Beyond Human Moderation: The Case for Automated Al Validation in Educational Assessment

Executive Summary

As artificial intelligence transforms educational assessment, the critical question is not whether AI can mark student work, but who—or what—ensures AI marking is trustworthy, fair, and accountable. While current practice relies on human moderation layers to oversee AI decisions, this white paper proposes a more systematic approach: automated AI validation layers that continuously benchmark AI performance against expert human standards.

This paper examines the limitations of existing human moderation approaches and presents the case for automated validation systems that offer greater consistency, transparency, and scalability while maintaining human expertise at the core of assessment standards.

Key Lessons from the RM Compare AI Validation Layer POC

- Small, expert consensus sets established through RM Compare are highly effective as gold standards for calibrating AI assessment, even across diverse subjects and response types.
- Iterative validation cycles—with continuous comparison and parameter adjustment—consistently improve AI scoring accuracy, fairness, and reliability.
- Real-time discrepancy analysis ensures that "difficult to score" items always receive human attention, further strengthening trust and system robustness.
- Scalable, proactive monitoring supports thousands of assessments with limited ongoing human resource, delivering both efficiency and quality.
- Adaptive consensus-building means the gold standard itself matures over time, maintaining relevance as curriculum, marking practice, or student performance evolve.
- The combination of human-anchored and AI-validated workflow is key: rather than replacing expertise, RM Compare amplifies it, resulting in fundamentally stronger trust and adoption across educators and awarding bodies.

1. Introduction: The Trust Challenge in AI Assessment

Educational institutions worldwide are rapidly adopting AI for assessment, driven by promises of efficiency, consistency, and reduced workload. Research shows that 86% of education organisations now use generative AI, with student adoption rising by 26 percentage points in just one year. However, this rapid adoption has exposed a fundamental challenge: ensuring AI decisions are fair, accurate, and trustworthy.

Current regulatory frameworks, including guidance from Ofqual in the UK, mandate that "a human assessor must review all the work in its entirety and determine the mark it warrants, regardless of the outcomes of an AI tool". This human-in-the-loop approach represents the dominant model for Al marking oversight, but it may not be sufficient for the scale and complexity of modern educational assessment.

Blog: Who is Assessing the AI that is Assessing Students?

2. The Current State: Human Moderation Layers How Human Moderation Works

The prevailing approach to AI assessment oversight involves human moderators who:

- Sample a subset of Al-marked scripts for review
- Check for obvious errors, inconsistencies, or bias 🚣
- Adjust marks where significant discrepancies occur 📑
- Provide anecdotal feedback on AI performance

Limitations of Human Moderation

While human oversight is essential, this approach faces several critical limitations:

Scale and Consistency Challenges: Human moderation struggles with large datasets, often reviewing only small samples while systematic issues may remain undetected across thousands of unmarked scripts.

Reactive Rather Than Proactive: Human moderation typically identifies problems after marking is complete, requiring time-consuming corrections rather than preventing errors from occurring.

Variable Standards: Different human moderators may apply inconsistent criteria, introducing the very subjectivity Al was intended to eliminate.

Resource Intensity: As assessment volumes grow, maintaining comprehensive human oversight becomes increasingly expensive and logistically challenging. \$

Limited Bias Detection: Subtle algorithmic bias affecting specific student populations may be difficult for human reviewers to identify, particularly when reviewing small samples.

3. The Alternative: Automated Al Validation Layers

Defining Automated Validation

An automated AI validation layer is a systematic framework that continuously benchmarks AI assessment decisions against a gold standard of expert human judgments. Rather than replacing human expertise, it amplifies and systematises it, ensuring every AI decision meets predetermined accuracy and fairness thresholds.

Core Components of Validation Systems

Gold Standard Benchmarking: A carefully curated set of assessment items marked by expert humans serves as the permanent reference point for Al calibration. 📊

Continuous Monitoring: Every AI-generated score is automatically compared against expected human performance, with discrepancies immediately flagged. 👀

Systematic Correction: When errors are detected, the system initiates structured correction protocols, adjusting AI parameters and reprocessing affected assessments.



Transparency and Auditability: All validation decisions are logged, creating a comprehensive audit trail that demonstrates fairness and accountability.

RM Compare – Producing Gold Standard Benchmarking

At the heart of trustworthy automated assessment lies the concept of a "gold standard"—a reference set of marks that can be relied upon as a definitive measure of quality. RM Compare is uniquely placed to provide this benchmark through its worldleading Adaptive Comparative Judgement (ACJ) system. Rather than relying on singlerater scores or rigid marking rubrics, RM Compare gathers judgements from multiple expert assessors, using a dynamic algorithm to compare pairs of student work. By iteratively building consensus, RM Compare creates a statistically robust, representative rank order that authentically reflects collective professional expertise and mitigates personal bias.

