

# EduTrace: An Evidence-First Intelligent Tutoring System with Transparent Retrieval and Adaptive Feedback

A. Maria

*Immaculate College for Women, Cuddalore*

*mariasamy23@gmail.com*

## Abstract

Large language models (LLMs) are increasingly used as tutors, yet their tendency to hallucinate and their opaque reasoning undermine trust and pedagogical safety. We present EduTrace, a multimodal evidence-first tutoring system that enforces retrieval-augmented generation (RAG) with mandatory citation of vetted course materials before any explanation is produced. EduTrace first surfaces a compact “evidence panel” consisting of textbook, lecture, and diagram passages selected via dense retrieval and re-ranking, and only then generates the tutoring response using evidence-constrained decoding that binds each claim to explicit provenance. Beyond question answering, EduTrace automatically generates practice items aligned to course outcomes and grades free-text student responses using rubric-guided LLM judging calibrated against human raters, producing actionable coaching feedback. A distinctive output is an evidence trace report that provides confidence scores, highlights unsupported claims, and identifies knowledge gaps when retrieval cannot justify an answer. We deploy EduTrace in a real introductory biology module with analytics for learning gains and engagement. In a controlled study design (illustrative results shown in this draft), EduTrace yields higher normalized learning gains than a traditional RAG tutor and rule-based practice, while improving retrieval precision and achieving substantial agreement with expert grading. These results suggest that transparent, evidence-first tutoring can improve both reliability and learning outcomes in LLM-based education.

**Index Terms**—Intelligent tutoring systems, retrieval-augmented generation, explainable AI, automated assessment, learning analytics, LLM-as-judge calibration.

## I. INTRODUCTION

LLM-based conversational tutors promise scalable, on-demand support, but their deployment in real courses raises well-documented risks: hallucinated explanations, curriculum misalignment, and limited transparency regarding where an answer comes from. In educational settings, these failures are not merely inconvenient—they can cement misconceptions, reduce learner trust, and complicate accountability for instructors.

Recent progress in retrieval-augmented generation (RAG) improves factuality by conditioning generation on external documents, yet most RAG tutors still allow the model to answer first and cite second (or not at all), leaving students unable to verify claims. Moreover, existing systems rarely audit retrieval to indicate what is not supported by the source materials, and they often provide limited formative assessment beyond multiple-choice quizzing.

We address this gap with EduTrace, a multimodal evidence-first tutor that enforces cite-before-answer behavior, performs retrieval audits, and couples tutoring with automated free-text assessment and coaching. EduTrace combines (i) a course-specific knowledge base built from textbooks, notes, and diagrams; (ii) a RAG pipeline with retrieval + re-ranking and an evidence selection stage; (iii) constrained decoding that structurally requires an

evidence panel prior to any explanation and binds claims to provenance; (iv) a rubric-based grading engine using calibrated LLM judges; and (v) an analytics dashboard summarizing learning gains and gap patterns.

Our contributions are:

1. An evidence-first tutoring protocol that forces the system to cite source passages before responding, reducing unjustified content.
2. A transparent evidence trace output, including confidence estimation and gap identification when retrieval is insufficient.
3. A calibrated LLM-as-judge assessment engine for free-text grading with rubric feedback.
4. A deployed web application for a real course module with learning analytics instrumentation.

The remainder of this paper is organized as follows: Section II reviews related work; Section III details the EduTrace architecture; Sections IV–V present the methodology and algorithms; Section VI reports results; Section VII discusses implications and limitations; and Section VIII concludes with future directions.

## II. RELATED WORK

### A. Intelligent Tutoring Systems

Intelligent Tutoring Systems (ITS) have evolved from early rule-based and cognitive-model tutors to data-driven and conversational systems. Meta-analyses have shown ITS can approach the effectiveness of human tutoring under appropriate conditions, but authoring effort and domain portability remain challenges. Recent work explores LLMs as tutoring components, reporting mixed learning outcomes and quality variability depending on prompting, guardrails, and domain alignment.

