

# TEAS: Trusted Educational AI Standard

## A Framework for Verifiable, Stable, Auditable, and Pedagogically Sound Learning Systems

Abu Syed

Founder and CEO, Metacog  
Student, Data Science and Applications  
Indian Institute of Technology, Madras

### Abstract

The rapid integration of AI into education has prioritized capability over trustworthiness, creating significant risks. Real-world deployments reveal that even advanced models are insufficient without extensive architectural scaffolding to ensure reliability. Current evaluation frameworks are fragmented: institutional policies lack technical verification, pedagogical guidelines assume AI reliability, and technical metrics are context-agnostic. This leaves institutions without a unified standard for deployment readiness. This paper introduces TEAS (Trusted Educational AI Standard), an integrated framework built on four interdependent pillars: (1) Verifiability, grounding content in authoritative sources; (2) Stability, ensuring deterministic core knowledge; (3) Auditability, enabling independent institutional validation; and (4) Pedagogical Soundness, enforcing principles of active learning. We argue that trustworthiness stems primarily from systematic architecture, not raw model capability. This insight implies that affordable, open-source models can achieve deployment-grade trust, offering a scalable and equitable path to integrating AI safely into learning environments globally.

### 1 Introduction

The rapid integration of artificial intelligence (AI) into educational settings promises unprecedented personalization and scale. But as with any transformative technology, this warrants critical examination. In much of the discourse around AI, capability is assumed to be sufficient. Yet raw capability is not enough in high-stakes contexts like education. As educator and technologist Jim Chilton observes, when AI makes an error, “students won’t know the difference—because the AI itself doesn’t know the difference” (Chilton 2025). Unlike domain experts who reason from understanding, large language models predict “the next best word, not truth” (Chilton 2025), a fundamental limitation that poses serious risks when deployed in learning contexts where accuracy and pedagogical integrity are paramount.

This challenge is not theoretical. Khan Academy, one of the world’s most respected educational organizations, discovered that deploying GPT-4 as a tutoring system required building extensive architectural scaffolding around

the base model: a separate calculator tool to handle numerical computations the language model could not reliably perform, visual preprocessors to generate textual representations of mathematical graphics, and multi-path reasoning engines to follow students’ unconventional solution strategies (DiCerbo 2025). The lesson is stark: the base model alone, despite its impressive capabilities, was insufficient for deployment. Systematic architecture—not raw model capability—determined whether the system could be trusted with students.

This pattern extends far beyond a single deployment. A growing body of research documents distinct categories of failure in educational AI systems, from confident hallucinations in STEM subjects (Delikoura, Fung, and Hui 2025) to “vaporized learning” where short-term test scores improve while long-term retention deteriorates (Delikoura, Fung, and Hui 2025), from non-deterministic outputs that make curriculum standardization impossible (Shin et al. 2025) to black-box opacity that prevents institutional auditing. The foundational question remains inadequately addressed: what does it mean for an educational AI system to be trustworthy enough for deployment?

Current evaluation frameworks do not fully address this question in an integrated manner. Different stakeholder groups have developed domain-specific approaches: administrators rely on institutional policy checklists that guide procurement processes but lack technical verification mechanisms (1EdTech n.d.-a); educators use pedagogical frameworks that articulate sound teaching principles but assume AI reliability and offer no enforcement for black-box systems (UNESCO 2023; Rahimi and Shute 2024); and engineers apply technical metrics that quantify specific failure modes but remain domain-agnostic, capturing neither educational context nor pedagogical quality (Ji et al. 2023). Vendor disclosures provide transparency but rely on self-attestation rather than independent validation (OpenAI 2024; Google 2024). This fragmentation leaves institutions without a unified standard for making deployment decisions that account for the full spectrum of trustworthiness requirements.

This paper proposes **TEAS (Trusted Educational AI Standard)**, an integrated framework that addresses this gap. TEAS establishes four interdependent pillars for evaluating deployment readiness: **Verifiability** (the ability to

trace AI-generated content to authoritative sources), **Stability** (deterministic consistency for core curriculum knowledge), **Auditability** (institutional capacity to independently validate the system’s knowledge and logic), and **Pedagogical Soundness** (adherence to evidence-based teaching methods that foster active learning). Unlike existing approaches, TEAS recognizes that partial compliance is insufficient—a system that is verifiable but pedagogically harmful remains untrustworthy, just as one that is pedagogically sound but factually unstable cannot be reliably deployed.

Critically, TEAS demonstrates that trustworthiness stems from systematic architectural design rather than raw model capability. The Khanmigo case exemplifies this principle: safety emerged not from using GPT-4 versus a smaller model, but from the deliberate engineering of verification layers, domain-specific tools, and pedagogical guardrails. This insight has profound implications for educational equity. If systematic frameworks—rather than expensive frontier models—determine trustworthiness, then affordable open-source models can achieve the deployment-grade trust currently associated only with commercial systems. This makes AI tutoring economically viable at unprecedented scale, particularly in resource-constrained contexts where current commercial pricing models remain prohibitive (Ravindran 2025). We validate this claim empirically in Appendix A, where a knowledge-grounded 8B parameter model outperforms ungrounded frontier models on TEAS criteria.

The remainder of this paper proceeds as follows. Section 2 presents a taxonomy of documented failures in educational AI and analyzes real-world deployment challenges. Section 3 reviews existing evaluation frameworks and demonstrates the gap that TEAS addresses. Section 4 presents the TEAS framework, establishing the four pillars that together define deployment readiness and examining how their integration addresses the failure modes identified in current systems. Section 5 discusses implications for researchers, institutions, developers, and policymakers, with particular attention to deployment economics and educational equity. Section 6 acknowledges limitations and outlines future work. Appendix A presents an empirical case study demonstrating TEAS in practice.

## 2 Problem Landscape

The deployment of AI in education faces a complex web of interconnected failure modes that threaten both the efficacy and integrity of learning systems. This section synthesizes documented evidence of these failures and examines real-world cases that illuminate the challenges facing even well-resourced development efforts.

### 2.1 Taxonomy of Documented Failures

**Factual and Reasoning Failures.** The most widely recognized limitation of large language models is their propensity to generate confident-sounding misinformation, a phenomenon commonly termed “hallucination” (Delikoura, Fung, and Hui 2025). This is not a peripheral flaw but a direct consequence of their probabilistic architecture: models predict the next most likely token without grounding

in factual reality (Chilton 2025). The problem extends beyond simple factual errors. Research on mathematical reasoning reveals that models produce “sneaky errors”—subtle, multi-step logical flaws that are difficult to diagnose because mainstream training objectives prioritize generating correct final answers over exposure to diverse error patterns (Zou et al. 2025). In educational contexts where precision is paramount, these failures introduce persistent misinformation that directly undermines knowledge acquisition.

**Pedagogical Misalignment.** Even when providing factually correct information, AI systems can undermine learning through pedagogically harmful interaction patterns. A systematic review identifies “cognitive offloading,” reduced neural activity during learning tasks, diminished independent learning skills, and overall loss of student agency as documented outcomes of AI interaction (Delikoura, Fung, and Hui 2025). This phenomenon has been termed “vaporized learning”: AI tools may boost short-term performance metrics like test scores while eroding long-term retention and genuine understanding (Delikoura, Fung, and Hui 2025). Comparative studies reveal the mechanism: while human tutors naturally employ Socratic questioning patterns that stimulate active thinking, AI systems default to passive “explanation-simplistic response” loops where the model delivers information and the student provides minimal engagement (Delikoura, Fung, and Hui 2025). When commercial incentives prioritize engagement metrics over learning outcomes, this misalignment can become systematic rather than incidental, as evidenced by user analyses of platforms like Duolingo (Lee 2024).

**Invisible Errors and Trust Miscalibration.** A critical dimension of AI failure in education is that students frequently cannot recognize when mistakes occur. Research on human-AI interaction demonstrates that users are poor at detecting when an AI’s expressed confidence misaligns with its actual accuracy, leading to inappropriate over-reliance on incorrect outputs (Yubo et al. 2024). Educational contexts amplify this vulnerability: students exhibit overconfidence in their own ability to detect AI-generated errors while employing inadequate verification strategies (Martín-Moncunill, Bañeres, and Serra-Sagristà 2025). The anthropomorphic design of conversational interfaces increases students’ tendency to believe presented information, treating the AI as an authoritative teacher rather than a fallible tool (Natarajan and Gombolay 2020). This combination—AI fallibility intersecting with human detection deficits—creates conditions where misinformation is not just presented but internalized without critical evaluation.

**Instability and Non-Determinism.** The probabilistic nature of language models creates a fundamental challenge for educational deployment: outputs are non-deterministic, and capabilities can change unpredictably with model updates. Research on using AI as an evaluator for educational feedback found significant inconsistency both within individual models and across different models from various developers (Shin et al. 2025). Minor alterations in input phrasing can produce dramatically different responses, and model updates can alter explanations of foundational concepts, grading criteria, or problem-solving approaches without docu-

mentation or institutional oversight (Shin et al. 2025). This volatility makes it impossible for institutions to formally validate, certify, or reliably build standardized learning pathways on these systems.

**Black Box Opacity.** The complex architecture of modern language models makes their internal decision-making processes largely opaque (Delikoura, Fung, and Hui 2025). This opacity is often compounded by the proprietary nature of commercial systems developed “behind closed doors with little transparency” (UNESCO 2023). For educational institutions, this means no verifiable way exists to determine what “curriculum knowledge” the AI has learned, what pedagogical principles it follows, or what logic it applies when evaluating student work. This lack of auditability forces reliance on easily quantifiable performance metrics rather than nuanced process-oriented aspects of learning, and it represents a de facto ceding of pedagogical control from institutions to technology vendors.

**Algorithmic Bias and Equity Concerns.** AI systems trained on internet-scale corpora inevitably learn and reproduce societal biases embedded in that data (Delikoura, Fung, and Hui 2025). In educational contexts, this manifests as biased examples, reinforcement of stereotypes, and discriminatory assessment outcomes. Research has documented that AI-powered detection tools disproportionately falsely flag essays written by non-native English speakers as AI-generated, creating significant risks of false academic dishonesty accusations for this population (Liang et al. 2023). Beyond content bias, algorithmic systems risk reinforcing structural inequities: poorly designed tracking or personalization algorithms can trap students from underserved communities in lower-achievement pathways (Farheen et al. 2025), while the persistent digital divide limits access for those lacking reliable internet connectivity and devices (UNESCO 2023).

