# TeachLM: Post-Training LLMs for Education Using Authentic Learning Data

**Janos Perczel**[*1], **Jin Chow**[†1], and **Dorottya Demszky**[‡2]

[1]Polygence
[2]Stanford University

October 7, 2025

## Abstract

The promise of generative AI to revolutionize education is constrained by the pedagogical limits of large language models (LLMs). A major issue is the lack of access to high-quality training data that reflect the learning of actual students. Prompt engineering has emerged as a stopgap, but the ability of prompts to encode complex pedagogical strategies in rule-based natural language is inherently limited. To address this gap we introduce TeachLM – an LLM optimized for teaching through parameter-efficient fine-tuning of state-of-the-art models. TeachLM is trained on a dataset comprised of 100,000 hours of one-on-one, longitudinal student-tutor interactions maintained by Polygence, which underwent a rigorous anonymization process to protect privacy. We use parameter-efficient fine-tuning to develop an authentic student model that enables the generation of high-fidelity synthetic student–tutor dialogues. Building on this capability, we propose a novel multi-turn evaluation protocol that leverages synthetic dialogue generation to provide fast, scalable, and reproducible assessments of the dialogical capabilities of LLMs. Our evaluations demonstrate that fine-tuning on authentic learning data significantly improves conversational and pedagogical performance – doubling student talk time, improving questioning style, increasing dialogue turns by 50%, and greater personalization of instruction.

## 1 Introduction

In his seminal 1984 study, educational psychologist Benjamin Bloom demonstrated that one-on-one tutoring can yield learning gains two standard deviations above those achieved through traditional classroom instruction [1]. Given the high cost of personalized tutoring, the advent of generative AI has raised hopes of scaling effective one-on-one learning to students worldwide [2–6]. Yet despite the rapid adoption of AI tools such as ChatGPT, Gemini, and Claude by millions of learners [7, 8], these technologies have so far failed to realize that promise [9]. For instance, a recent University of Pennsylvania study found that unfettered access to GPT-4 for math tutoring can harm educational outcomes [10]. Similarly, an MIT study reported that participants who used LLMs exhibited significantly reduced brain connectivity and struggled to quote from answers they wrote just minutes earlier [11].

---

[*]janos@polygence.org (Corresponding author)

[†]jin@polygence.org
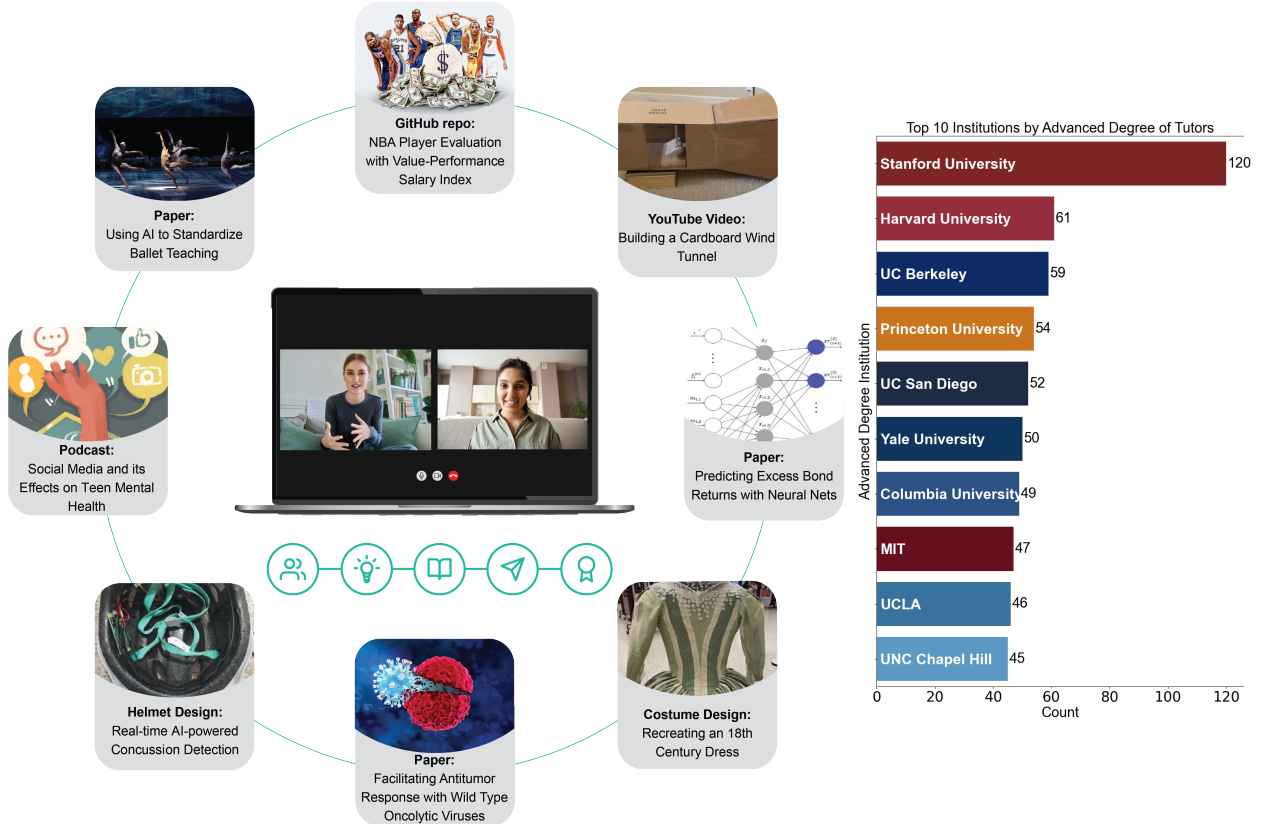
[‡]ddemszky@stanford.edu

