.com



www.NBME.org

ADVANCING ASSESSMENT, SUPPORTING OPTIMAL CARE

Through research and collaboration, NBME is evolving how we evaluate and support learners, with a focus on applying new technology to develop assessments that measure and build the knowledge and skills needed to provide optimal, effective care to all.



©2025 National Council on Measurement in Education (NCME)

Order copies of this and other NCME proceedings from:

National Council on Measurement in Education (NCME)
520 S. Walnut St. Box 2388
Bloomington, IN 47402
USA
Tel: +1-812-245-8096
ncme@ncme.org

ISBN 979-8-218-84229-1

Preface



Introduction

The inaugural NCME-sponsored Artificial Intelligence in Measurement and Education Conference (AIME-Con) brought together an interdisciplinary community of experts working at the intersection of artificial intelligence (AI), educational measurement, assessment, natural language processing, learning analytics, and technological development. As AI continues to transform education and assessment practices, this conference provided a critical platform for fostering cross-disciplinary dialogue, sharing cutting-edge research, and exploring the technical, ethical, and practical implications of AI-driven innovations in measurement and education. By bringing together experts from varied domains, the conference fostered a rich exchange of knowledge to enhance the collective understanding of AI's impact on educational measurement and evaluation.

Conference Theme - Innovation and Evidence: Shaping the Future of AI in Educational Measurement

The NCME-Sponsored AIME-Con focused on how rigorous measurement standards and innovative AI applications can work together to transform education. With sessions spanning summative large-scale assessment, formative classroom assessment, automated feedback, and informal learning tools, this conference fostered both the advancement and evaluation of AI technologies that are effective, reliable, and fair.

The National Council on Measurement in Education

The **National Council on Measurement in Education** is a community of measurement scientists and practitioners who work together to advance theory and applications of educational measurement to benefit society. A professional organization for individuals involved in assessment, evaluation, testing, and other aspects of educational measurement, our members are involved in the construction and use of standardized tests; new forms of assessment, including performance-based assessment; program design; and program evaluation. Learn more about NCME, including our goals and our leadership, at www.ncme.org. We are grateful to the NCME.

NCME Special Interest Group on Artificial Intelligence in Measurement and Education

The **AIME SIGIMIE** seeks to advance the theoretical and applied research into AI of educational measurement by bringing together data scientists, psychometricians, education researchers, and other interested stakeholders. The SIGIMIE will discuss current practices in using Generative AI, approaches to evaluate their precision/accuracy, and areas where more foundational research is required into the way we test and measure educational outcomes. This group seeks to create a strong professional identity and intellectual home for those interested in the use of AI in many areas, including automated scoring, item evaluation, validity studies, formative feedback, and generative AI for automated item generation.

Proposal Requirements and Review Process for Work-in-Progress Papers

AIME-Con invited submissions of Working Papers, which were submissions of up to 1,000 words in the form of a structured summary. This format was designed for work-in-progress or pilot studies. Working Papers required a title, short abstract, and followed a structured format with the following headings:

- Background
- Aims
- Sample(s)
- Methods
- Results (or Anticipated Results)
- Conclusions (or Anticipated Contributions)
- References, tables, and figures included as needed.

Submissions were evaluated by members of the review committee using a rubric that evaluated the following dimensions:

- **Relevance and community impact:** pertinence to the AI in measurement and education community, and potential contribution to current discussions and challenges in the field
- **Significance and value:** scholarly merit or practical importance of the work, and potential impact on theory, practice, or policy
- **Methodological rigor:** coherence and appropriateness of the proposed methods, techniques, and approaches; and soundness of the overall research design
- **Quality of expected outcomes:** whether the proposed analysis and interpretation methods are appropriate, and the potential contribution to knowledge in the field
- **Feasibility and timeline:** the realistic likelihood that the proposed work can be completed by the conference date

For the purposes of this conference, “AI” was defined broadly to include rule-based methods, machine learning, natural language processing, and generative AI/large language models. Reviewers provided constructive feedback and overall recommendations to ensure that accepted sessions reflected both scholarly merit and practical value to the AI in measurement and education community.

Organizing Committee

NCME Leadership

Amy Hendrickson, Ph.D. (President)
Rich Patz, Ph.D. (Executive Director)

Conference Chairs

Joshua Wilson, University of Delaware
Christopher Ormerod, Cambium Assessment
Magdalen Beiting Parrish, Federation of American Scientists

Proceedings Chair

Nitin Madnani, Duolingo

Proceedings Committee

Jill Burstein, Duolingo
Polina Harik, NBME

Program Committee

Conference Chairs

Joshua Wilson, University of Delaware
Christopher Ormerod, Cambium Assessment
Magdalen Beiting Parrish, Federation of American Scientists