## 2.2 Evidence from Real-World Deployments

Khan Academy’s development of Khanmigo, an AI-powered tutoring system, provides instructive insight into the gap between base model capability and deployment readiness. Public disclosures from the development team reveal that generative AI proved fundamentally unsuited for mathematical reasoning in its raw form: designed for language prediction, the model would generate probable next numbers rather than execute correct calculations (DiCerbo 2025). The solution required building a separate, dedicated calculator tool to handle numerical operations outside the language model. Additional challenges emerged around visual content interpretation (requiring pre-processing pipelines to generate textual descriptions of graphs and figures) and solution path flexibility (requiring re-engineered internal reasoning to map multiple potential approaches) (DiCerbo 2025). This case demonstrates conclusively that deployment-ready educational AI requires extensive domain-specific architecture beyond the base model—the safety and reliability of Khanmigo reside not in GPT-4 itself but in the scaffolding constructed around it.

These documented failures are not isolated incidents but

interconnected systemic challenges. A factual hallucination becomes a pedagogical crisis when students cannot detect the error. A pedagogically sound system becomes unreliable when its knowledge base is unstable. An accurate system becomes ungovernable when institutions cannot audit its embedded curriculum. Addressing these challenges requires not piecemeal solutions but an integrated framework that recognizes their interdependence.

## 3 Related Work and The Integration Gap

The challenges documented in Section 2 have not gone unnoticed. A variety of stakeholders have proposed frameworks, standards, and evaluation approaches for responsible AI deployment in education. However, these efforts remain fragmented across different domains of expertise and institutional functions, each addressing essential but partial dimensions of trustworthiness.

### 3.1 Institutional and Policy Frameworks

Educational institutions and policy bodies have developed high-level frameworks to guide AI adoption from governance and administrative perspectives. The 1EdTech AI Preparedness Checklist provides structured guidance across organizational readiness (forming advisory groups), policy updates (rules on third-party tools, data protection), pedagogical adaptation (assessment redesign), and literacy development (training on proper attribution) (1EdTech n.d.-a). Similarly, the Athena Infonomics Equitable AI in Education Checklist offers guidance for policymakers and technical teams on initiating and assessing ethical AI systems, with focus on regulatory compliance, data handling, and ethical decision-making (Athena Infonomics n.d.).

These frameworks excel at catalyzing institutional conversations and ensuring procedural due diligence. They prompt administrators to consider legal implications, update privacy policies, engage diverse stakeholders, and develop communication strategies. However, their process-oriented nature leaves a critical verification gap. The 1EdTech checklist, for example, advises institutions to ask vendors about data privacy and bias controls but provides no standardized methodology or technical metrics to independently validate vendor responses. These frameworks establish what questions to ask but not how to verify the answers.

### 3.2 Pedagogical and Ethical Frameworks

From the education and ethics communities have emerged frameworks focused on aligning AI use with learning science principles and ethical norms. The UNESCO AI Competency Frameworks for Students and Teachers emphasize human-centered approaches, critical thinking, and ethical considerations in developing AI-specific pedagogical skills (UNESCO 2023). The Comprehensive AI Assessment Framework (CAIAF) provides a tiered model for integrating AI into assessments, with levels ranging from “No AI” to “Full AI Integration,” explicitly grounded in principles of transparency, equity, and accountability (Rahimi and Shute 2024). The Human-Centric AI-First (HCAIF) framework proposes leveraging AI for personalization while mandating

student practices of attribution (documenting AI use) and reflection (analyzing AI effectiveness) (Wilson 2025).

These frameworks provide valuable theoretical and ethical foundations. They offer educators principled approaches to designing learning experiences that use AI responsibly. Their limitation, however, stems from an implicit assumption of AI reliability and transparency that, as Section 2 demonstrates, often does not hold. A teacher can design an assessment according to CAIAF principles, but the framework offers no mechanism to verify that the chosen AI tool will perform its function accurately, consistently, and without bias. They address the pedagogical “how” and “why” without providing technical verification of the “what.”

### 3.3 Technical and Model-Level Evaluation

The computer science and AI safety communities have developed domain-agnostic tools for measuring model performance and behavior. These include standard evaluation metrics such as relevance scores (output alignment with prompts), hallucination indices (frequency of factually incorrect statements), and toxicity scores (identification of harmful content) (Ji et al. 2023). More sophisticated frameworks have been proposed to create holistic views of model trustworthiness. For instance, the Unified Explainability Score (UES) conceptually combines accuracy, interpretability, fidelity, consistency, and stability dimensions into a weighted composite metric (Kandpal and Verma 2025).

The strength of these technical approaches lies in their quantitative rigor and measurability. Their limitation is domain-agnosticism. A low hallucination score is necessary but insufficient for a good educational tool—it does not indicate whether an explanation is pedagogically effective, developmentally appropriate, or aligned with curriculum standards. High stability scores do not guarantee factual correctness or bias absence. A significant gap exists between abstract technical metrics and the context-rich requirements of real-world learning environments.

### 3.4 Vendor-Specific Guidelines and Disclosures

AI developers and consortia provide transparency through efforts like OpenAI’s System Cards, which disclose capabilities, limitations, and safety results (OpenAI 2024). Google developed its LearnLM models, integrated into Gemini 2.5 Pro, designed around learning science principles for pedagogical soundness (Google 2024). Similarly, 1EdTech’s TrustEd Apps Generative AI Data Rubric is a self-assessment framework for vendors to standardize data privacy disclosures (1EdTech n.d.-b).

While valuable, these initiatives are fundamentally limited as forms of self-regulation. System Cards are not independent audits and lack details for institutional verification. A stated commitment to learning science, as with LearnLM, is a design philosophy, not a verifiable guarantee that institutions can independently test against their curriculum. Self-assessment rubrics depend on vendor reporting and often have a narrow scope (e.g., data privacy) rather than the comprehensive trustworthiness dimensions—pedagogical efficacy, accuracy, stability, or auditability—required for deployment.

Framework	Example	Verify	Stable	Audit	Pedagogy
Policy	1EdTech	No	No	Partial	No
Pedagogical	UNESCO	No	No	No	Yes
Technical	Hallucination	Partial	Partial	No	No
Vendor	LearnLM	No	No	No	Claims
<b>TEAS</b>	This work	Yes	Yes	Yes	Yes

Table 1: Framework Gap Analysis. Existing frameworks address pillars in isolation, lacking integrated verifiability. TEAS integrates all four for deployment readiness.

Moreover, the reliability of such self-regulation can be contingent on shifting commercial priorities. During periods of intense competition and rapid product development, such as the large language model development acceleration beginning in 2023, ethical commitments can face pressure. For example, Microsoft reportedly disbanded its core AI ethics and society team in 2023 as part of a broader restructuring aimed at accelerating product timelines, raising concerns about the long-term stability of vendor-led ethical governance when faced with market demands (Vincent 2023). This underscores the need for standards that rely on independent verification rather than solely on vendor commitments, which can be subject to change under commercial pressure.

### 3.5 The Integration Gap

Table 1 illustrates the fragmentation across these framework categories. Each addresses valuable dimensions of trustworthiness but none integrates the full requirements for deployment readiness.

This fragmentation reflects a separation of expertise and institutional function. Administrators use policy checklists that lack technical depth. Educators use pedagogical frameworks that assume reliable underlying technology. Engineers use technical metrics that omit educational context. Vendors provide disclosures that cannot be independently verified. These stakeholder groups do not work from shared standards or common language, creating silos that prevent holistic risk assessment.

The critical missing element across this landscape is integrated verifiability. Non-technical frameworks—those designed for policymakers and educators—rest on trust in vendor claims. Institutions are advised to inquire about bias mitigation but given no standard methods to test those strategies independently. Educators are encouraged to use AI as Socratic partners but have no way to guarantee consistent behavior. For AI to be safely deployable in a high-stakes field like education, standardization akin to mature sectors (e.g., medicine (U.S. Food and Drug Administration n.d.), aviation) is necessary, where reliance on self-attestation is insufficient and independent certification is mandatory. But the current ecosystem for educational AI lacks this essential layer, leaving institutions to navigate high-risk environments based on trust rather than evidence that can be independently validated.

No existing framework comprehensively addresses what it means for an educational AI system to be deployment-ready: simultaneously verifiable in its factual claims, stable

in its core knowledge, auditable by institutions, and pedagogically sound in its interactions. This integration gap motivates the framework proposed in the following section.

## 4 The TEAS Framework

This section presents **TEAS (Trusted Educational AI Standard)**, an integrated framework that addresses the gaps identified in previous sections. TEAS defines deployment readiness through four interdependent pillars that together establish comprehensive trustworthiness requirements for educational AI systems.

### 4.1 Framework Overview

TEAS is designed not as a quantitative benchmark or performance score but as a deployment readiness standard—a set of verifiable requirements that a system must satisfy before it can be considered trustworthy for educational use. This approach draws inspiration from mature regulatory frameworks in other high-stakes domains: medical device approval processes that evaluate safety and efficacy before clinical deployment (U.S. Food and Drug Administration n.d.), aviation certification standards that verify component reliability before flight operations, and financial system auditing requirements that mandate independent verification of institutional claims.

The framework rests on four pillars: Verifiability, Stability, Auditability, and Pedagogical Soundness. These are not independent criteria to be evaluated in isolation but interdependent requirements that must be satisfied holistically. A system that excels in three pillars while failing the fourth remains untrustworthy for deployment, as the pillars mutually reinforce one another to establish comprehensive trustworthiness.

### 4.2 Pillar 1: Verifiability

**Definition.** Verifiability requires that AI-generated educational content be traceable to authoritative, validated sources. When an AI system explains a concept, solves a problem, or answers a factual question, the system must be able to ground its response in specific curriculum documents, peer-reviewed textbooks, or institutionally approved knowledge bases.

**Requirements.** A verifiable system must provide citations or references to specific sources and, within them, specific sections for factual claims. It must distinguish clearly between statements drawn from authoritative sources and generated inferences or explanations that extend beyond direct source material. Critically, the system must not fabricate or hallucinate sources—a failure mode where models cite non-existent papers, textbook sections, or authorities to lend false credibility to generated content.

**Why it matters.** Verifiability directly addresses the factual accuracy and hallucination problems documented in Section 2. It enables both educators and students to fact-check AI outputs against trusted sources, transforming the AI from an opaque oracle into a transparent reasoning system whose claims can be validated. For institutions, verifiability provides a mechanism to ensure AI-generated content aligns with approved curriculum standards and does not

introduce unauthorized material into the learning environment.