Figure 1: Left: Illustration of the Polygence program and project outcomes. Students meet online with tutors pursuing or holding advanced degrees (PhD, MD, JD, MBA, etc.) to take projects from ideation, to exection, to showcasing. Project outcomes range from academic papers to creating podcasts to engineering physical devices. Topics range from AI to cancer biology to sport analytics. Right: Top 10 institutions represented by the advanced degrees pursued or held by Polygence tutors.

A fundamental issue with LLMs is that they have been optimized to act as "helpful assistants" that maximize productivity and minimize cognitive labor [12–20]. This contrasts with the natural friction that expert teachers introduce into learning (for example, by withholding the right answer and prompting students to first attempt a response) [21–23]. Effective tutoring also requires dynamic adaptation to learners' states of mind as opposed to one-size-fits-all instructional designs [24]. This tendency toward friction minimization and sycophantic behavior—prioritizing compliance over pedagogy—is systematically encoded in model parameters through supervised fine-tuning [12] and reinforcement learning from human feedback (RLHF) [25, 26]. These processes rely on datasets produced by human annotators instructed to provide responses that maximize completeness while minimizing the number of conversational turns [15].

Off-the-shelf LLMs can be steered towards improved pedagogy to a limited extent through prompt engineering [27–32], but prompting alone cannot resolve the underlying issues. No finite set of rules or instructions, however sophisticated, can capture the complexity and nuance of high-quality human pedagogy, which necessarily adapts to diverse learners, study contexts, and student goals [33]. We encountered these limitations first-hand in our attempts to build a project-based tutor—called Polypilot—with GPT-4, where iterative prompt refinement led to an endless cycle of increasingly elaborate instructions in response to new scenarios and edge cases. Such peda-

gogical constraints persist even in the most recent education-focused LLMs, including Anthropic's Learning Mode [28], OpenAI's Study Mode [29], and Google's Guided Learning [30] (integrating LearnLM [31]). For example, when confronted with student confusion, these models typically default to rephrasing the problem rather than diagnosing its underlying source. Similar challenges arise in prompt-engineered student simulators, which typically lack the authenticity and diversity needed to represent the full spectrum of learning personas [34].