This gold standard is not only reliable but also adaptable across a wide range of assessment domains, including essays, creative outputs, oracy, and more. As a result, the RM Compare benchmark can be used to calibrate and validate Al-based marking at scale, ensuring that every automated decision is continually anchored to proven human standards. By acting as the keystone for AI validation layers, RM Compare underpins fairness, transparency, and trust in next-generation educational assessment.

Blog: Building Trust: From "Ranks to Rulers" to On-Demand Marking

4. Comparative Analysis: Human Moderation vs. Automated Validation

Aspect	Human Moderation	Automated Validation
Coverage	Sample-based, partial	Complete dataset coverage
Consistency	Variable across moderators	Uniform application of standards
Error Detection	Reactive, post-marking	Proactive, real-time
Bias Identification	Limited, subjective	Systematic, quantifiable
Scalability	Resource-constrained	Highly scalable
Transparency	Anecdotal reporting	Comprehensive audit trails
Cost Efficiency	High ongoing costs	Lower marginal costs at scale



Figure 1: Automated Validation Process

5. Benefits of Automated Al Validation

Enhanced Fairness and Bias Reduction

Automated validation systems excel at detecting subtle patterns of bias that human reviewers might miss. By systematically comparing AI decisions across different student populations and question types, these systems can identify and correct discriminatory patterns before they affect student outcomes.

Research indicates that students increasingly perceive AI-driven assessment as fairer than human-only evaluation, particularly when the underlying processes are transparent. Automated validation enhances this perception by providing demonstrable evidence of systematic fairness checks.

Improved Efficiency and Resource Allocation

While establishing a gold standard requires initial human input—often as few as 100 expertly marked items—this investment provides perpetual calibration value. Once established, the validation layer can process unlimited volumes of assessment without proportional increases in human oversight requirements.

This efficiency gain is particularly valuable given research showing that AI can reduce marking time and costs by up to 60%, while automated validation ensures these efficiency gains don't compromise quality or fairness.

Transparency and Trust Building

Automated validation transforms assessment from a "black box" process into a transparent, auditable system. Every validation decision is recorded, creating a comprehensive trail that demonstrates:

- How AI decisions align with human expertise
- Where corrections were made and why
- The ongoing accuracy and fairness of the system

This transparency directly addresses stakeholder concerns about AI accountability, providing evidence-based confidence rather than requests for trust.

Continuous Improvement and Learning

Unlike static human moderation processes, automated validation enables continuous system refinement. Each validation cycle provides data on AI performance, enabling systematic improvements to accuracy, fairness, and consistency over time.

6. Addressing Concerns and Limitations

The Human Element 🐵

Automated validation doesn't eliminate human expertise—it systematises and amplifies it. Human experts remain essential for:

- Establishing initial gold standards
- Defining fairness criteria and assessment objectives
- Reviewing complex edge cases
- Making policy decisions about acceptable performance thresholds

Technical Implementation Challenges 🧩

Implementing automated validation requires careful attention to:

- Data Quality: Gold standard items must be representative and expertly assessed. 🤽
- Algorithm Transparency: Validation processes should be explainable to stakeholders
- System Integration: Validation must integrate seamlessly with existing assessment workflows
- Regular Calibration: Gold standards may need periodic review and updating

Regulatory and Policy Considerations

Current regulatory frameworks emphasise human oversight, and automated validation systems must be designed to complement rather than replace regulatory requirements. This may involve:

- Demonstrating equivalence or superiority to human moderation
- Providing audit trails that satisfy regulatory scrutiny
- Maintaining human accountability for final assessment decisions

7. Implementation Roadmap

Phase 1: Proof of Concept

- Establish gold standard datasets for specific assessment types
- Develop validation algorithms and metrics <
- Conduct pilot testing with limited scope 6

Phase 2: Systematic Deployment

- Scale validation systems across broader assessment programs 🕐
- Integrate with existing assessment platforms
- Train stakeholders on validation outputs and interpretation (

Phase 3: Continuous Improvement

- Analyse validation data for system refinement
- Expand to new assessment types and contexts
- Share best practices and standards across institutions