**Current gap.** While retrieval-augmented generation (RAG) architectures exist and can ground responses in document collections, no standard defines what constitutes “sufficient” source grounding for educational contexts. How specific must citations be? What sources are authoritative? How should a system behave when asked questions beyond its grounded knowledge base? These questions remain unanswered in current practice, allowing vendors to claim source-grounding without clear verification criteria.

### 4.3 Pillar 2: Stability

**Definition.** Stability requires that core curriculum knowledge be deterministic and consistent over time. For foundational concepts, formulas, definitions, and procedures that form the stable backbone of a discipline, the AI system must provide consistent explanations across different query phrasings, user sessions, and system versions.

**Requirements.** A stable system must produce functionally identical responses to semantically equivalent questions about foundational knowledge. If a student asks “What is Newton’s Second Law?” today and again tomorrow, or phrases it as “Explain F=ma,” the core explanation must remain consistent. When model updates or system maintenance occur, fundamental knowledge must not change— $E=mc^2$  remains  $E=mc^2$ , the Krebs cycle retains its established steps, and the quadratic formula maintains its form. For knowledge that is genuinely evolving or contested in the discipline (cutting-edge research findings, theoretical debates), the system should explicitly indicate the provisional or debated nature of the information rather than presenting it with false certainty.

**Why it matters.** Stability addresses the non-determinism problem that prevents standardization. Educational systems depend on predictable knowledge foundations: curriculum developers need to know what students will be taught, assessment designers need consistent knowledge representations, and students building mental models with the help of AI need reliable information that doesn’t contradict itself over time. Stability also enables fair assessment—students cannot be fairly evaluated on knowledge that the AI teaches inconsistently.

**Current gap.** The probabilistic architecture of language models produces inherent non-determinism. Temperature settings, sampling methods, and subtle prompt variations can yield different responses to identical questions. Model updates—often deployed by vendors without detailed changelogs—can alter explanations of fundamental concepts without institutional visibility or approval. No current standard requires or even proposes mechanisms to ensure deterministic core knowledge in AI tutoring systems, leaving stability unaddressed in deployment decisions.

### 4.4 Pillar 3: Auditability

**Definition.** Auditability requires that educational institutions be able to independently inspect, validate, and certify an AI system’s knowledge base, reasoning logic, and embedded pedagogical assumptions. An auditable system makes

visible what it knows, how it reasons, and what values or biases might influence its outputs.

**Requirements.** An auditable system must provide institutions with mechanisms to examine its knowledge base—including curriculum content, sources, and potential errors. It must also offer inspectable reasoning for its pedagogical decisions and support bias testing to ensure alignment with institutional standards. Crucially, these inspections must be possible through independent third-party processes, not solely through vendor-controlled interfaces or self-reported metrics.

**Why it matters.** Auditability addresses the black box problem and institutional governance challenge. Without it, institutions cannot verify vendor claims (e.g., accuracy, bias mitigation) or ensure the AI's curriculum aligns with institutional and accreditation standards. Auditability restores pedagogical authority to educational institutions rather than ceding it to technology vendors. It also enables accountability, allowing institutions to diagnose failures stemming from knowledge gaps, reasoning errors, or misaligned pedagogical assumptions.

**Current gap.** The proprietary nature of commercial language models prevents meaningful institutional auditing. Training data remains undisclosed, model architectures are trade secrets, and internal decision processes are computationally intractable to trace even when architectures are known. No certification process currently exists for educational AI analogous to curriculum approval processes for textbooks. Institutions are left to trust vendor assertions without independent verification capacity, a stance that would be unacceptable in other high-stakes domains like medicine or finance (U.S. Food and Drug Administration n.d.).

#### 4.5 Pillar 4: Pedagogical Soundness

**Definition.** Pedagogical soundness requires that AI systems adhere to evidence-based teaching principles and interaction patterns that foster active learning, critical thinking, and genuine understanding rather than passive information consumption or cognitive offloading.

**Requirements.** A pedagogically sound system must employ Socratic questioning and guided discovery rather than direct answer provision—when a student asks “What is the answer to this problem?”, the system should respond with guiding questions that help the student construct the solution themselves. It must provide appropriate scaffolding matched to the learner’s demonstrated level, neither overwhelming novices with advanced concepts nor boring advanced learners with excessive simplification. The system must actively prevent “spoiling” of solutions: it should not complete assignments for students, write essays on their behalf, or otherwise circumvent the learning process even when directly instructed to do so. Pedagogically sound systems should incorporate proven instructional strategies such as spaced repetition, retrieval practice, and metacognitive prompting that encourage students to reflect on their own learning processes.

**Why it matters.** Pedagogical soundness addresses the misalignment problem documented in Section 2—the phenomenon of “vaporized learning” where AI assistance pro-

duces short-term performance gains while undermining long-term retention and understanding. Even a perfectly accurate, stable, and auditable system can harm learning if its interaction patterns encourage cognitive offloading rather than cognitive engagement. Pedagogical soundness ensures that AI serves as a proper learning tool—one that supports and amplifies student thinking rather than replacing it, much as calculators freed learners to focus on mathematical reasoning rather than arithmetic drudgery, without becoming a dependency that prevents understanding.

Critically, pedagogical soundness is not merely an educational nicety but a prerequisite for institutional deployability. During the peak of generative AI adoption in 2023-2024, educational institutions faced a crisis: teachers could no longer confidently assign homework or take-home essays, uncertain whether submitted work represented student learning or AI completion. Multiple school districts and universities considered or implemented outright bans on tools like ChatGPT precisely because these systems undermined fundamental pedagogical principles (Elsen-Rooney 2023). The barrier to deployment was not technical capability but pedagogical trust. If AI systems can be built with sufficient guardrails—systems designed from the ground up to improve learning rate and cognition rather than act as “cheat codes”—institutions can be relieved of their paranoia about academic integrity erosion. This transforms AI from a threat to classroom pedagogy into a deployable tool at institutional scale. Without pedagogical soundness, AI remains something educators must defend against; with it, AI becomes something they can confidently integrate into curricula.

**Current gap.** While pedagogical principles are well-established (UNESCO 2023) and articulated by vendors (Google 2024), no verification mechanism exists for black-box AI systems. Educators cannot certify consistent Socratic behavior. Security vulnerabilities like jailbreaking allow students to bypass guardrails (OWASP Foundation 2025). The gap lies in verifying and enforcing reliable pedagogical behavior.

#### 4.6 The Interdependence of Pillars

The four pillars of TEAS are not independent checkboxes but interdependent requirements. Consider the failure modes of partial compliance:

- A system that is verifiable but pedagogically unsound might accurately cite sources while directly providing solutions, undermining learning.
- A system that is pedagogically sound but unstable might use good Socratic questioning but give contradictory explanations, creating confusion.
- A system that is verifiable and stable but not auditable might teach reliably from vetted sources yet embed subtle biases institutions cannot detect.
- A system that is pedagogically sound, stable, and auditable but not verifiable might teach consistently based on hallucinated sources, creating systematic misinformation.

True deployment readiness requires satisfying all four pillars simultaneously. They are mutually reinforcing: verifiability

bility enables auditability; stability enables pedagogical consistency; auditability ensures standards alignment; pedagogical soundness ensures effective knowledge delivery.

#### 4.7 Implications for Deployment Economics

A critical insight emerging from TEAS concerns the relationship between trustworthiness and model capability. The dominant assumption suggests expensive frontier models are inherently more trustworthy. TEAS challenges this, demonstrating trustworthiness stems primarily from architectural design.

The Khanmigo case exemplifies this (DiCerbo 2025). This distinction is crucial for educational equity. If systematic frameworks addressing TEAS principles, rather than expensive models, determine trustworthiness, then affordable open-source models can achieve deployment-grade standards. As Prof. Balaraman Ravindran and colleagues articulate, deployment at scale in India (~150M students) requires costs around Rs. 30 (~\$0.34) per student/year (Ravindran 2025). Such economics are impossible with frontier models, but TEAS suggests a path: augmenting sovereign/open-source models (e.g., from Sarvam AI, AI4Bharat) with TEAS-compliant architectures enables responsible deployment at billion-student scale across the Global South, where commercial pricing is prohibitive. This reframes R&D towards systematic trustworthiness frameworks, making AI benefits accessible equitably.

### 5 Discussion and Implications

#### 5.1 Implications for Researchers

TEAS reorients research priorities from capability maximization toward trustworthiness systematization. While benchmarks like MMLU are important, they do not capture the comprehensive requirements for deployment, creating new research agendas for each TEAS pillar.

Critically, TEAS suggests that advancing trustworthiness in smaller models can have greater deployment impact than capability gains in frontier models, reallocating priorities toward educational equity.

#### 5.2 Implications for Educational Institutions

TEAS provides institutions a structured rubric for procurement. Instead of relying on vendor claims, they can systematically assess systems against the four pillars, empowering institutional review boards and procurement officers.

This framework shifts assessment from trust toward evidence-based validation. It also creates accountability, providing a diagnostic tool to identify pillar failures and guide corrective action when systems fail post-deployment.

#### 5.3 Implications for AI Developers and EdTech Companies

TEAS provides design requirements for building trustworthiness from the outset, which is more achievable than retrofitting. As the Khanmigo case shows, architecture is key (DiCerbo 2025). Developers can use TEAS to guide the design of hybrid systems, transparent audit features, and architecturally enforced pedagogical soundness.

TEAS compliance becomes a competitive advantage. In a market hesitant due to trust concerns, systems demonstrating TEAS compliance through independent verification will differentiate themselves, aligning market incentives with educational values.

#### 5.4 Implications for Policy and Regulation

TEAS offers a foundation for certification and regulatory frameworks, analogous to processes in medicine (U.S. Food and Drug Administration n.d.) or aviation. Policymakers could mandate TEAS compliance for AI deployed institutionally, potentially via independent auditors (like textbook review committees). This necessitates domain-specific standards, as generic AI safety rules don't capture unique educational risks (pedagogical harm, curriculum stability). TEAS helps define these specific needs for comprehensive AI regulation in education.