Post-training frontier models on domain-specific data has recently led to rapid progress toward human-level performance across a range of domains, including coding, law, and science [35–38]. This progress has been made possible by the availability of copious amounts of high-quality training data generated by human annotators who adhere to clearly defined standards of excellence [39]. We expect similar progress to be possible in education given sufficient availability of training data and well-defined success metrics [40]. Recently, Google's LearnLM team demonstrated that supervised fine-tuning of LLMs on synthetic data can improve their performance on a range of education-related benchmarks [33]. Post-training LLMs for education is especially important given that modeling effective pedagogical behavior inherently entails *both* an expert teacher *and* an authentic learner. Without access to a realistic student model—either a high-fidelity simulator or real students, the latter being hard to scale and ethically problematic—benchmarking candidate teacher models is severely limited.

A persistent challenge for post-training educational models is the scarcity of authentic learning data from human teachers and students due to logistical barriers, privacy protections, and concerns about data quality [41–50]. Moreover, human annotators cannot reliably simulate the active learning processes of students or replicate expert teacher practices without engaging with real learners, making the on-demand collection of such data particularly difficult. To address these limitations, researchers from MIT, Carnegie Mellon University, and Cornell University have launched the National Tutoring Observatory [51], with support from the Gates Foundation, the Chan Zuckerberg Initiative [52], and the National Science Foundation [53]. The initiative aims to collect and open-source one million teacher–student interactions to inform the development of AI tutoring tools [51]. While this represents an important step toward alleviating the critical shortage of educational training data, additional work will be required to realize the full potential of post-training models in education.

In this preliminary report, we present a case study on post-training LLMs for education, drawing on data from the Polygence platform [54] consisting of over 100,000 hours of one-on-one, project-based tutoring sessions between PhD-level tutors and students across more than 150 subjects (Fig. 1). Data was curated consistent with the platform's terms of use and privacy policy, reflecting participant opt-outs, and underwent a rigorous anonymization process to protect privacy. Using this dataset, we fine-tune a high-fidelity student model to benchmark frontier LLMs on six multi-turn conversational and pedagogical evaluations. We also fine-tune a teacher model, TeachLM, and demonstrate that it significantly outperforms off-the-shelf models on these benchmarks. Our main contributions are as follows:

1. We develop a pipeline for transcribing, diarizing, and cleaning single-track audio recordings to produce high-quality dialogical data for post-training.

2. We show that student data enables the training of authentic student models, which are essential for scalable and reproducible evaluation of LLMs' pedagogical capabilities.

3. We benchmark off-the-shelf LLMs against human tutors across six education-focused evaluations, highlighting systematic differences in conversational and engagement metrics.

4. We demonstrate that parameter-efficient fine-tuning of state-of-the-art LLMs on authentic learning data substantially improves their pedagogical performance.

We conclude by outlining the limitations of our current approach and identifying next steps for refining the post-training process and evaluating its efficacy.

# 2 Background: Improving LLM Pedagogy

Addressing the pedagogical limitations of LLMs through both post-training and prompt engineering has been an active area of research. Below we review a few of these efforts.

## 2.1 LearnLM: Google's Fine-Tuned Model

A pioneering effort by Google's LearnLM team [55] has focused on improving LLMs for education by post-training. Their early efforts focused on the targeted collection of human tutoring data and the generation of synthetic data under the guidance of education experts [33]. While the on-demand collection of human contractor data from impersonating students proved too noisy for training, their synthetic data allowed the training of an LLM optimized for pedagogy, called LearnLM [33]. LearnLM demonstrated improvements across a range of measurable benchmarks, such as guiding students to answers, promoting engagement, or identifying misconceptions [33]. While synthetic data alone cannot fully capture the nuances of effective human pedagogy and authentic student learning, these improvements highlight the promise of using domain-specific data for post-training LLMs for education. Recently, the LearnLM Team at Google has shifted its focus from *post-hoc fine-tuning* of models to optimizing it for *pedagogical instruction following* of teacher- or developer-defined prompts [56]. While acknowledging the shortcomings of prompting, they explained their decision citing the 'prohibitively difficult' challenge of defining ideal AI tutoring behavior, instead leaving it to teachers and developers to decide on the desired behavior [56]. They also highlighted the cost and overhead of maintaining fine-tuned models while base models are developing rapidly [56].

