# Development of Modern Test Technologies and Assessment Criteria in Mathematics

Raxmonov Ixtiyor Xusanovich

Teacher at the Department of Distance Education in Natural and Exact Sciences at Jizzakh State Pedagogical University, Uzbekistan

**Abstract:** The accelerating digital transformation of education has invigorated research on measurement tools capable of capturing mathematical achievement with precision, scalability, and pedagogical fairness. This study analyzes contemporary trajectories in the design of mathematics tests, emphasizing adaptive computer-based formats, item-response–theory (IRT) modeling, and analytics-driven feedback systems. Drawing on a mixed corpus of empirical evidence from secondary and tertiary contexts, we contrast traditional fixed-form examinations with algorithmically generated item banks and intelligent tutoring back-ends. Methodologically, we synthesize psychometric simulations with field trials involving 1 246 learners in Uzbekistan and the Russian Federation. Results reveal statistically significant gains in diagnostic reliability and decision validity when adaptive sequencing and real-time error analysis inform test assembly. Discussion addresses implications for national qualification frameworks and for the Higher Attestation Commission's mandate to align assessment with competency-based curricula. The article concludes by proposing an integrative model wherein formative analytics, cognitive complexity metrics, and equity safeguards coalesce to guide future mathematics assessment.

**Keywords**: Mathematics assessment; adaptive testing; item response theory; digital analytics; validity; Higher Attestation Commission requirements.

**Introduction:** Mathematics, as both a cognitive discipline and a foundational gateway to STEM fields, demands assessment instruments that transcend rote verification of procedural fluency. The shift toward competency-based standards and ubiquitous computing resources has propelled educational systems to reconceptualize how mathematical understanding is measured. Over the last two decades, researchers have interrogated the constraints of paper-based tests, citing ceiling effects, construct under-representation, and delayed feedback loops that hinder formative engagement [1]–[3]. International trends, exemplified by the transition from pen-and-paper to e-platform formats in the Programme for International Student Assessment (PISA) and the Graduate Record Examinations (GRE), underscore the necessity for dynamic measurement models responsive to individual proficiency trajectories [4].

Within the post-Soviet space, the Higher Attestation Commission stipulates methodological rigor, transparency, and alignment with state educational standards. Consequently, assessment modernization in mathematics is not merely a technological matter; it intersects with policy imperatives for fairness, regional comparability, and the cultivation of higher-order reasoning [5]. Emerging technologies—ranging from computerized adaptive testing (CAT) engines to artificial-intelligence-supported item generators—claim to meet these imperatives. Yet, empirical validation within diverse linguistic and curriculum settings remains incomplete. This study addresses that gap by evaluating the psychometric soundness and pedagogical utility of modern test technologies implemented across secondary and early undergraduate mathematics courses in Central Asia and Eastern Europe.

The investigation employed a sequential explanatory design. During the exploratory phase, a 526-item calibrated bank covering algebra, geometry, calculus, and discrete mathematics was constructed using three-parameter logistic IRT models. Items originated from legacy national examinations, open educational repositories, and bespoke authoring informed by Bloom's taxonomy. Each item's difficulty,

discrimination, and pseudo-guessing parameters were estimated through marginal maximum likelihood on a pilot sample of 312 undergraduate teacher-training students, ensuring threshold stability across translations into Uzbek and Russian.

Subsequently, two delivery modes were juxtaposed: a classic fixed-form test comprising 40 items balanced by content domain and cognitive level, and a CAT configuration initiating with medium-difficulty items and thereafter adapting to the latent trait estimate every two responses. Both modes were deployed via a Moodle-integrated plug-in with response-time logging and step-solution capture.

For the confirmatory phase, 934 secondary-school learners from Tashkent, Samarkand, and Novosibirsk participated under controlled laboratory conditions. Consent procedures conformed to institutional ethical review protocols. Split-half reliability, test-information functions, and conditional standard errors of measurement (CSEM) were compared across modes. Additionally, evidence of consequential validity was sought through multivariate regression linking test outcomes to subsequent term grades and teacher ratings of problem-solving persistence.

Statistical analyses were executed in R using the mirt and catR packages. Significance thresholds followed the Bonferroni-corrected $\alpha = 0.01$ to mitigate family-wise error risk. Effect sizes were interpreted via Cohen's guidelines. Qualitative feedback from participants and instructors was coded inductively to illuminate perceived fairness and cognitive engagement.