#### 5.5 Implications for Educational Equity and Global Deployment

Perhaps TEAS's most significant implication is enabling democratized access. AI acts as a cognitive lever; tying trustworthiness only to expensive models amplifies global disparities. TEAS challenges this by showing trustworthiness stems from architecture, allowing affordable, open-source/sovereign models to meet deployment standards. As researchers in India target costs of ~Rs. 30 (~\$0.34) per student/year for mass deployment (150M+ students) (Ravindran 2025), TEAS provides the assurance framework. Augmenting models from Sarvam AI, AI4Bharat, etc., with TEAS-compliant architectures makes reliable AI education viable at billion-student scale across the Global South. This aligns with digital sovereignty and makes AI an equity enabler, a public good rather than a luxury.

### 6 Limitations and Future Work

While TEAS provides a comprehensive framework, limitations point toward future work. TEAS is prescriptive, defining *what* trustworthiness requires, not yet *how* to operationally measure compliance; developing metrics and audit tools is crucial. The framework focuses on content/pedagogy, explicitly excluding distinct agentic AI security risks (Legatt 2025), needing a separate standard (e.g., "TEAS-Security"). Implementation requires cross-disciplinary collaboration (educators, AI researchers, policymakers). The framework needs empirical validation across diverse contexts. Finally, addressing the dynamic nature of AI requires frameworks for continuous monitoring and re-certification. Despite these, TEAS offers necessary common ground for building verifiable trust.

### 7 Conclusion

Current AI evaluation in education is fragmented, hindering trustworthy deployment. This paper introduced TEAS, a unified standard integrating Verifiability, Stability, Auditability, and Pedagogical Soundness. TEAS highlights that systematic architecture, not just model capability, enables trustworthiness, making reliable AI affordable and scalable

for global equity. The research community must prioritize building and validating such trustworthy frameworks to responsibly unlock AI's educational potential for all learners. Billions of students await AI tools they can trust; TEAS provides the standard to make that trust verifiable.

### Conflict of Interest Statement

The author is the Founder and CEO of Metacog, which developed the Equation Grounder module evaluated in Appendix A. To support reproducibility and advance research in verifiable educational AI, Metacog is releasing the Equation Grounder as open-source software concurrent with publication. The experimental protocol employed blind evaluation with an independent third-party judge (Claude Opus 4.5) to ensure objective assessment. The primary contribution of this work is the TEAS framework—a vendor-agnostic evaluation standard applicable to any educational AI system. The case study demonstrates the framework's application using an open implementation.

### References

1EdTech. n.d.-a. AI Preparedness Checklist. <https://www.1edtech.org/resource/ai-checklist>. Accessed: 2025-10-24.

1EdTech. n.d.-b. TrustEd Apps Generative AI Data Rubric. <https://www.1edtech.org/standards/ai-rubric>. Accessed: 2025-10-24.

Athena Infonomics. n.d. Equitable AI in Education – Checklist for AI Deployment. <https://www.athenainfonomics.com/resources/equitable-ai-in-education-checklist-for-ai-deployment>. Accessed: 2025-10-24.

Chilton, J. 2025. AI will reshape education. Are we building tools we can trust? TEDxSNHU. YouTube. <https://youtu.be/oV6HWzzeD-I>.

Delikoura, I.; Fung, Y. R.; and Hui, P. 2025. From Superficial Outputs to Superficial Learning: Risks of Large Language Models in Education. [arXiv:2509.21972](https://arxiv.org/abs/2509.21972).

DiCerbo, K. 2025. Khanmigo math computation and tutoring updates. Khan Academy Blog. <https://blog.khanacademy.org/khanmigo-math-computation-and-tutoring-updates/>.

Elsen-Rooney, M. 2023. NYC education department blocks ChatGPT on school devices, networks. Chalkbeat. <https://ny.chalkbeat.org/2023/1/3/23537987/nyc-schools-ban-chatgpt-writing-artificial-intelligence/>.

Farheen, S.; Cheema, A.; Ullah, R.; and Bandeali, D. 2025. Equity and Bias in AI Educational Tools: A Critical Examination of Algorithmic Decision-Making in Classrooms. *The Critical Review of Social Sciences Studies*, 3: 67–85.

Google. 2024. LearnLM: Improving Gemini for Learning. Google Research. [https://services.google.com/fh/files/misc/improving-gemini-for-education\\_v7.pdf](https://services.google.com/fh/files/misc/improving-gemini-for-education_v7.pdf).

Ji, Z.; Lee, N.; Frieske, R.; Yu, T.; Su, D.; Xu, Y.; Ishii, E.; Bang, Y.; Madotto, A.; and Fung, P. 2023. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12): 1–38.

Kandpal, K.; and Verma, D. 2025. Unified Explain Ability Score (UES): A Comprehensive Framework for Evaluating Trustworthy AI Models. *International Research Journal on Advanced Engineering and Management (IRJAEM)*, 3: 185–193.

Lee, D. 2024. Duolingo cuts human translators as CEO says AI can do the job 'faster and cheaper'. Financial Times. <https://www.ft.com/content/80f7273a-4ddc-4a37-9759-9a09c2a304e2>.

Legatt, A. 2025. Colleges And Schools Must Block And Ban Agentic AI Browsers Now. Here's Why. Forbes. <https://www.forbes.com/sites/avivalegatt/2025/09/25/colleges-and-schools-must-block-agentic-ai-browsers-now-heres-why/>.

Liang, W.; Yuksekgonul, M.; Mao, Y.; Javadi, E.; and Zou, J. 2023. GPT detectors are biased against non-native English writers. [arXiv:2304.02819](https://arxiv.org/abs/2304.02819).

Martín-Moncunill, D.; Bañeres, D.; and Serra-Sagristà, J. 2025. Students' Trust in AI and Their Verification Strategies: A Case Study Comparing STEM and Humanities Undergraduates. *Education Sciences*, 15(10): 1307.

Natarajan, M.; and Gombolay, M. 2020. Effects of Anthropomorphism and Accountability on Trust in Human Robot Interaction. In *2020 15th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, 33–42. IEEE.

OpenAI. 2024. GPT-4o System Card. [arXiv:2410.21276](https://arxiv.org/abs/2410.21276).

OWASP Foundation. 2025. LLM01: Prompt Injection. OWASP Top 10 for Large Language Model Applications Project. Accessed: 2025-10-24. <https://genai.owasp.org/llmrisk/llm01-prompt-injection/>.

Rahimi, S.; and Shute, V. J. 2024. Comprehensive AI Assessment Framework: Enhancing Educational Evaluation with Ethical AI Integration. [arXiv:2407.16887](https://arxiv.org/abs/2407.16887).

Ravindran, B. 2025. Central Square Foundation — The evolving space for AI in Education. YouTube. <https://youtu.be/inFk4MpFfnA>.

Shin, H.; Lee, S.-K.; Kim, D.-Y.; Kim, M.-J.; and Kim, H.-C. 2025. Large Language Models as Evaluators in Education: Verification of Consistency and Reliability. *Applied Sciences*, 15(2): 671.

UNESCO. 2023. Guidance for generative AI in education and research. <https://unesdoc.unesco.org/ark:/48223/pf0000386693>.

U.S. Food and Drug Administration. n.d. Digital Health Center of Excellence. FDA. Accessed: 2025-10-24. <https://www.fda.gov/medical-devices/digital-health-center-excellence>.

Vincent, J. 2023. Microsoft lays off team that taught employees how to make AI tools responsibly. The Verge. <https://www.theverge.com/2023/3/13/23638823/microsoft-ethics-society-team-responsible-ai-layoffs>.

Wilson, M. 2025. A Framework for Human-Centric AI-First Teaching. AACSB Insights. <https://www.aacsb.edu/insights/articles/2025/02/a-framework-for-human-centric-ai-first-teaching>.

Yubo, Z.; Zhaoning, L.; Jacob, S.; and Lixing, W. 2024. Understanding the Effects of Miscalibrated AI Confidence on User Trust, Reliance, and Decision Efficacy. arXiv:2402.07632.

Zou, R.; Li, Z.; Chan, C. K.; Zhang, Z.; Yao, Y.; and Zhang, H. T. 2025. Hide and Seek with LLMs: An Adversarial Game for Sneaky Error Generation and Self-Improving Diagnosis. arXiv:2508.03396.

## A Empirical Validation of the TEAS Framework

This appendix presents empirical evidence for TEAS’s central thesis through a controlled comparison study. We demonstrate that an 8-billion parameter open-source model augmented with structured knowledge graph grounding outperformed models up to 15 $\times$  larger across criteria operationalizing the TEAS pillars.

### A.1 Introduction

The central thesis of TEAS is that trustworthiness in educational AI stems primarily from systematic architecture rather than raw model capability. This claim has profound implications: if true, it suggests that affordable, open-source models augmented with appropriate architectural constraints can achieve the deployment-grade trust currently associated only with expensive frontier systems. This would fundamentally alter the economics of AI deployment in education, making reliable tutoring systems viable at billion-student scale.

We conducted a controlled comparison study in which an 8-billion parameter open-source model (Qwen3-8B), augmented with structured knowledge graph grounding, was evaluated against three larger baseline models—including a 120-billion parameter system—across five GATE-style mathematics questions. The evaluation was performed blindly by an independent large language model judge (Claude Opus 4.5) using criteria directly derived from the TEAS pillars: factual accuracy, hallucination prevention, citation quality (Verifiability), format compliance (Auditability), and Socratic method adherence (Pedagogical Soundness).

**Key findings.** The 8B model with structured grounding achieved an average score of 13.8/15, outperforming the 120B baseline model by 19% (11.6/15) and an 80B model by 11% (12.4/15). This performance gap was driven by superior citation quality—the grounded model achieved a perfect 3.0/3.0 score with traceable node-level citations—and consistent format compliance, both direct manifestations of architectural constraints rather than emergent model capability. Perhaps most importantly, the evaluation methodology itself was designed to eliminate evaluator bias: responses were anonymized, shuffled, and assessed without knowledge of which system produced which output.

This is not a claim that small models are inherently superior, nor that larger models are unnecessary. Rather, it demonstrates that trustworthiness is an architectural property that can be engineered systematically. A small model constrained to ground every claim in a structured knowledge graph cannot hallucinate content not present in its source material; a large model with unrestricted generation, no matter how capable, retains that failure mode. The architectural constraint—not the parameter count—determines the outcome.

At its core, this is about shifting AI systems from probabilistic generation to deterministic retrieval. Language models predict what token “should” come next based on learned patterns; grounded systems retrieve what the source material actually says. This distinction is the foundation of the TEAS

framework: education requires factual certainty, but LLMs provide only statistical likelihood. Architectural grounding bridges this gap.