## 2.2 PolyPilot: Polygence's Prompt-Engineered Tutor

We first experienced the limitations of prompt engineering through a high-conviction internal product experiment in early 2024, when a dedicated engineering and product team at Polygence [54] built Polypilot —- a dynamically prompt-engineered project-based tutor using GPT-4. PolyPilot was deployed in production after many months of user feedback and iterative development, and we collected detailed feedback from more than 70 engaged users (see Appendix A).

Given that Polygence specializes in high-quality project-based tutoring (the success of which can be judged by students' ability to deliver showcaseable artifacts), we quickly recognized the limitations of our prompt-engineered tutor and its ability to steer students towards a project outcome. We invested several months of deliberate effort to iteratively refine stage-dependent prompts and we combined the LLM with sophisticated user interfaces and constantly gathered user feedback to guide our experimentation. However, we soon reached the conclusion that the gap between a human tutor and GPT-4 was simply too large to be closed by prompt engineering. Even relatively simple tasks, such as varying the number and placement of questions or avoiding "wall-of-text" responses to better mirror human tutors, proved inconsistent with prompting. We continued our experimentation with GPT-4o, o1, Gemini 1.5, and Gemini 2.0 Flash, and progressively introduced complex RAG-based approaches to provide LLMs with high-quality examples. However, by early

2025 we concluded that, despite the rapid improvement of LLMs, prompt engineering of off-the-shelf models continues to be insufficient to build a product that can meet the standard of human tutoring.

## 2.3 Anthropic's Learning Mode, OpenAI's Study Mode, and Google's Guided Learning

In recent months, large model developers have released educational LLMs that have been prompt-engineered to improve their pedagogical behavior. Examples include Anthropic's Claude Sonnet 4 with *Learning Mode* [28], OpenAI's GPT 5 with *Study Mode* [29], and Google's Gemini 2.5 Pro with *Guided Learning* [30] (which integrates LearnLM [31]). None of these customized LLMs are available through application programming interfaces (APIs), making rigorous, large-scale, and repeatable evaluation of these models difficult. Nonetheless, to get a directional sense of their behavior, we manually tested all three LLMs via multi-turn conversations dozens of times with different learning scenarios (see Appendix B for more details).

Based on our experimentation, our informal assessment of these educational modes is that their pedagogical capabilities remain rudimentary, and are noticeably constrained by the limitations of rule-based prompt engineering. For example, we found that these models often miss students' learning context, default to multiple-choice-style questioning rather than asking open questions, give away answers, fail to appropriately address confusion, and struggle with verbosity. We also found that models from different providers exhibit remarkably similar behavioral patterns – highlighting both the limited power of prompt engineering to overwrite the pedagogy-agnostic principles learned from large-scale training data and the likely consequence of major AI labs sourcing data from the same vendors (e.g., Scale AI, Surge AI, and Mercor [57]).

# 3 Curating Authentic Learning Data for Post-training

After experiencing the limitations of prompt engineering while building PolyPilot, our focus shifted toward post-training LLMs on authentic learning data from Polygence.

## 3.1 Authentic Learning Dataset

The Polygence dataset comprises more than 100,000 hours of one-on-one, longitudinal student–tutor interactions in multiple modalities. This data was collected through the Polygence platform [54] consistent with the platform's terms and conditions while honoring opt-outs (see Section 3.2) and underwent a rigorous anonymization process to protect privacy. The Polygence platform supports an online program for one-on-one project-based learning under the guidance of expert, PhD-level US-based mentors. Students join Polygence to engage in projects covering a wide array of subjects both in STEM and in the Humanities. Illustrative projects include building a cardboard wind tunnel [60], sewing a historically accurate 19th-century dress [61], writing a paper on using AI to standardize ballet teaching [62], recording a podcast about the neuroscience of dementia [63], and creating a helmet with an AI-powered concussion detection system [64] (see Fig. 1). The projects are overseen by PhD-level experts from top US research institutions. The top 10 institutions represented are shown in Fig. 1.