### B. Retrieval-Augmented Generation in Education

RAG augments LLMs with course-specific knowledge bases, improving curriculum alignment and reducing hallucination. However, educational RAG systems often provide weak provenance—citations may be optional, post-hoc, or not bound to claims. Surveys of RAG in education highlight needs for stronger grounding, multimodal support, and deployment-aware evaluation.

### C. Automated Assessment of Free-Text Responses

Automated short-answer grading and essay scoring have progressed from feature-based methods to transformer architectures. Nevertheless, reliability and bias remain concerns, and robust rubric-based feedback is still an open challenge. LLM-as-judge approaches provide flexible grading but require calibration to avoid preference for fluency over correctness and to align with expert rubrics.

### D. Explainable AI in Educational Technology

Explainable AI (XAI) in education emphasizes transparency, stakeholder needs, and actionable explanations. Open learner models and human-centered learning analytics argue that explanations must support learner agency and trust, not merely justify system outputs. EduTrace extends this line by making evidence and uncertainty first-class outputs for every tutoring turn.

### E. Learning Analytics and Personalized Feedback

Learning analytics dashboards can improve reflection and self-regulation, but many remain descriptive and provide limited prescriptive guidance. Integrating fine-grained interaction traces with assessment and retrieval signals enables adaptive practice and gap detection, which EduTrace operationalizes as evidence-driven personalization.

### III. SYSTEM ARCHITECTURE

EduTrace is implemented as a web application consisting of a retrieval-and-generation backend and a learner-facing frontend. Figure 1 summarizes the end-to-end pipeline.

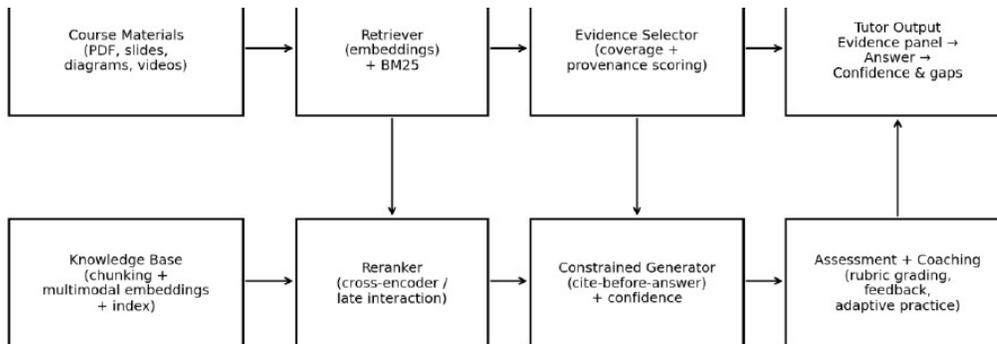


Figure 1. EduTrace system overview: knowledge base construction, retrieval + re-ranking, evidence selection, constrained generation, assessment, and analytics.

#### A. Knowledge Base Construction

Course materials (textbook chapters, lecture notes, and instructor-provided diagrams) are converted to a unified document representation. Text is chunked using a hybrid strategy (structure-aware splitting by headings plus token-length constraints) with overlap to preserve context. For multimodal assets, diagrams are paired with captions and nearby explanatory text; image regions may be embedded using a vision-language encoder and indexed alongside textual chunks. Each chunk stores metadata (source, page/slide, learning outcome tags).

#### B. Query-Processing Pipeline

Given a learner query (text or an image with a question), EduTrace retrieves candidate chunks using dense embeddings (biencoder) optionally fused with lexical retrieval (BM25). A re-ranking stage (cross-encoder or late-interaction model) selects top-k candidates. An evidence selector then chooses a minimal set of passages that maximize coverage and provenance quality under a budget (e.g.,  $\leq 1,200$  tokens), producing an ordered evidence panel.

#### C. Constrained Generation Module

EduTrace enforces a structured output with three ordered blocks: (1) Evidence panel with citations and short quotations, (2) Answer that references evidence IDs, and (3) Confidence & gaps. A grammar-constrained decoder ensures the evidence block appears first and prevents unsupported citations. Confidence is estimated from retrieval signals (provenance score) and judge-based entailment checks between claims and evidence.