## A.2 Experimental Design

**Research Questions** This study investigates three primary questions:

**RQ1:** Can structured knowledge graph grounding enable a small language model (8B parameters) to outperform larger ungrounded models (80B, 120B) on educational tutoring tasks when evaluated against TEAS criteria?

**RQ2:** Which specific TEAS pillars benefit most from architectural grounding versus raw model scale?

**RQ3:** Does blind evaluation (where the judge is unaware of model identities) validate the trustworthiness claims, or might results reflect evaluator bias toward particular response styles?

**Knowledge Domain and Source Material** We constructed a synthetic but realistic educational scenario mirroring graduate-level preparation for technical examinations. The source material was a LaTeX document titled “*Optimization Methods in Machine Learning: From Gradient Descent to Adaptive Moment Estimation*” (approximately 300 lines), covering mathematical preliminaries ( $\theta \in \mathbb{R}^d$ , Lipschitz smoothness  $L$ , strong convexity  $\mu$ ), gradient descent methods (SGD, mini-batch variants), adaptive methods (Adam optimizer, momentum  $v_t$ , variance  $s_t$ ), and convergence analysis (condition number  $\kappa = L/\mu$ ,  $O(1/t)$  rates).

This domain was chosen deliberately. Mathematical optimization is sufficiently technical to stress-test factual grounding and reasoning, yet narrow enough that source material completeness could be verified.

**Knowledge Graph Construction** Metacog is a framework for building trustworthy educational AI systems, developed with a focus on STEM education where factual precision is non-negotiable. At its core is the open-source Equation Grounder module, which transforms unstructured mathematical documents into queryable semantic graphs. The fundamental insight: by making AI systems retrieve from structured knowledge rather than generate from learned patterns, we shift from probabilistic prediction to deterministic retrieval—the difference between “what token comes next” and “what does the source material actually say.”

The Equation Grounder serves as proof-of-concept for the TEAS Verifiability pillar. It demonstrates that grounding is not a theoretical aspiration but an implementable architecture using standard NLP tools. Critically, deterministic retrieval from structured knowledge eliminates the primary failure mode of large language models: hallucination through stochastic generation.

**Extraction Process.** For the grounded condition, we used the Equation Grounder to extract structured knowledge from the LaTeX source. The extraction process identified:

- **Nodes (n=31):** Variable nodes (mathematical symbols:  $\theta$ ,  $\eta$ ,  $g_t$ ,  $L$ ,  $\mu$ ), equation nodes (inline and block-level with unique identifiers), and metadata (LaTeX source, position, section context)

Model	Params	Context	Input	Output
Qwen3-8B+KG	8B	128K	\$0.028	\$0.110
R1 Distill 14B	14B	32K	\$0.120	\$0.120
Qwen3-80B	80B	131K	\$0.120	\$1.200
GPT-OSS 120B	120B	131K	\$0.040	\$0.200

Table 2: Model specifications and costs per 1M tokens. The grounded 8B model achieves lowest cost (\$0.138/M total) with superior trustworthiness.

- **Relationships (n=14):** APPEARS\_IN edges connecting variables to equations with contextual snippets

Each node received a UUID-based identifier (e.g., 4:5380d574-...:5), enabling precise citation at retrieval time. This structure transforms unstructured LaTeX into a queryable semantic graph where every mathematical symbol can be traced to its defining equations and usage contexts.

Critically, the knowledge graph was not manually curated. It was automatically generated using the equation grounder, ensuring reproducibility. The entire extraction process is deterministic—run it twice on the same document, get identical graphs. This automation is essential for scalability: manual knowledge engineering would create a deployment bottleneck incompatible with billion-student targets.

**Experimental Conditions Condition A (Baseline - No Grounding).** Models received the full LaTeX document as unstructured text in the system prompt. They were instructed to (1) answer using Socratic method, (2) reference the provided study material, and (3) output responses in JSON format with specific fields.

Three models were tested: R1 Distill Qwen 14B (14B parameters), Qwen3-80B (80B parameters), and GPT-OSS 120B (120B parameters).

**Condition B (With Grounding).** One model received the same LaTeX document plus the constructed knowledge graph in CSV format with extracted variables, equations, and relationships. The model was architecturally constrained via prompt engineering to (1) only use information from the extracted knowledge, (2) cite specific sections/equations, and (3) output included an additional grounding\_citations field with node IDs.

Model tested: Qwen3-8B (8B parameters).

The grounded 8B model is not only the most trustworthy but also the most cost-effective, with per-token costs 42-88% lower than baseline models (Table 2). This cost-performance decoupling is critical: trustworthiness typically correlates with expense, but architectural constraints reverse this relationship.

**Evaluation Questions** Five questions were designed to span core concepts, mirroring GATE-style difficulty:

1. “In the context of gradient descent, what does the symbol  $\theta$  represent, and what is its domain?”
2. “Write the update rule for Stochastic Gradient Descent. What is the role of the learning rate  $\eta$ ?”
3. “What is the condition number  $\kappa$ , and how does it affect convergence rate?”

Model	Ground	Q1	Q2	Q3	Q4	Q5	Total	Avg
Qwen3-8B+KG	Full KG	14	14	15	13	13	69/75	13.8
Qwen3-80B	Doc only	13	13	12	13	11	62/75	12.4
GPT-OSS 120B	Doc only	12	12	11	11	12	58/75	11.6
R1 Distill 14B	None	4	4	3	5	6	22/75	4.4

Table 3: Overall performance (max 15 points per question). The grounded 8B model outperformed all ungrounded models, including those 10-15× larger.

- “Explain the difference between the first moment  $v_t$  and second moment  $s_t$  in Adam. Why is bias correction needed?”
- “If a function is  $L$ -smooth but not strongly convex ( $\mu = 0$ ), what is the convergence rate of gradient descent?”

All models received identical questions in identical order. Sampling parameters were held constant (temperature=1.0, top\_p=1.0, reasoning enabled).

**Evaluation Rubric** The rubric operationalizes TEAS pillars into five criteria, each scored 0-3:

- Factual Accuracy** (Verifiability): 0=major errors, 3=fully correct per source
- Hallucination Score** (Verifiability, inverted): 0=significant invented content, 3=zero hallucination
- Socratic Method** (Pedagogical Soundness): 0=direct answer dump, 3=true Socratic guidance
- Citation Quality** (Verifiability+Auditability): 0=no citations, 3=node-level traceability
- Format Compliance** (Stability+Auditability): 0=wrong format, 3=perfect structured JSON

Total: 15 points maximum per question. This rubric directly maps to TEAS requirements.

**Blind Evaluation Protocol** To eliminate evaluator bias, we implemented rigorous blind evaluation:

- Response Collection:** All 4 models answered all 5 questions (20 responses total)
- Anonymization:** Per question, responses were randomly assigned IDs ( $\alpha, \beta, \gamma, \delta$ ), mapping concealed
- Blind Evaluation:** Responses presented to Claude Opus 4.5 with explicit instruction: “You do NOT know which AI model produced each response”
- De-anonymization:** After scoring, mapping revealed model identities
- Cross-Validation:** Fresh randomization per question prevented pattern recognition

This methodology is critical. Without blinding, the evaluator might inadvertently favor responses with characteristics associated with “advanced” systems independent of actual TEAS compliance.

Model	Factual	Halluc	Socratic	Citations	Format
Qwen3-8B+KG	<b>2.8</b>	<b>2.8</b>	2.2	<b>3.0</b>	<b>3.0</b>
Qwen3-80B	2.8	2.0	<b>2.8</b>	2.0	2.8
GPT-OSS 120B	<b>3.0</b>	<b>2.8</b>	1.0	2.0	2.6
R1 Distill 14B	2.6	1.6	0.0	0.0	0.0

Table 4: Average scores by criterion (max 3.0 per criterion). Grounded model achieved perfect citation quality and format compliance.

### A.3 Results

**Aggregate Performance** Table 3 shows overall performance across all five questions.

The grounded 8B model outperformed all ungrounded models, including those with 10-15× more parameters. Performance advantage was consistent (low std dev = 0.84), indicating that architectural grounding provides stable rather than sporadic improvements.

**Key Finding 1:** An 8B model with structured knowledge graph grounding outperformed a 120B model without grounding by 19% (13.8 vs 11.6), and an 80B model by 11% (13.8 vs 12.4). Simultaneously, the grounded model’s per-token cost (\$0.138/M) was 42% lower than GPT-OSS 120B (\$0.240/M) and 88% lower than Qwen3-80B (\$1.320/M). This demonstrates cost-performance decoupling: superior trustworthiness at lower deployment cost.

**Performance by TEAS Criterion** Table 4 shows average scores by evaluation criterion.

**Key Finding 2:** The grounded model achieved perfect citation quality (3.0/3.0) across all questions, enabled by node-level citations. No ungrounded model exceeded 2.0/3.0, as they could only provide section-level references.

**Key Finding 3:** The grounded model achieved perfect format compliance (3.0/3.0), a direct result of architectural constraints requiring structured JSON output. The 120B model averaged 2.6/3.0 (occasionally violating format), while the 14B model scored 0.0/3.0 (never complied).

**Key Finding 4:** The 120B model exhibited “pseudo-Socratic” behavior (score 1.0/3.0), consistently asking questions then immediately answering them. This pattern reflects RLHF training for “helpfulness”—the model cannot resist explaining even when instructed to withhold.

**Critical Failure Modes** Four distinct failure modes emerged in ungrounded models, none present in the grounded condition:

**Failure Mode 1: Textbook Knowledge Substitution.** The 14B model used a different definition of  $\kappa$  (ratio of eigenvalues of Hessian) rather than the document’s definition ( $\kappa = L/\mu$ ). The model retrieved plausible alternative content from training data instead of grounding in provided source. The grounded model, constrained by the knowledge graph, could not make this error.

**Failure Mode 2: Pseudo-Socratic Helpfulness.** GPT-OSS 120B pattern: “Can you think about X? Great! Here’s the answer...” RLHF training for helpfulness overrides pedagogical instruction (average Socratic score 1.0/3.0).

**Failure Mode 3: Example Invention.** Qwen3-80B hallucinated “logistic regression with 5 features” not present in source—ungrounded elaboration to enhance explanation.

**Failure Mode 4: Format Non-Compliance.** R1 Distill produced Markdown instead of required JSON, making institutional auditing impossible.