This dataset represents a distinct opportunity to evaluate the efficacy of training LLMs on authentic learning data because of its unique characteristics:
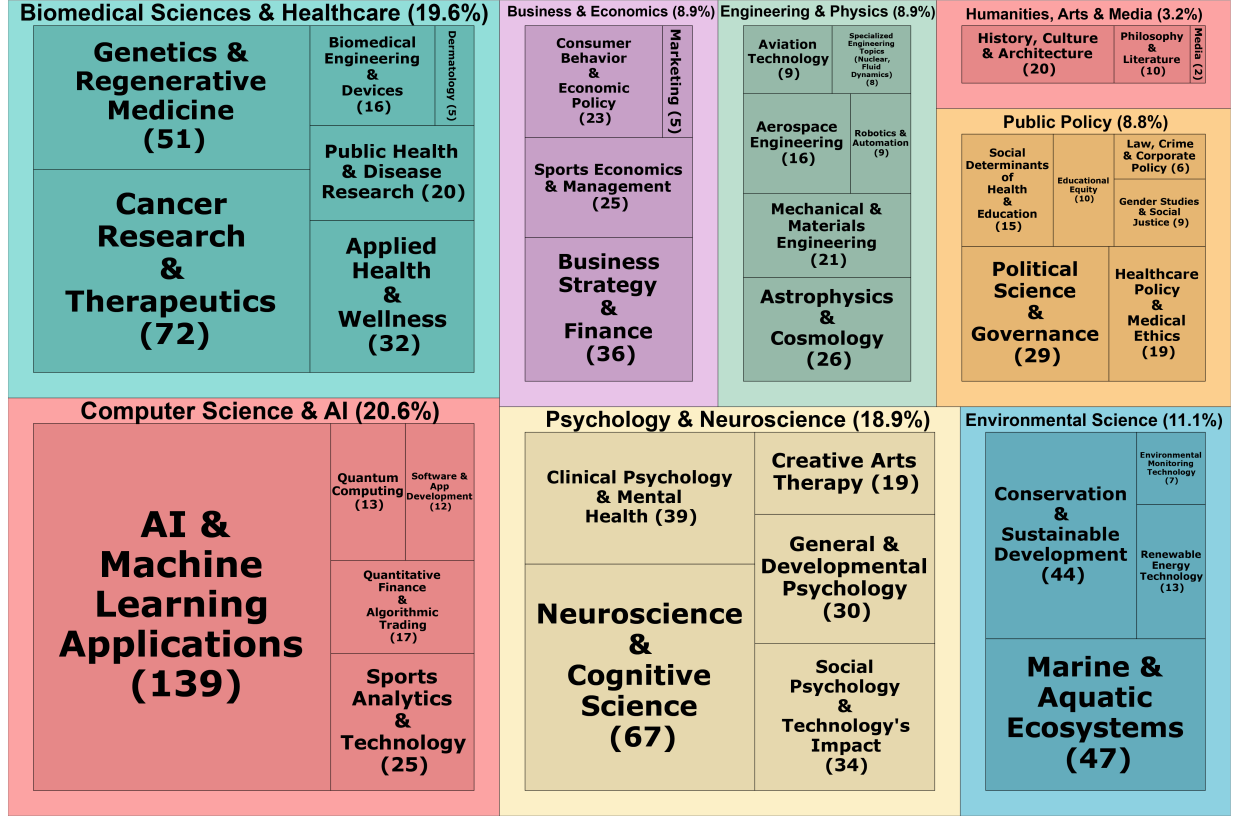
Figure 2: A squarified hierarchical map of the distribution of project topics based on a random sample of $n = 1,000$ Polygence projects. To create this map, we used a customized version of K-means clustering of project topics based on Anthropic's Clio framework [58] and the open-source Kura library [59]. The size of each box is proportional to the relative frequency of each topic or topic cluster.

- **Longitudinal interactions**: Projects typically last 4–6 months and capture the entire learning process, including the development of student–tutor relationships over time—a key driver of instructional effectiveness [65–67].

- **Full personalization**: Each project is tailored to the student's academic needs and goals.

- **Multi-modal exchanges**: Each project's dataset encompasses meeting transcripts, shared documents, chat, and other modalities of student–tutor interaction.

- **Outcome-oriented projects**: Over 80% of completed projects culminate in a showcaseable artifact, such as an academic paper, video, podcast, physical prop, or GitHub repository (see Fig. 1).

- **Alignment with student AI usage**: Approximately 80% of the tutoring activities overlap with the top 10 student use cases of AI reported by OpenAI [7].