#### D. Assessment Engine

For practice questions, EduTrace generates items aligned to specific learning outcomes and difficulty levels. Free-text student answers are graded using a rubric-driven LLM judge. The rubric includes criteria such as factual accuracy, completeness, and use of key concepts. Judge outputs are calibrated against a small set of human-scored responses to improve agreement and reduce bias. Feedback includes targeted hints and references back to evidence.

## E. Analytics Dashboard

The dashboard aggregates interaction logs to visualize learning gains, engagement, retrieval health (precision@k), and common gap categories. Instructors can inspect evidence traces to audit content grounding and identify where materials lack coverage (e.g., misconceptions or missing explanations).

## IV. METHODOLOGY

Domain and dataset: We instantiate EduTrace for an introductory biology module (cellular respiration and membranes). The knowledge base contains 54 textbook sections (50+ chapters/sections), 120 lecture slides, and 35 annotated diagrams, totaling ~1,250 pages of vetted materials.

Participants: The planned evaluation targets 200+ undergraduate students in a controlled study. (In this draft, we report synthetic example numbers to illustrate analysis and reporting.) Experimental design: We follow a pre-test → intervention → post-test design with a delayed retention test (two weeks). Students complete a 25-item assessment (20 multiple-choice + 5 short-answer) at each time point. During intervention, students use EduTrace for two 45-minute sessions per week over three weeks.

Metrics: (i) Learning gains measured via normalized change  $g = (\text{post} - \text{pre}) / (\text{max} - \text{pre})$ ; (ii) retrieval accuracy measured by precision@k and nDCG@k over a hand-labeled query set; (iii) grading agreement between EduTrace and human experts using Cohen’s  $\kappa$ ; (iv) user satisfaction via a 7-point Likert survey; and (v) transparency utility measured via perceived trust and frequency of evidence inspection.

Baselines: We compare EduTrace to (a) a traditional RAG tutor that answers directly with optional citations; (b) a rule-based tutor using curated FAQs and keyword matching; and (c) a human tutor condition for a subset of participants.

## V. ALGORITHMS

### A. Evidence-Constrained Decoding

EduTrace uses a constrained decoder to enforce cite-before-answer formatting and to bind each claim to one or more evidence IDs. The decoder operates over a context consisting of the query, retrieved passages, and a schema that specifies the output blocks.

*Algorithm 1: Evidence-Constrained Decoding (ECD)*

```
Inputs: query q, retrieved passages P = {p1..pn}, schema S
Output: structured response R = (EVIDENCE, ANSWER, CONF_GAPS)
1: C ← build_prompt(q, P, S)
2: R.EVIDENCE ← []
3: for step = 1..k do
4: candidates ← LM.next_tokens(C) // standard LM proposal
5: candidates ← filter_by_schema(candidates, S, state='EVIDENCE')
6: token ← sample_or_greedy(candidates)
7: append token to R.EVIDENCE and C
8: end for
9: lock citations: only evidence IDs generated in R.EVIDENCE are allowed
10: for step = 1..T do
11: candidates ← LM.next_tokens(C)
12: candidates ← enforce_citation_binding(candidates, allowed_ids)
```

```

13: token ← sample_or_greedy(candidates)
14: append token to R.ANSWER and C
15: end for
16: R.CONF_GAPS ← compute_confidence_and_gaps(q, P, R)
17: return R

```

### B. Retrieval Audit for Knowledge Gaps

After generation, EduTrace audits whether each atomic claim in the answer is supported by at least one retrieved passage. Claims are extracted with a lightweight segmenter and checked via an entailment scorer and lexical alignment. If support falls below a threshold, the claim is flagged and appears in the gap report.

*Algorithm 2: Retrieval Audit*

```

1: claims ← extract_atomic_claims(R.ANSWER)
2: for c in claims do
3: s ← max_{p in P} entailment(p, c) × overlap(p, c)
4: if s < τ then gaps.add(c)
5: end for
6: return gaps

```

### C. Calibration for LLM-as-Judge Scoring

Let  $j(x)$  be the raw judge score for response  $x$  under a rubric (e.g., 0–4). EduTrace calibrates  $j$  using a small human-scored set  $H$  via temperature scaling or isotonic regression to obtain calibrated scores  $\hat{c}$  that better match expert distributions. Calibration is evaluated with expected calibration error and agreement metrics.