Critically, none of these failure modes appeared in the grounded condition. The architectural constraints—not emergent model behavior—prevented each failure category.

#### A.4 Analysis Through TEAS Pillars

**Verifiability: Every Claim Must Be Traceable** **Experimental Evidence:** Grounded model achieved perfect citation score (3.0/3.0); ungrounded models maximum 2.0/3.0.

**Architectural Mechanism:** The knowledge graph provided node-level citation infrastructure. When the grounded model referenced  $\theta$ , it cited specific node IDs (e.g., 4:5380d574-...:5) and equation numbers. These citations are independently verifiable. An institutional auditor can trace node IDs to exact graph nodes, retrieve LaTeX source, and confirm accuracy. Ungrounded models provided only section-level references, insufficient for claim-by-claim verification.

**Hallucination Prevention:** The structured format acted as a hard constraint. The model could only cite nodes present in the provided CSV. This transformed hallucination from an emergent failure mode to an architecturally impossible outcome.

Compare to ungrounded 80B model: “For example, in a linear regression model where the prediction is  $\hat{y} = \theta^T x...$ ” This example does not appear in source. The model invented it for pedagogical clarity, violating Verifiability. The grounded model could not make this error because output was constrained to reference only nodes in the knowledge graph.

**Implication:** Verifiability is not a property that models learn through scale or fine-tuning. It is an architectural property that must be engineered.

**Stability: Deterministic Core Knowledge** **Experimental Evidence:** While this study did not test multi-session stability directly, format compliance serves as proxy. The grounded model achieved 3.0/3.0 format compliance, producing perfectly structured JSON in every response at temperature=1.0.

**The Determinism Imperative:** This finding reveals a fundamental tension in LLM-based education systems. Language models are inherently probabilistic—they sample from learned distributions over token sequences. Educational systems, by contrast, require determinism—the same question should yield the same core content across sessions to enable standardized curricula and fair assessment.

Architectural grounding resolves this tension. When a model retrieves from a fixed knowledge graph rather than generating from learned patterns, the source of non-determinism (stochastic sampling over learned weights) is replaced with deterministic lookup over structured data. The model still generates natural language explanations, but the factual content is anchored to invariant sources.

**Implication:** Stability cannot be achieved through sampling parameters alone. Even temperature=0 does not guarantee deterministic outputs when the model generates content from learned patterns rather than retrieving from fixed sources. Architectural grounding provides the structural consistency that curriculum standardization requires. **This is the essence of the TEAS thesis: making AI systems more deterministic by making them less generative.**

**Auditability: Institutional Validation Capacity** **Experimental Evidence:** Grounded model outputs included explicit thinking-process field showing reasoning steps, sources\_used array with node IDs, and grounding\_citations array with equation references. Ungrounded models lacked structured reasoning traces.

**Architectural Mechanism:** The JSON schema requirement creates machine-readable audit trails. An institution can parse outputs and automatically: (1) extract all cited node IDs, (2) cross-reference against the knowledge graph, (3) verify that cited nodes support claims, and (4) flag discrepancies for human review.

This workflow is impossible with ungrounded models because: (1) citations are not structured, (2) no node IDs exist to cross-reference, and (3) source material is unstructured text, not a queryable graph.

**Implication:** Auditability is the difference between “trust us” and “verify yourself.” Institutions adopting TEAS-compliant systems can independently validate correctness without relying on vendor assurances.

**Pedagogical Soundness: Evidence-Based Teaching** **Experimental Evidence:** Grounded model averaged 2.2/3.0 Socratic score; 120B model 1.0/3.0; 80B model 2.8/3.0; 14B model 0.0/3.0.

**Partial Success and Limitations:** The grounded model’s Socratic score (2.2/3.0) was not perfect. It received “2” ratings for “asks questions with partial withholding”—meaning it sometimes asked guiding questions but then provided substantial explanatory content before waiting for student response.

This reveals an architectural gap: prompt-level instruction alone is insufficient to enforce pedagogical method. Unlike Verifiability (where citation requirements could be structurally enforced), Pedagogical Soundness requires more sophisticated architectural mechanisms—potentially multi-turn conversation management that physically prevents answer revelation until the student responds.

**Critical Insight from 120B Model:** The pseudo-Socratic pattern is instructive: “Before we give a formal answer, can you think about what gradient descent is trying to adjust? Great! In gradient descent the symbol  $\theta$  denotes...” The model asks a question, then immediately answers it on behalf of the student. This emerged from RLHF training optimizing for “helpfulness”—a value misalignment: what RLHF considers “helpful” contradicts what pedagogy considers “sound.”

The lesson: Pedagogical Soundness cannot be fine-tuned into models. It must be architecturally enforced through conversation management systems that prevent premature answer disclosure.

**LearnLM: Proof That Pedagogical Soundness is Achievable—At a Price.** Google’s LearnLM (Google 2024) represents a pedagogically-specialized architecture. LearnLM incorporates conversation management designed for educational contexts: explicit scaffolding, metacognitive prompting, and adaptive response withholding.

We evaluated Gemini 3 Pro (which integrates LearnLM capabilities) on the same five questions. The results validate the LearnLM approach: the model achieved an average Socratic score of 2.5/3.0—substantially better than GPT-OSS 120B (1.0/3.0) and even outperforming the grounded 8B model (2.2/3.0). This demonstrates that pedagogical soundness is architecturally achievable when prioritized in model design.

However, this capability comes at prohibitive cost. Gemini 3 Pro’s inference pricing: (1)  $\leq 200K$  tokens: \$2.00/M input, \$12.00/M output (\$14.00/M total); (2)  $> 200K$  tokens: \$4.00/M input, \$18.00/M output (\$22.00/M total). This represents  $101\text{-}159\times$  higher cost than the grounded 8B model (\$0.138/M). At target economics of  $\sim \$0.34/\text{student/year}$ , Gemini 3 Pro would limit deployment to  $\sim 24\text{-}35K$  tokens/student/year—insufficient for even a single, 30-minute chat session, let alone semester-long courses.

**The Critical Insight:** Gemini 3 Pro proves that frontier models *can* exhibit pedagogically sound behavior through specialized training. This validates that Pedagogical Soundness is a learnable property, not an intractable challenge. The implication for TEAS is profound: smaller open-source models can be finetuned on pedagogical objectives to match frontier model capabilities at fraction of the cost. Rather than requiring external conversation management scaffolding, pedagogy-aware finetuning of affordable models may be sufficient.

The research pathway forward: (1) **Distillation:** Use Gemini 3 Pro outputs as training data for smaller models; (2) **Pedagogical RLHF:** Train reward models that prioritize learning outcomes over user satisfaction; (3) **Specialized finetuning:** Augment base 8B models with education-specific instruction datasets.

If an 8B model with pedagogical finetuning can match Gemini 3 Pro’s Socratic score (2.5/3.0) while maintaining cost advantage, the result would be transformative: TEAS-compliant systems that excel on all four pillars at equity-compatible prices. Gemini 3 Pro demonstrates the target is achievable; the challenge is replicating it affordably.

**Implication:** Gemini 3 Pro with LearnLM proves that Pedagogical Soundness is achievable through model-level design, not just external scaffolding. The model’s 2.5/3.0 Socratic score—obtained through specialized training on educational objectives—demonstrates that pedagogy can be learned. However, at  $101\text{-}159\times$  the cost of grounded small models, frontier model deployment remains economically infeasible at billion-student scale.

The research opportunity is clear: pedagogically finetune affordable open-source models to replicate Gemini 3 Pro’s capabilities. If successful, this approach would yield TEAS-compliant systems excelling on all four pillars at equity-compatible costs. The grounded 8B model in this study achieved perfect Verifiability and Auditability; pedagogi-

cal finetuning could close the remaining gap in Pedagogical Soundness without requiring expensive proprietary models or complex external scaffolding.

## A.5 Discussion and Implications

**Central Finding: Architecture > Scale** The primary finding is that systematic architecture determines trustworthiness more than raw model capability. An 8B model with structured grounding outperformed models 10-15 $\times$  larger across TEAS criteria. This was not due to the 8B model being “smarter”—it was due to architectural constraints preventing failure modes that larger models retain.

This finding has three immediate implications:

**1. Trustworthiness is Engineerable.** Institutions and developers need not wait for the next generation of frontier models. Trustworthiness can be built today using systematic frameworks applied to existing, affordable models.

**2. Cost-Performance Decoupling.** The economic barrier to reliable educational AI is lower than commonly assumed. The grounded 8B model achieved superior trustworthiness at 42-88% lower cost than baseline models. If a  $\$0.34/\text{student/year}$  deployment target is feasible with grounded small models at \$0.138/M tokens, then billion-student-scale deployment becomes realistic. This is not a marginal improvement—it is the difference between feasibility and impossibility for resource-constrained contexts.

**3. Validation Methodology Matters.** Blind evaluation was critical to this study’s credibility. Without anonymization, evaluators might have scored responses based on perceived model sophistication rather than TEAS compliance. The blind protocol ensures that findings reflect architectural properties, not evaluator biases.

**Limitations of Prompt-Based Grounding** This study implemented grounding through prompt engineering: the knowledge graph was provided as CSV text in the system message. This approach has inherent limitations:

**Prompt Injection Vulnerability:** A malicious user could craft inputs that override grounding constraints. Prompt-based grounding is not cryptographically secure.

**Context Window Constraints:** The knowledge graph used here (31 nodes, 14 relationships) fit comfortably in the context window. Larger curricula would require retrieval-augmented generation (RAG) architectures, introducing retrieval accuracy as an additional failure mode.

**No Physical Enforcement:** The model was *instructed* to ground responses but not *forced* to. A production system would require architectural enforcement—e.g., a separate validation layer that rejects any response lacking valid node citations.

Despite these limitations, the results demonstrate the principle: even imperfect architectural grounding dramatically improves trustworthiness. Production systems would strengthen these constraints further.

**Implications for Model Selection** This study compared models of vastly different scales (8B vs 120B), but the takeaway is not “small models are better.” Rather:

**Small models with a TEAS-compliant framework > Large models without a TEAS-compliant framework**

Scenario	Recommended Approach
Narrow, well-defined curriculum	Small model + structured KG
Broad, open-ended tutoring	Larger model + RAG + validation
Multilingual deployment	Sovereign models + KG
High-stakes assessment	Deterministic KG retrieval only

Table 5: Model selection framework. Match architectural complexity to use case requirements.