To understand the distribution of topics and activities covered in Polygence sessions, we apply K-means clustering to the session transcripts using Anthropic's Clio framework [58] and a customized version of the open-source Kura library [59]. First, we focus on a high-level overview of the diverse
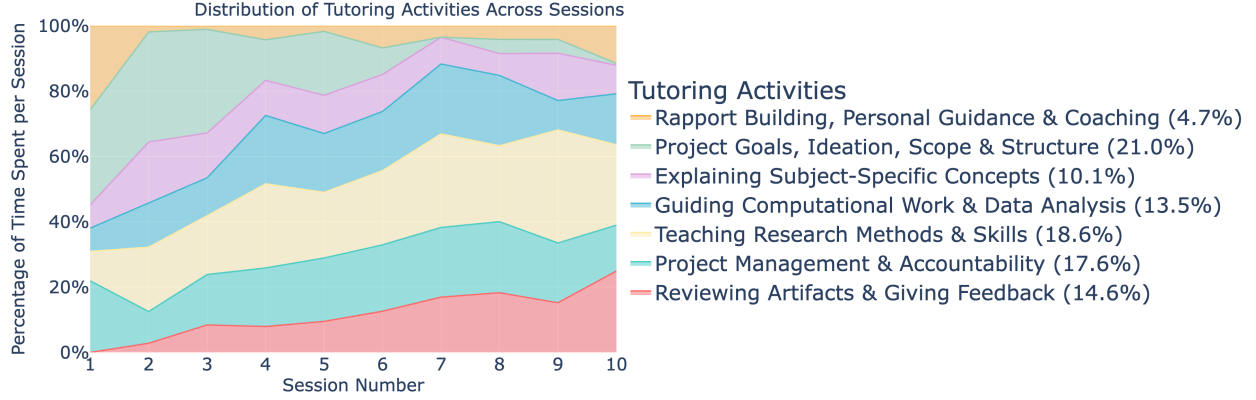
Figure 3: Tutoring activity overview for $n = 195$ completed 10-session projects. Each 1-hour session is segmented into 5-minute chunks, analyzed individually, and hierarchically clustered into 4 levels using using Anthropic's Clio framework [58] and a customized version of the open-source Kura library [59]. The distribution of the 7 top-level tutoring categories are shown across sessions. We find a 78% overlap with the top 10 student usage categories of ChatGPT reported by OpenAI [7].

| Main Category | Subcat 1 | Subcat 2 | Subcat 3 | Subcat 4 | Subcat 5 |
|---|---|---|---|---|---|
| Rapport Building, Personal Guidance & Coaching | Building Rapport & Personal Connection (56.8%) | Providing Encouragement & Emotional Support (20.9%) | Advising on College & Career Pathways (15.5%) | Guiding Student Reflection & Goal Setting (6.8%) | |
| Project Goals, Ideation, Scope & Structure | Structuring & Outlining Project Deliverables (37.8%) | Developing & Refining Methodology (29.6%) | Brainstorming & Refining Project Ideas (17.2%) | Defining & Refining Scope, Goals & Topics (15.4%) | |
| Explaining Subject-Specific Concepts | Explaining STEM Concepts (40.2%) | Explaining AI/ML & Data Science Concepts (26.3%) | Explaining Social Science & Humanities Concepts (18.2%) | Explaining Foundational Math & Statistical Theory (8.6%) | Explaining Business, Finance & Marketing Concepts (6.7%) |
| Guiding Computational Work & Data Analysis | Guiding Coding & Debugging (39.6%) | Guiding Data Analysis & Interpretation (37.6%) | Assisting with Technical Setup & Tools (13.6%) | Guiding Data Sourcing & Preparation (9.2%) | |
| Teaching Research Methods & Skills | Guiding Academic Writing & Structure (48.9%) | Guiding Literature Review & Source Analysis (26.6%) | Instruction on Citation, Formatting & Ethics (17.2%) | Developing Arguments & Integrating Evidence (5.2%) | Coaching Presentation & Public Speaking Skills (2.2%) |
| Project Management & Accountability | Scheduling Sessions & Coordinating Logistics (46.9%) | Managing Timelines & Deadlines (22.6%) | Assigning & Clarifying Tasks (14.3%) | Guiding Final Submission Process (12.4%) | Navigating Platform & Admin Procedures (3.9%) |
| Reviewing Artifacts & Giving Feedback | Providing General Writing Feedback (42.4%) | Collaborative Editing & Word-Smithing (23.9%) | Refining Paper & Project Structure (16.1%) | Giving Feedback on Presentations & Visuals (8.8%) | Reviewing Overall Project Progress (8.8%) |