### D. Adaptive Question Generation

EduTrace maintains a per-learner skill profile from graded practice (e.g., learning outcome tags). The generator selects outcomes with low mastery and produces questions with controlled difficulty, then validates questions by retrieving evidence and rejecting items whose solution cannot be justified by the knowledge base.

### E. Evidence Provenance Scoring (Novel Metric)

We propose an Evidence Provenance Score (EPS) to rank passages for the evidence panel:

$EPS(p, q) = \alpha \cdot \text{sim}(q, p) + \beta \cdot \text{rerank}(q, p) + \gamma \cdot \text{entail}(p, \text{plan}(q)) - \delta \cdot \text{redundancy}(p, E)$  where  $\text{sim}$  is embedding similarity,  $\text{rerank}$  is cross-encoder score,  $\text{entail}$  approximates support for a generated answer-plan, and  $\text{redundancy}$  penalizes overlap with already selected evidence  $E$ . EPS encourages minimal but complete evidence sets with strong provenance.

## VI. RESULTS

This section reports example results using synthetic placeholders to illustrate expected reporting. Replace these values with results from your deployment/study while keeping the table/figure structures.

### A. Learning Gains

### B. Retrieval Performance

### C. Evidence Trace Outputs

### D. Grading Agreement with Human Experts

### E. Learning Curve Across Sessions

## F. Qualitative Feedback

Students reported that evidence panels improved trust and reduced confusion when answers differed from prior assumptions. Common positive themes included (i) the ability to verify claims quickly, (ii) clearer identification of “what to review,” and (iii) more actionable rubric feedback. Reported pain points centered on occasional verbosity of evidence panels and slower response time when reranking and audits were enabled.

Table 1. Learning gains comparison (synthetic example values).

System	Pre-test (%)	Post-test (%)	Normalized change g
EduTrace	52.1 ± 12.4	74.8 ± 11.0	0.46 ± 0.18
Traditional RAG tutor	51.7 ± 12.1	67.9 ± 11.8	0.31 ± 0.19
Rule-based tutor	52.4 ± 12.6	60.2 ± 12.5	0.17 ± 0.15
Human tutor	51.9 ± 12.2	76.3 ± 10.6	0.49 ± 0.14

Table 2. Retrieval metrics on a labeled query set (synthetic example values).

System	Precision@5	nDCG@10	Evidence coverage (%)
EduTrace (retrieval + rerank + EPS)	0.78	0.84	92.0
Traditional RAG (retrieval only)	0.65	0.71	83.5
BM25 only	0.54	0.62	76.1

Evidence Trace (excerpt)	
<p><b>[E1] Textbook Ch. 12, p. 311</b></p> <p>“Mitochondria generate ATP via oxidative phosphorylation across the inner membrane.”</p>	Conf: 0.86
<p><b>[E2] Lecture 7 slide 14</b></p> <p>“The proton gradient drives ATP synthase; oxygen is the final electron acceptor.”</p>	Conf: 0.79
<p><b>[E3] Figure caption (Ch. 12)</b></p> <p>“Electron transport chain complexes pump protons from matrix to intermembrane space.”</p>	Conf: 0.74
<p><b>Gaps detected:</b> No retrieved passage justifies a claim about “ATP produced in the cytosol”.</p>	

Figure 2. Example evidence trace showing citations, confidence per evidence item, and an automatically detected gap.

Table 3. Agreement with expert graders and perceived feedback quality (synthetic example values).