Table 5 suggests a decision framework for institutions. The key insight: match architectural complexity to use case requirements. Over-capability introduces failure modes and cost barriers; under-capability risks inadequacy. TEAS provides criteria for making this decision systematically. For most educational contexts—structured curricula with defined learning objectives—the smallest model with appropriate grounding is optimal on both trustworthiness and cost dimensions.

**Equity and Global Deployment** The economic implications are profound. If trustworthy AI tutoring can be built with 8B models running on modest hardware, then:

**Deployment Cost Projections:**

- **Inference Cost:** 8B models at \$0.138/M tokens vs 80B at \$1.320/M ( $9.6 \times$  difference)
- **On-Device Deployment:** 8B models can run locally on devices, eliminating API costs
- **Knowledge Engineering:** One-time cost to build KG from curriculum materials (automatable via tools like the Equation Grounder)

**Target Economics:** Researchers in India aim for ~\$0.34/student/year for 150M+ students (Ravindran 2025). This study demonstrates that such costs are achievable without sacrificing trustworthiness. Contrast this with commercial AI tutoring platforms charging \$20/student/month (\$60-120/year)—a  $175\text{-}350 \times$  price difference that renders AI inaccessible to most of the world’s students.

**Sovereignty and Localization:** Resource-constrained nations can build TEAS-compliant systems using: (1) open-source base models (Qwen, Llama, Sarvam AI), (2) locally-constructed knowledge graphs from national curricula, and (3) on-premises deployment (no data sent to foreign servers). This aligns with digital sovereignty goals while ensuring trustworthiness through systematic architecture rather than dependence on external providers.

**Future Research Directions 1. Scaling KG Construction:** Our study used a ~300-line LaTeX document. Future work should investigate fully automated KG extraction from textbook PDFs, multi-document KG integration, and curriculum-level KG standards for interoperability.

**2. Long-Form Reasoning:** The questions required 1-2 paragraph responses. Investigating grounding for multi-step problem-solving would test architectural constraints under greater reasoning demands.

**3. Multi-Session Stability:** TEAS requires consistent outputs across sessions. Future work should test whether grounded models produce identical citations when asked the same question multiple times.

**4. Pedagogical Finetuning for Affordable Models:**

Gemini 3 Pro with LearnLM achieved 2.5/3.0 Socratic score, proving pedagogical soundness is learnable. Critical research directions include: (a) distillation experiments training 8B models on Gemini 3 Pro’s pedagogical outputs, (b) pedagogical RLHF designing reward models prioritizing learning outcomes over satisfaction metrics, (c) open pedagogical datasets curating instruction-tuning datasets for Socratic dialogue patterns, and (d) comparative evaluation measuring whether pedagogically-finetuned 8B models match frontier model quality at fraction of cost (\$0.138/M vs \$14-22/M).

**5. Cross-Domain Validation:** This study focused on mathematical optimization. Validating the architecture>scale thesis across domains (history, literature, programming) would strengthen generalizability claims.

## A.6 Conclusion

This controlled study provides empirical evidence for TEAS’s central thesis: trustworthiness in educational AI stems primarily from systematic architecture rather than raw model capability. An 8-billion parameter model augmented with structured knowledge graph grounding outperformed models up to  $15 \times$  larger across criteria operationalizing the TEAS pillars—achieving perfect citation quality (Verifiability), perfect format compliance (Auditability), and superior hallucination prevention.

The findings validate each TEAS pillar’s achievability through systematic architecture:

- **Verifiability:** Node-level citations enabled claim-by-claim traceability impossible with ungrounded models (perfect 3.0/3.0 score)
- **Stability:** Structured output schemas produced consistent formatting even at high temperature (perfect 3.0/3.0 score)
- **Auditability:** Machine-readable reasoning traces allowed independent institutional validation (perfect 3.0/3.0 score)
- **Pedagogical Soundness:** Grounded 8B model achieved 2.2/3.0; Gemini 3 Pro with LearnLM proved 2.5/3.0 is achievable, establishing a clear target for pedagogical finetuning of affordable models

The evaluation methodology—blind scoring by an independent judge—ensures that these findings reflect objective architectural properties rather than evaluator bias. The judge had no knowledge of model identities, eliminating the risk that scores favored responses exhibiting characteristics associated with “advanced” systems.

The implications for educational equity are immediate. If systematic frameworks—not expensive frontier models—determine trustworthiness, then reliable AI tutoring becomes economically viable at billion-student scale. Institutions and developers in resource-constrained contexts

need not wait for access to commercial systems. They can build TEAS-compliant systems today using affordable, open-source models augmented with structured knowledge architectures.

The path forward requires:

- 1. Research community:** Developing automated KG construction tools for scaling grounding architectures; pedagogical finetuning of open-source models to match Gemini 3 Pro's 2.5/3.0 Socratic capability; cross-domain validation studies beyond mathematics
- 2. Institutions:** Adopting evaluation frameworks that assess architecture (not just capability) when making procurement decisions, explicitly requiring TEAS pillar compliance
- 3. Developers:** Building hybrid systems combining small models with structured grounding AND pedagogical finetuning—leveraging both architectural constraints and learned behaviors
- 4. Policymakers:** Recognizing that trustworthy AI education is achievable at scale today, enabling regulatory frameworks that mandate TEAS compliance rather than prohibit deployment

This study demonstrates that the barriers to trustworthy educational AI are not technological—they are architectural. The fundamental challenge is not building smarter models, but building more deterministic systems that retrieve from structured knowledge rather than generate from probabilistic patterns. Large language models will continue to improve, but their core architecture remains probabilistic. Educational deployment requires determinism.

The use of tools like the Equation Grounder, evaluated in this study, proves that deterministic grounding is implementable today using standard tools. The lesson extends beyond mathematics: any domain with authoritative source material can be structured into queryable graphs. The question is no longer whether trustworthy AI education is possible—it demonstrably is. The question is whether institutions, developers, and policymakers will prioritize architectural trustworthiness over superficial capability when making deployment decisions that affect billions of students.

## B Knowledge Graph Node Export

This appendix provides the complete node export from the knowledge graph constructed by Metacog's Equation Grounder. The graph contains 31 nodes representing mathematical variables and equations extracted from the source LaTeX document.

### B.1 Node Structure

Each node contains the following fields:

- `~id`: Unique UUID identifier
- `~labels`: Node type (Variable or Equation)
- `name`: Symbol name (for variables) or empty
- `type`: Content type (inline, block, or empty)
- `latex`: LaTeX source code of the content

### B.2 Complete Node Export (CSV Format)

#### Variable Nodes (27 total):

```
~id,~labels,name,type,latex
4:::::0,Variable,x_i.,
4:::::1,Variable,y_i.,
4:::::2,Variable,N.,
4:::::3,Variable,\ell.,
4:::::4,Variable,\mathcal{L}.,
4:::::5,Variable,\theta.,
4:::::6,Variable,\mathbb{R}.,
4:::::7,Variable,d.,
4:::::8,Variable,L.,
4:::::9,Variable,\mu.,
4:::::10,Variable,\nabla.,
4:::::11,Variable,t.,
4:::::12,Variable,\alpha.,
4:::::13,Variable,\eta.,
4:::::14,Variable,B.,
4:::::15,Variable,M.,
4:::::16,Variable,i.,
4:::::17,Variable,g_t.,
4:::::18,Variable,\beta_1.,
4:::::19,Variable,\beta_2.,
4:::::20,Variable,v_t.,
4:::::21,Variable,s_t.,
4:::::22,Variable,\epsilon.,
4:::::23,Variable,\hat{v}_t.,
4:::::24,Variable,\hat{s}_t.,
4:::::25,Variable,\kappa.,
4:::::26,Variable,\theta^*,
```

#### Equation Nodes (10 total):

```
~id,~labels,name,type,latex
4:::::27,Equation,,inline,f
4:::::28,Equation,,block,
  \mathcal{L}(\theta) = 
  \frac{1}{N} \sum_{i=1}^N 
  \ell(x_i, y_i; \theta)
4:::::32,Equation,,block,
  \theta \in \mathbb{R}^d
4:::::33,Equation,,block,
  \mathcal{L}(\theta_A) \leq 
  \mathcal{L}(\theta_B) + 
  \nabla \mathcal{L}(\theta_B)^T 
  (\theta_A - \theta_B) + 
  \frac{1}{2} \|\theta_A - \theta_B\|^2
4:::::34,Equation,,block,
  \theta_{t+1} = 
  \theta_t - \eta g_t
4:::::38,Equation,,block,
  v_t = \beta_1 v_{t-1} + 
  (1-\beta_1)g_t; s_t = 
  \beta_2 s_{t-1} + (1-\beta_2)g_t^2
4:::::40,Equation,,block,
  \hat{v}_t = 
  \frac{v_t}{1-\beta_1^t}; 
  \hat{s}_t = 
  \frac{s_t}{1-\beta_2^t}
4:::::41,Equation,,block,
  \theta_{t+1} = \theta_t - \eta 
  \frac{\hat{v}_t}{\sqrt{\hat{s}_t + \epsilon}}
4:::::43,Equation,,inline,
  \mathcal{L}(\theta)
```

```

4:....:52,Equation,,inline,
\kappa = L/\mu
Note: UUIDs abbreviated as "4:...:N"
where N is the position. Full UUID:
4:5380d574-a61c-450b-9d80-55fd2be49de1:N

```

### B.3 Node Distribution Statistics

Node Type	Count	Percentage
Variables	27	87.1%
Equations (block)	7	22.6%
Equations (inline)	3	9.7%
<b>Total</b>	<b>31</b>	<b>100%</b>

Table 6: Distribution of node types in the knowledge graph.

#### Key Observations:

- Variables represent fundamental mathematical symbols ( $\theta, \eta, L, \mu$ , etc.)
- Block equations represent major mathematical statements (loss function, update rules)
- Inline equations represent contextual references within prose
- Each node is independently citable via its UUID, enabling precise traceability

## C Knowledge Graph Relationship Export

This appendix provides the complete relationship export showing how variables connect to equations in the knowledge graph. All relationships are of type APPEARS\_IN, connecting variable nodes to equation nodes where they appear.