Reviewers

Ketan , University of Massachusetts, Amherst
Hope Adegoke, University of North Carolina
Tazin Afrin, NBME
Ernest Amoateng, Western Michigan University
Kylie Anglin, University of Connecticut
Sergio Araneda, Caveon
Meirav Attali, Fordham University
Nurseit Baizhanov
Lee Becker, Pearson
Beata Beigman Klebanov, ETS
Ummugul Bezirhan, Boston College
Janet Shufor Bih Epse Fofang, University of Pittsburgh
Peter Bodary, University of Michigan School of Kinesiology
Brad Bolender, Finetune by Prometric
Jill Burstein, Duolingo
Hye-Jeong Choi, HumRRO
Jinmin Chung, Univ. of Iowa
Christina Cipriano, Yale University
Lisa Clark, City University of New York
Victoria Delaney, San Diego State University
Onur Demirkaya, Riverside Insights
Scott Elliot, SEG Measurement
Andrew Emerson, National Board of Medical Examiners
Mingyu Feng, WestEd
Taiwo Feyijimi, University of Georgia
Carla Firetto, Arizona State University
Jonathan Foster, University at Albany
Samantha Goldman, The University of Kansas
Chad Green, Loudoun County Public Schools
Joe Grochowalski, College Board
Yi Gui, The University of Iowa
Aysegul Gunduz, University of Alberta
Hongwen Guo, ETS Research Institute
Yage Guo, Center for Applied Linguistics
Gulsah Gurkan, Pearson
Suhwa Han, Cambium Assessment
Michael Hardy, Stanford University
Qiwei He, Georgetown University
Alexander Hoffman, AleDev Research & Consulting
Ruikun Hou, Technical University of Munich

Ruiping Huang, University of Illinois Chicago
Yue Huang, Measurement Incorporated
Hiu Ching Hung, Friedrich-Alexander-Universität Erlangen-Nürnberg
HUIMIN JIAO
Jamie Jirout, University of Virginia
Ji Yoon Jung, Boston College
Olasunkanmi Kehinde, Norfolk State University
YoungKoung Kim, The College Board
Becky King, University of Pittsburgh
Miryeong Koo, University of Illinois at Urbana-Champaign
Aakash Kumar, Texas A&M University
Alexander Kwako, Cambium Assessment
Brandon LeBeau, WestEd
Hansol Lee, Stanford University
Arun Balajiee Lekshmi Narayanan, University of Pittsburgh
Hongli Li, Georgia State University
Tianwen Li, University of Pittsburgh
Li Liang
Boyuan LIU, Department of Educational Psychology, The Chinese University of Hong Kong
Chen Liu, UC Merced
Will Lorie
Susan Lottridge, Cambium Assessment
Max Lu, Harvard University
Yi Lu, Federation of State Boards of Physical Therapy
Wenchao Ma, University of Minnesota
Henry Makinde, University of North Carolina - Greensboro
Mike Maksimchuk, Kent Intermediate School District
Salih Mansur, Touro University of New York
Jamie Mikeska, ETS
Mubarak Mojoyinola, University of Iowa
Wesley Morris, Vanderbilt University
Tim Moses, Buros Center for Testing
William Muntean, National Council of State Boards of Nursing
Mariel Musso, University of Granada- CONICET
Supraja Narayanaswamy, Acelero Inc.
Lynn Nguyen, Fruits eTutoring
Tram-Anh Tran Nguyen, University of Massachusetts, Amherst
Chunling Niu, The University of the Incarnate Word
Kai North, Cambium Learning Group, Inc.
Teresa Ober, ETS
Maria Oliveri, Purdue University
Christopher Ormerod, Cambium Assessment
Jay Parkes, University of New Mexico
Hallie Parten, University of Virginia
Katie Pedley, Pearson
Benjamin Pierce, University of Pittsburgh
Andrew Potter, Arizona State University
Sonya Powers, Edmentum
Ricardo Primi, Universidade São Francisco
Sarah Quesen, WestEd
Ruchi Sachdeva, Pearson

Fariha Hayat Salman, American University in Dubai
Lydia Scholle-Cotton, Queen's University (Kingston, ON, Canada)
Qingzhou Shi, Northwestern University
Jinnie Shin, University of Florida
Anthony Shiver, Law School Admission Council
Stephen Sireci, University of Massachusetts Amherst
Anastasia Smirnova, San Francisco State University
Xiaomei Song, Case Western Reserve University School of Medicine
Kayden Stockdale, Virginia Tech
Caitlin Tenison, ETS
Danielle Thomas, Carnegie Mellon University
Zewei Tian, University of Washington
Nhat Tran, University of Pittsburgh
FELIPE Valentini, Graduate School of Psychology, Universidade São Francisco
Marcus Walker, National Commission on Certification of Physician Assistants
Cole Walsh, Acuity Insights
Huanxiao Wang, University of Pennsylvania
Yun-Han Weng, Ohio State University
Joshua Wilson, University of Delaware
Sirui Wu, University of British Columbia
Hyesun You, University of Iowa
Meltem Yumsek Akbaba, Ministry of National Education, Turkey
Diego Zapata-Rivera, ETS
Dake Zhang, Rutgers University
Jiayi (Joyce) Zhang, University of Pennsylvania
Liang Zhang, University of Georgia
Ting Zhang, American Institutes for Research
Lauren Zito, WGU Labs