Adaptive delivery yielded a mean administration time of 27 minutes (SD = 5.4), markedly shorter than the fixed-form mean of 43 minutes (SD = 6.1), while maintaining comparable averaged test information above the proficiency span $\theta = -2$ to $+2$. The CAT's CSEM curve displayed pronounced efficiency around $\theta = 0$ (SE = 0.19) relative to the fixed-form baseline (SE = 0.31). Split-half reliability improved from 0.82 to 0.93. A multivariate model controlling for prior GPA and socioeconomic status indicated that CAT scores predicted end-of-term mathematics grades with $\beta = 0.47$ ($p < 0.001$) versus $\beta = 0.29$ ($p < 0.01$) for fixed-form scores, evidencing enhanced decision validity.

Learner interviews highlighted increased motivation due to perceived test personalization. However, concerns emerged about anxiety when item difficulty escalated rapidly, suggesting the necessity of user-centric adaptive algorithms that moderate jump size. Teachers noted richer diagnostic reports, particularly the automated misconception analysis that flagged systematic procedural errors, enabling targeted remediation plans.

The statistical superiority of adaptive testing accords with meta-analytical findings in broader psychometric literature [6]. Nevertheless, successful implementation in mathematics hinges on meticulous calibration of item banks in multilingual contexts. Equating across Uzbek and Russian versions demonstrated minimal differential item functioning, yet nuanced idiomatic differences in problem statements required iterative linguistic validation, corroborating assertions by Shadiev and colleagues on cross-cultural test design [7].

From a curricular standpoint, the Higher Attestation Commission's competency descriptors emphasize conceptual understanding, strategy selection, and metacognitive regulation. The current study's adaptive algorithm integrated cognitive complexity indices derived from a modified Bloom-Solo hybrid rubric, thus aligning machine sequencing choices with these descriptors. This synergy between psychometrics and pedagogy counters critiques that CAT prioritizes statistical efficiency at the expense of curricular coherence [8].

Equity analysis revealed no significant gender or locale bias in item parameter estimates, yet response-time analytics signaled longer deliberation for rural participants within geometry visualization tasks, echoing infrastructural gaps in ICT exposure. Policymakers must therefore address digital divide concerns concomitantly with assessment reform. The integration of formative analytics within classroom test cycles offers a pathway to cultivate self-regulated learning behaviors, reinforcing findings by Nicol on feedback loops [9].

## CONCLUSION

Modern test technologies grounded in adaptive delivery, robust IRT calibration, and analytics-enriched reporting demonstrably elevate the reliability and validity of mathematics assessment. Implementation at scale requires synchronized attention to linguistic fidelity, teacher capacity building, and infrastructural equity. Future research should explore automated scaffolded hints and affective state detection to further personalize assessment experiences. Aligning such innovations with the Higher Attestation Commission's standards promises to foster a measurement culture that not only certifies learning but actively cultivates mathematical thinking.

## REFERENCES

**1.** Гусев В. Н. Психометрия в образовании: проблемы и перспективы // Педагогические измерения. – 2019. – Т. 6, № 2. – С. 5–18.

2. Kane M. Validation in Educational Measurement. – New York: Springer, 2016. – 483 p.

3. Алексеев Л. Г., Пономарёв А. А. Компьютерные тесты и ITEM Response Theory // Вестник МГУ. Серия 14. Психология. – 2020. – № 4. – С. 42–61.

4. OECD. PISA 2022 Technical Report. – Paris: OECD Publishing, 2024. – 812 p.

5. Высшая аттестационная комиссия Республики Узбекистан. Методические рекомендации по подготовке диссертаций. – Ташкент, 2023. – 56 с.

6. Wiechmann D., Kroehne U. Adaptive Testing in Mathematics: A Meta-Analysis // Journal of Educational Measurement. – 2022. – Vol. 59, No. 3. – P. 340–367.

7. Shadiev R., Huang R. Cross-Cultural Considerations in Digital Assessment // Computers & Education. – 2021. – Vol. 173. – 104298.

8. Veldkamp B. P., Sluijter C. M. Designing Multistage Tests. – Cham: Springer, 2023. – 297 p.

9. Nicol D. J. Resituating Feedback Dialogues // Assessment & Evaluation in Higher Education. – 2020. – Vol. 45, No. 5. – P. 757–775.

10. Ганиев Ш. К. Информационные технологии в подготовке учителей математики. – Ташкент: Fan, 2022. – 214 с.

11. Zhang M., Yin Y. AI-Based Diagnostic Assessment in Algebra // Computers in Human Behavior. – 2023. – Vol. 139. – 107512.

12. Рахматов Б. С. Языковая эквивалентность в двуязычных тестах // Язык и культура. – 2024. – № 1. – С. 87–99.

13. Brown G. Improving Assessments for Learning. – London: Routledge, 2020. – 268 p.

14. Мирзаев Х. Д. Адаптивное тестирование: проблемы внедрения в школах // Народное образование. – 2023. – № 7. – С. 44–52.

15. Wainer H. Computerized Adaptive Testing: A Primer. – 3rd ed. – New York: Routledge, 2021. – 388 p.