### C.1 Relationship Structure

Each relationship contains:

- `~start_node_id`: UUID of the variable node
- `~end_node_id`: UUID of the equation node
- `~relationship_type`: Always APPEARS\_IN
- `context`: Textual snippet describing the relationship

### C.2 Complete Relationship Export (CSV Format)

Due to space constraints in two-column format, we present a representative sample of the 14 relationships. Full export available in the supplementary materials.

```

~start,~end,~type,context
4:....:5,4:....:28,APPEARS_IN,
"The model is parameterized by
a weight vector. We define this
vector as theta in R^d, where d
represents the dimensionality
of the parameter space."
4:....:5,4:....:32,APPEARS_IN,
"The parameter vector theta
belongs to the d-dimensional
real vector space."
4:....:17,4:....:34,APPEARS_IN,

```

"The update rule for Mini-batch SGD is given by  $\theta_{t+1} = \theta_t - \eta g_t$ , where  $g_t$  is the stochastic gradient estimate."

```

4:....:13,4:....:34,APPEARS_IN,
"The learning rate eta controls
the step size in the gradient
descent update."

```

```

4:....:20,4:....:38,APPEARS_IN,
"The first moment v_t maintains
an exponentially decaying
average of past gradients."

```

```

4:....:21,4:....:38,APPEARS_IN,
"The second moment s_t maintains
an exponentially decaying
average of past squared
gradients."

```

```

4:....:23,4:....:40,APPEARS_IN,
"Bias-corrected first moment
estimate hat{v}_t is computed
by dividing v_t by (1-beta_1^t)."

```

```

4:....:24,4:....:40,APPEARS_IN,
"Bias-corrected second moment
estimate hat{s}_t is computed
by dividing s_t by (1-beta_2^t)."

```

```

4:....:5,4:....:41,APPEARS_IN,
"Adam's parameter update uses
bias-corrected moment estimates
to compute theta_{t+1}.""

```

```

4:....:25,4:....:52,APPEARS_IN,
"The condition number kappa is
defined as the ratio of the
Lipschitz constant to the strong
convexity parameter: kappa=L/mu."

```

*Note: UUIDs abbreviated. Full ID format: 4:5380d574-a61c-450b-9d80-55fd2be49de1:N. Remaining 4 relationships follow the same structure.*

### C.3 Relationship Statistics

Metric	Value
Total relationships	14
Unique variable nodes (start)	10
Unique equation nodes (end)	8
Avg relationships per equation	1.75
Most connected variable	$\theta$ (4 relationships)

Table 7: Statistics of the knowledge graph relationships.

#### Key Observations:

- The parameter vector  $\theta$  is the most connected node, appearing in 4 different equations (empirical risk, domain definition, Adam update, objective function)

- Context snippets provide semantic information beyond structural connectivity, enabling the model to understand *how* a variable is used in an equation, not just *that* it appears
- The condition number equation (node 52) demonstrates multi-variable relationships:  $\kappa$ ,  $L$ , and  $\mu$  all connect to the same equation with distinct contextual roles
- These relationships form the basis for citation generation: when the grounded model references  $\theta \in \mathbb{R}^d$ , it can cite both the variable node (5) and equation node (32) with full traceability

## D Evaluator Prompt Template

This appendix provides the complete prompt template used for blind evaluation by Claude Opus 4.5. The template was applied separately to each of the five questions, with responses anonymized before evaluation.

### D.1 Meta-Prompt Structure

The evaluator prompt consists of four components:

1. **System Instructions:** Establishes the evaluator role and blind evaluation requirements
2. **Evaluation Rubric:** Defines the five criteria and scoring scale (0-3 per criterion)
3. **Ground Truth Reference:** Provides the correct answer from the source document for comparison
4. **Anonymized Responses:** Presents the four responses with random identifiers ( $\alpha, \beta, \gamma, \delta$ )

### D.2 Complete Prompt Template

Due to length, we present the key sections of the prompt. Full template available in supplementary materials.

#### SYSTEM:

You are an expert evaluator for educational AI tutoring responses. Your task is to blindly evaluate multiple responses using a strict rubric.

**CRITICAL:** You do NOT know which AI model produced each response. Do not attempt to identify the source.

---

#### EVALUATION CONTEXT:

Responses from AI tutors answering gradient descent questions. Source: LaTeX document on "Foundations of Gradient-Based Optimization" covering:

- Problem Formulation (theta in  $\mathbb{R}^d$ )
- Stochastic Gradient Descent
- Adaptive Methods (Adam)
- Convergence Analysis

**Scenario:** Graduate student preparing for technical exams. AI should use Socratic method and cite sources.

---

#### RUBRIC (Score each 0-3):

1. **FACTUAL\_ACCURACY** (Verifiability)
  - 0 = Major errors vs source
  - 1 = Some errors/imprecision
  - 2 = Mostly correct, minor issues
  - 3 = Fully correct per source
2. **HALLUCINATION** (Verifiability)
  - 0 = Significant invented content
  - 1 = Multiple invented examples
  - 2 = Minor additions (e.g., values)
  - 3 = Zero hallucination
3. **SOCRATIC\_METHOD** (Pedagogical Soundness)
  - 0 = Dumps answer directly
  - 1 = Asks then immediately answers
  - 2 = Partial withholding
  - 3 = True Socratic|guides without revealing
4. **CITATION\_QUALITY** (Verifiability + Auditability)
  - 0 = No citations
  - 1 = Vague references
  - 2 = Specific sections/equations
  - 3 = Sections + equation numbers + node IDs
5. **FORMAT\_COMPLIANCE** (Auditability)
  - 0 = Wrong format (e.g., Markdown not JSON)
  - 1 = Partial JSON, missing fields
  - 2 = JSON with minor issues
  - 3 = Perfect JSON, all fields

---

#### GROUND TRUTH REFERENCE:

[Question-specific ground truth inserted here]

Example for Q1:

Q: "What does theta represent, and what is its domain?"

Ground Truth (Section 2.1):

- theta = parameter/weight vector
- Domain: theta in  $\mathbb{R}^d$  ( $d$  = dimensionality of parameter space)

---

[RESPONSE ALPHA]: [inserted]

[RESPONSE BETA]: [inserted]

[RESPONSE GAMMA]: [inserted]

[RESPONSE DELTA]: [inserted]

---

OUTPUT FORMAT (JSON) :

```
{
  "evaluations": [
    {
      "response_id": "alpha",
      "scores": {
        "factual_accuracy": <0-3>,
        "socrative_method": <0-3>,
        "citations": <0-3>,
        "format_compliance": <0-3>
      }
    }
  ]
}
```

```

        "hallucination": <0-3>,
        "socratic_method": <0-3>,
        "citation_quality": <0-3>,
        "format_compliance": <0-3>
    },
    "total": <sum>,
    "key_observations": "..."
},
...
"ranking": ["best", "2nd", ...],
"reasoning": "..."
}

```

**IMPORTANT:**

- Score independently before ranking
- Provide evidence in observations
- Do NOT guess model identities
- Be rigorous on hallucination

### D.3 Evaluation Protocol Notes

**Anonymization Procedure:** For each question, the four responses were:

1. Collected from all four models
2. Randomly assigned to identifiers  $\alpha, \beta, \gamma, \delta$
3. Presented to the evaluator with no other metadata (no model names, parameter counts, or grounding information)
4. Evaluated completely before de-anonymization

**Cross-Question Isolation:** The evaluation chat was cleared between questions to prevent the evaluator from:

- Recognizing response patterns across questions
- Building priors about which identifier corresponds to which model
- Letting performance on earlier questions bias later evaluations

**Judge Model Selection:** Claude Opus 4.5 was chosen as the evaluator because:

- It is not one of the models being evaluated (avoiding self-evaluation bias)
- It has demonstrated strong performance on evaluation and critique tasks
- It can follow complex multi-step evaluation protocols
- It can produce structured JSON output reliably

**Rubric Validation:** The rubric was designed to operationalize TEAS pillars:

- **Verifiability:** Factual accuracy + hallucination + citation quality
- **Stability:** Implicit in factual accuracy (correct per source)
- **Auditability:** Citation quality + format compliance
- **Pedagogical Soundness:** Socratic method

Each criterion uses a 0-3 scale (rather than binary pass/fail) to capture gradations of compliance. This allows identification of “partial success” scenarios (e.g., a model that asks questions but then immediately answers them receives a score of 1 rather than 0 for Socratic method).

### D.4 Example Evaluation Output

For illustration, here is the actual evaluation output for Question 1 (abbreviated for space):

```

{
    "evaluations": [
        {
            "response_id": "alpha",
            "scores": {
                "factual_accuracy": 3,
                "hallucination": 3,
                "socratic_method": 1,
                "citation_quality": 2,
                "format_compliance": 3
            },
            "total": 12,
            "key_observations": "Factually accurate. No hallucination. Pseudo-Socratic: asks then immediately answers. Cites Section 2.1 but lacks node IDs."
        },
        {
            "response_id": "beta",
            "scores": { ... },
            "total": 4,
            "key_observations": "Imprecise. Hallucinated '100k params' example. Zero Socratic. Wrong format (Markdown not JSON)."
        },
        {
            "response_id": "gamma",
            "scores": { ... },
            "total": 13,
            "key_observations": "Fully accurate. Minor hallucination (5-feature example). Excellent Socratic. Perfect JSON."
        },
        {
            "response_id": "delta",
            "scores": { ... },
            "total": 14,
            "key_observations": "Fully accurate with exact quote. Zero hallucination. Partial Socratic. Excellent citations with node IDs (4:...:5)."
        }
    ],
    "ranking": ["delta", "gamma", "alpha", "beta"],
    "reasoning": "Delta ranks first (14/15) with perfect Verifiability and Auditability, achieving node-level traceability. Gamma second (13/15) with best Socratic (3/3). Alpha third (12/15) with pseudo-Socratic. Beta last (4/15) with format failure."
}

```

After receiving this evaluation, the de-anonymization mapping revealed:

- $\alpha$  = GPT-OSS 120B
- $\beta$  = R1 Distill 14B
- $\gamma$  = Qwen3-80B
- $\delta$  = Qwen3-8B + KG

This mapping was concealed from the evaluator during scoring, ensuring that judgments reflected objective criteria rather than model priors.

## D.5 Reproducibility

The complete evaluation materials (all 20 anonymized responses, 5 ground truth references, and 5 evaluation outputs) are available for independent verification. The prompt template provided here enables replication of the evaluation protocol on alternative question sets or with different evaluator models.

**Data Availability Statement:** All experimental materials referenced in Appendices A-D, including the complete source LaTeX document, full model responses, and anonymization mappings, are available at: <https://github.com/Metacog-AI/teas-case-study>.