Table of Contents

<i>Automated Item Neutralization for Non-Cognitive Scales: A Large Language Model Approach to Reducing Social-Desirability Bias</i> Sirui Wu and Daijin Yang	1
<i>AI as a Mind Partner: Cognitive Impact in Pakistan’s Educational Landscape</i> Eman Khalid, Hammad Javaid, Yashal Waseem and Natasha Sohail Barlas	14
<i>Detecting Math Misconceptions: An AI Benchmark Dataset</i> Bethany Rittle-Johnson, Rebecca Adler, Kelley Durkin, L Burleigh, Jules King and Scott Crossley	20
<i>Optimizing Opportunity: An AI-Driven Approach to Redistricting for Fairer School Funding</i> Jordan Abbott	25
<i>Automatic Grading of Student Work Using Simulated Rubric-Based Data and GenAI Models</i> Yiyao Yang and Yasemin Gulbahar	34
<i>Cognitive Engagement in GenAI Tutor Conversations: At-scale Measurement and Impact on Learning</i> Kodi Weatherholtz, Kelli Millwood Hill, Kristen Dicerbo, Walt Wells, Phillip Grimaldi, Maya Miller-Vedam, Charles Hogg and Bogdan Yamkovenko	40
<i>Chain-of-Thought Prompting for Automated Evaluation of Revision Patterns in Young Student Writing</i> Tianwen Li, Michelle Hong, Lindsay Clare Matsumura, Elaine Lin Wang, Diane Litman, Zhexiong Liu and Richard Correnti	49
<i>Predicting and Evaluating Item Responses Using Machine Learning, Text Embeddings, and LLMs</i> Evelyn Johnson, Hsin-Ro Wei, Tong Wu and Huan Liu	66
<i>Evaluating LLM-Based Automated Essay Scoring: Accuracy, Fairness, and Validity</i> Yue Huang and Joshua Wilson	71
<i>Comparing AI tools and Human Raters in Predicting Reading Item Difficulty</i> Hongli Li, Roula Aldib, Chad Marchong and Kevin Fan	84
<i>When Machines Mislead: Human Review of Erroneous AI Cheating Signals</i> William Belzak, Chenhao Niu and Angel Ortmann Lee	90
<i>Fairness in Formative AI: Cognitive Complexity in Chatbot Questions Across Research Topics</i> Alexandra Barry Colbert and Karen D Wang	98
<i>Keystroke Analysis in Digital Test Security: AI Approaches for Copy-Typing Detection and Cheating Ring Identification</i> Chenhao Niu, Yong-Siang Shih, Manqian Liao, Ruidong Liu and Angel Ortmann Lee	107
<i>Talking to Learn: A SoTL Study of Generative AI-Facilitated Feynman Reviews</i> Madeline Rose Mattox, Natalie Hutchins and Jamie J Jirout	117
<i>AI-Powered Coding of Elementary Students’ Small-Group Discussions about Text</i> Carla Firetto, P. Karen Murphy, Lin Yan and Yue Tang	125
<i>Evaluating the Reliability of Human–AI Collaborative Scoring of Written Arguments Using Rational Force Model</i> Noriko Takahashi, Abraham Onuorah, Alina Reznitskaya, Evgeny Chukharev, Ariel Sykes, Michele Flammia and Joe Oyler	135

<i>Evaluating Deep Learning and Transformer Models on SME and GenAI Items</i> Joe Betts and William Muntean	141
<i>Comparison of AI and Human Scoring on A Visual Arts Assessment</i> Ning Jiang, Yue Huang and Jie Chen	147
<i>Explainable Writing Scores via Fine-grained, LLM-Generated Features</i> James V Bruno and Lee Becker	155
<i>Validating Generative AI Scoring of Constructed Responses with Cognitive Diagnosis</i> Hyunjoo Kim	166
<i>Automated Diagnosis of Students' Number Line Strategies for Fractions</i> Zhizhi Wang, Dake Zhang, Min Li and Yuhan Tao	178
<i>Medical Item Difficulty Prediction Using Machine Learning</i> Hope Oluwaseun Adegoke, Ying Du and Andrew Dwyer	185
<i>Examining decoding items using engine transcriptions and scoring in early literacy assessment</i> Zachary Schultz, Mackenzie Young, Debbie Dugdale and Susan Lottridge	191
<i>Addressing Few-Shot LLM Classification Instability Through Explanation-Augmented Distillation</i> William Muntean and Joe Betts	197
<i>Identifying Biases in Large Language Model Assessment of Linguistically Diverse Texts</i> Lionel Hsien Meng, Shamyia Karumbaiah, Vivek Saravanan and Daniel Bolt	204
<i>Implicit Biases in Large Vision–Language Models in Classroom Contexts</i> Peter Baldwin	211
<i>Enhancing Item Difficulty Prediction in Large-scale Assessment with Large Language Model</i> Mubarak Mojinyinola, Olasunkanmi James Kehinde and Judy Tang	218
<i>Leveraging LLMs for Cognitive Skill Mapping in TIMSS Mathematics Assessment</i> Ruchi J Sachdeva and Jung Yeon Park	223