Table 1: Main tutoring activities, with subcategories and percentages listed across columns obtained via hierarchical K-means clustering.

set of topics covered by Polygence projects. We randomly sample $n = 1,000$ projects and cluster them based on their project topics. Fig. 2 shows that the majority of the topics cover computer science, biomedical topics, psychology and neuroscience. The most popular topic is AI and machine learning followed by cancer research.

Next, we focus on the various tutoring activities that are undertaken during the live sessions. We select a representative subset of $n = 195$ completed projects, which typically have 9-10 full sessions. We then divide each 1-hour session into 10–15 segments, each about 5 minutes long ($n = 24,587$

chunks in total) and recursively classify them to obtain both lower-level (more specific) and higher-level categories of topics and activities.

Fig. 3 highlights how the distribution of activities between the tutor and the student shifts dynamically as the project progresses across the ten sessions. Rapport and relationship building represents up to 25% of the first session and up to 10% of the last session and remains an important part of the mentoring process throughout the ten sessions (e.g. providing encouragement and emotional support). Ideation and setting project goals and scope are heavily prioritized upfront, but remain 10-20% of the time budget up to mid-way through the project – showing the highly iterative nature of this process. Subject-specific tutoring (e.g. explaining AI concepts) and technical guidance (e.g. coding and data analysis) account for about a quarter of all time spent on the project and is evenly distributed across the entire process. Reviewing artifacts (e.g. writing and presentation feedback) gradually increases in importance as the project progresses. These details highlight the highly dynamic and iterative nature of these projects as well as the richness of longitudinal tutor-student interactions that are difficult to capture in simple, static rules or heuristics. Table 1 gives a detailed breakdown of each high-level activity into more specific activities.

This detailed activity map of a representative sample of tens of thousands of 5-minute session fragments allows us to map the tutoring activity on the Polygence platform to the top 27 categories of ChatGPT usage by 18-24 year old students, as reported by OpenAI [7]. By introducing the OpenAI-defined categories into our clustering process, we find that approximately 78% of our data covers the top 10 ChatGPT use cases by students. These include starting papers/projects, brainstorming creative projects, exploring topics, editing writing, solving mathematical problems, conducting academic research, tutoring, and drafting essays (see Appendix C for more details).

## 3.2 Data privacy, User Consent, Anonymization, and Data Security

Data was collected and processed consistent with our Terms of Use and Privacy Policy, and reflects participant opt-outs. In addition to the rights granted through these policies, we secure consent for every recorded session with explicit, real-time notification of participants about recording at the start of each call. Given that some tutor profiles may be publicly associated with Polygence, we took certain steps to anonymize their identities in the transcripts and remove information that could personally identify them (PII) before training, including their personal background, education, life experiences, academic expertise, personal opinions, etc. This data anonymization was done on Polygence's internal servers. All subsequent experimentation and model evaluation was conducted exclusively on our service provider's enterprise-grade platforms. Those service providers are contractually required to maintain confidentiality, ensure appropriate data security protections for the data, and are prohibited from using Polygence's data for their own AI model training.

## 3.3 Transcription, Diarization, and Cleaning

Our training pipeline is built on audio from tutoring sessions, recorded with the explicit consent of all participants as outlined in our consent framework (Section 3.2). High-quality training data is a prerequisite for effective post-training. Prior work shows that quality outweighs quantity during fine-tuning (e.g., [68]). We apply a multi-step pipeline to clean and post-process audio and transcripts from tutoring sessions before text-based fine-tuning.