Blog: Fairness in Focus: The Al Validation Layer Proof of Concept Powered by RM

Compare

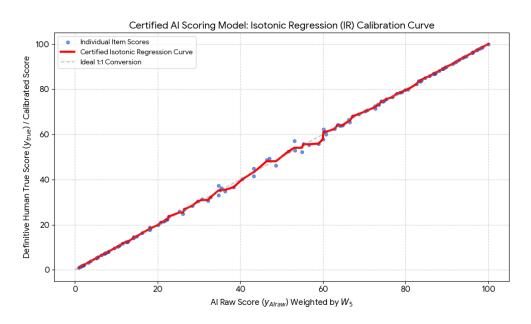


Figure 2: POC results

8. RM Compare AI Validation Layer POC Results and Learnings

The proof of concept (POC) for the RM Compare AI Validation Layer demonstrates how consensus-driven benchmarking delivers robust, reliable, and future-ready validation for automated assessment. By leveraging RM Compare's Adaptive Comparative Judgement (ACJ) platform, the POC established a gold standard not through static rubrics or single-rater judgment, but by synthesising expert consensus across multiple markers and diverse assessment types.

During calibration, the AI model marked a set of training items and compared its results to the established gold standard. This process was visualised through output curves that indicated the AI's ability to closely align with human judgement. Discrepancies were systematically reviewed, leading to iterative improvements in the AI's scoring reliability and fairness.

Operationally, the validation system is proactive: every new item is checked against the gold standard. Where the automated result matches the benchmark, the mark is awarded automatically. Where discrepancies arise, scripts are flagged for human review—ensuring the highest standards of fairness, and using each challenging case to further refine the system's accuracy. This approach also allows the model to remain adaptive, improving with each cycle rather than becoming rigid or outdated.

Key learnings from the POC include:

- The effectiveness of RM Compare's ACJ approach in creating a transparent, defensible reference standard for calibration.
- The importance of continuous monitoring and auditing, supporting scalable and consistent validation for thousands of assessments over time.
- Evidence that this combined human-anchored, AI-validated workflow fundamentally strengthens trust, scalability, and fairness in educational assessment.

Together, these results highlight why RM Compare serves as the keystone of this Al validation framework—ensuring every automated mark is fair, accountable, and upholds the confidence of educators and learners alike.

9. Future Directions and Research Opportunities

Collaborative Standards Development

The education sector would benefit from collaborative development of validation standards, ensuring interoperability and shared best practices across institutions and platforms.

Integration with Emerging Technologies

As AI assessment capabilities evolve, validation systems must adapt to handle:

- Multimodal assessments (text, audio, video)
- Complex reasoning and creativity evaluation
- Personalized and adaptive assessment formats

Long-term Impact Studies

Comprehensive research is needed to evaluate the long-term impact of automated validation on:

- Student outcomes and satisfaction
- Teacher workload and effectiveness
- System-wide assessment quality and fairness

10. Conclusion: A Call for Systematic Validation

The rapid adoption of AI in educational assessment demands equally sophisticated approaches to ensuring fairness, accuracy, and accountability. While human oversight remains essential, purely manual moderation approaches are insufficient for the scale and complexity of modern assessment systems.

Automated AI validation layers offer a path forward that maintains human expertise at the core while providing systematic, scalable, and transparent oversight of AI decisions. By continuously benchmarking AI performance against expert human standards, these systems can deliver the efficiency benefits of AI assessment while building demonstrable trust among all stakeholders.

The question is not whether we need better oversight of AI assessment—the question is whether we will develop systematic, evidence-based approaches to validation or continue to rely on resource-intensive, limited-scope human moderation.

Educational institutions, technology providers, and policymakers must collaborate to develop and implement automated validation standards that serve the needs of students, educators, and society. The future of fair, trustworthy AI assessment depends on our willingness to move beyond traditional approaches and embrace systematic validation as a core component of educational technology.

11. Call to Action

We invite stakeholders across the education sector to engage in this critical conversation:

- Educators and Administrators: Consider how automated validation could enhance fairness and transparency in your assessment systems
- Technology Providers: Invest in developing robust validation capabilities as core platform features
- Policymakers: Support research and standards development for AI validation in education
- Researchers: Contribute to the evidence base for validation effectiveness and best practices

The transformation of educational assessment is underway. By working together to develop systematic validation approaches, we can ensure this transformation serves the needs of all learners while maintaining the highest standards of fairness and accountability.

Want to learn more, contribute, or see a demonstration? We welcome educators, examiners, and policymakers to engage with us as we help shape the next era of fair and scalable assessment. Contact is here – <u>compare.rm.com</u>.