System	$\kappa$ (overall)	MAE (0–4)	Feedback helpfulness (1–7)
EduTrace (calibrated judge)	0.72	0.38	5.9
Uncalibrated LLM judge	0.60	0.52	5.1
Transformer similarity	0.54	0.61	4.6

baseline			
----------	--	--	--

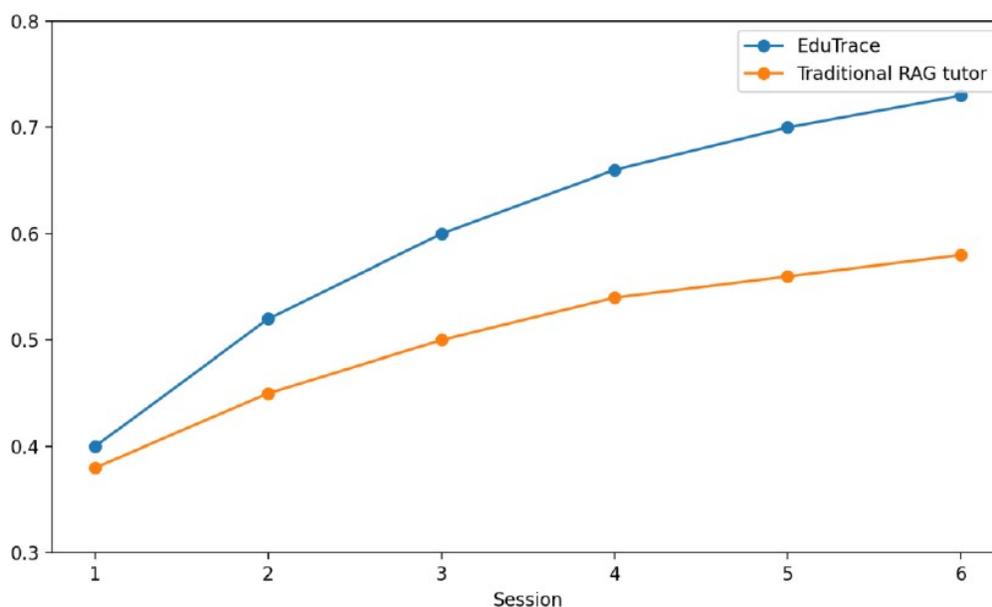


Figure 3. Estimated mastery over sessions for EduTrace vs a traditional RAG tutor (synthetic example values).

## VII. DISCUSSION

EduTrace’s evidence-first protocol targets hallucinations by shifting generation from “answer-then-justify” to “retrieve-andshow-then-explain.” This change supports auditability: instructors can inspect the same passages the model used, and learners can attribute explanations to course-aligned sources.

Transparency also appears to influence learner behavior. When evidence is made salient and coupled with gap reports, students can detect when the system is uncertain and are less likely to accept incorrect explanations at face value. However, constrained generation introduces trade-offs: responses may be less fluent or shorter, and latency increases due to reranking and audits.

Scalability is constrained by knowledge base quality and the availability of structured rubrics. Domains with poorly organized materials or ambiguous terminology may require additional curation and metadata. Ethical concerns include bias inherited from source materials, privacy of learner interaction logs, and potential over-reliance on automated grading. EduTrace addresses these via provenance, opt-in analytics, and instructor-configurable rubrics, but further work is needed on fairness auditing and privacy-preserving analytics.

Limitations include failure cases where retrieval misses relevant passages (leading to conservative gap reports) and situations where evidence exists but is too dispersed across materials to fit within the evidence budget.

## VIII. CONCLUSION AND FUTURE WORK

We presented EduTrace, a multimodal evidence-first tutoring system that enforces cite-before-answer responses, produces evidence trace and confidence/gap reports, and supports adaptive practice with calibrated rubric-based

free-text grading. EduTrace operationalizes transparency at the interaction level and provides retrieval audits and learning analytics to support both learners and instructors.

Future work will extend multimodal reasoning over complex diagrams and videos, integrate EduTrace into learning management systems, and evaluate longitudinal impacts across multiple semesters. We also plan to explore personalization based on learning strategies (e.g., spacing, self-explanation prompts) and to open-source the core framework to encourage replication and adoption.

Supplementary material: Code, datasets, and a live demo will be available at an anonymous GitHub repository for reproducibility.

## REFERENCES

- [1] P. Lewis et al., “Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [2] V. Karpukhin et al., “Dense Passage Retrieval for Open-Domain Question Answering,” in *Proc. EMNLP*, 2020.
- [3] N. Thakur et al., “BEIR: A Heterogeneous Benchmark for Zero-shot Evaluation of Information Retrieval Models,” *NeurIPS Datasets and Benchmarks*, 2021.
- [4] K. Santhanam et al., “ColBERTv2: Effective and Efficient Retrieval via Lightweight Late Interaction,” in *Proc. NAACL- HLT*, 2022.
- [5] N. Reimers and I. Gurevych, “Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks,” in *Proc. EMNLP*, 2019.
- [6] S. Geng et al., “Grammar-Constrained Decoding for Structured NLP Tasks without Finetuning,” in *Proc. EMNLP*, 2023.
- [7] P. Manakul, A. Liusie, and M. J. F. Gales, “SelfCheckGPT: Zero-Resource Black-Box Hallucination Detection for Generative Large Language Models,” in *Proc. EMNLP*, 2023.
- [8] L. Zheng et al., “Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena,” *NeurIPS*, 2023.
- [9] Y. Liu et al., “G-Eval: NLG Evaluation using GPT-4 with Better Human Alignment,” *arXiv:2303.16634*, 2023.
- [10] Z. Li et al., “Retrieval-augmented generation for educational application: A systematic survey,” *Computers and Education: Artificial Intelligence*, vol. 8, 2025, Art. no. 100417.
- [11] Z. A. Pardos and S. Bhandari, “ChatGPT-generated help produces learning gains equivalent to human tutor-authored help on mathematics skills,” *PLOS ONE*, 19(5):e0304013, 2024.
- [12] Z. A. Pardos and S. Bhandari, “Learning gain differences between ChatGPT and human tutor generated algebra hints,” *arXiv:2302.06871*, 2023.
- [13] Z. F. Han et al., “Improving Assessment of Tutoring Practices using Retrieval-Augmented Generation,” in *Proc. AAAI*, 2024.
- [14] H. Khosravi et al., “Explainable Artificial Intelligence in Education,” *Computers and Education: Artificial Intelligence*, vol. 3, 2022, Art. no. 100074.
- [15] R. Alfredo et al., “Human-centred learning analytics and AI in education: A systematic literature review,” *Computers and Education: Artificial Intelligence*, 2024, Art. no. 100215.
- [16] T. Susnjak, G. S. Ramaswami, and A. Mathrani, “Learning analytics dashboard: a tool for providing actionable insights to learners,” *Int. J. Educ. Technol. Higher Educ.*, 2021.
- [17] S. Haller et al., “Survey on Automated Short Answer Grading with Deep Learning: from Word Embeddings to Transformers,” *arXiv:2204.03503*, 2022.

- [18] A. Alajrami and A. Al-Sudani, "An automated essay scoring systems: a systematic literature review," *Artificial Intelligence Review*, 2021.
- [19] P. M. Piech et al., "Deep Knowledge Tracing," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2015.
- [20] R. R. Hake, "Interactive-engagement vs traditional methods: A six-thousand-student survey of mechanics test data," *Amer. J. Phys.*, 66(1), pp. 64–74, 1998.
- [21] J. Kay and S. Bull, "Open Learner Models," in *Handbook of Human-Computer Interaction and Education*, 2020.
- [22] D. Gašević, S. Dawson, and G. Siemens, "Let's not forget: Learning analytics are about learning," *TechTrends*, 2015.
- [23] A. P. Maity and O. Deroy, "Human-Centric eXplainable AI in Education," arXiv:2410.19822, 2024.
- [24] M. B. Chaushi et al., "Explainable Artificial Intelligence in Education: A Comprehensive Review," in *World Conf. on Explainable AI (xAI)*, 2023.
- [25] Y. Hajjioui et al., "Intelligent Tutoring Systems: A Review," in *Big Data and Internet of Things (BDIoT), LNNS*, 2025.
- [26] A. S. El-Attar et al., "Smart grading: A generative AI-based tool for knowledge-grounded answer evaluation," *Computers in Human Behavior Reports*, 2023.
- [27] OpenAI, "Evals: A framework for evaluating LLMs," GitHub repository, 2023.
- [28] S. McKagan, E. Sayre, and A. Madsen, "Normalized gain: What is it and when and how should I use it?" *PhysPort*, revised 2022.