**Data pre-filtering.** We retain projects with a single, uninterrupted tutor–student match. We utilize dual-track Zoom recordings (separate tracks per speaker), a feature that enables higher-fidelity transcription and reliable diarization.
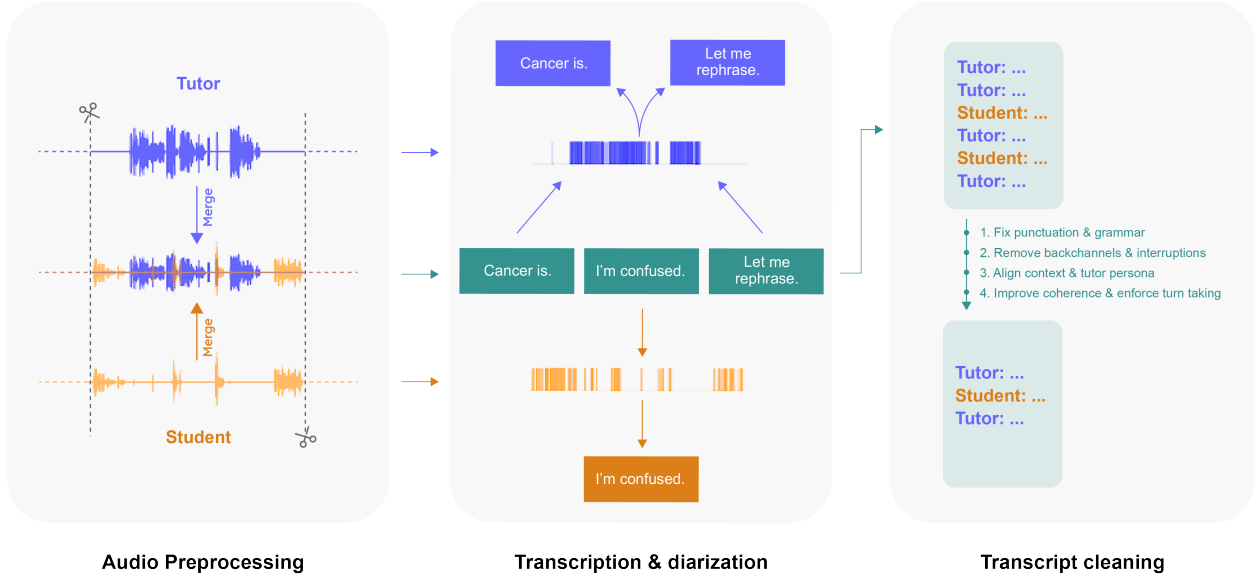
Figure 4: End-to-end transcript processing pipeline. Dual-track audio is merged and trimmed, then transcribed with high fidelity. Speaker activity masks enable accurate diarization, followed by a multi-step cleaning process (fix punctuation and grammar, remove backchannels and interruptions, align context and tutor persona, and improve coherence and enforce turn taking) to yield polished tutor–student transcripts.

**Audio transcription and diarization.** As illustrated in Fig. 4:

1. Assign single-speaker tracks to TUTOR or STUDENT via fuzzy matching of Zoom display names to known participants.

2. Merge tracks into one file and trim leading/trailing silence for all files to reduce transcription hallucinations from long gaps.

3. Transcribe the merged two-speaker audio with ElevenLabs' API [69], which yields high-fidelity text (including fillers) but unreliable speaker tags.

4. Derive per-speaker activity masks from the single-speaker tracks.

5. Attribute each utterance to the speaker with maximal temporal overlap; compile a diarized transcript with one speaker per statement.

6. Generate multiple transcript candidates and select the best using a reasoning model (Gemini 2.5 Pro), mitigating rare but consequential ASR lapses (e.g., dropped segments).

**Transcript cleaning.** Using a reasoning model (Gemini 2.5 Pro), we apply:

1. Normalization of punctuation and grammar.

2. Removal of backchannels that interrupt flow (e.g., "yeah", "gotcha", "uh-huh").

3. Context adjustments and anonymization by removing or reframing references to human tutor identity, personal anecdotes, comments related to physical embodiment, platform names, and program-specific details.

4. Coherence smoothing and enforced turn-taking by merging consecutive same-speaker utterances.

Finally, we conduct human spot checks to ensure transcripts meet the quality bar for post-training.

# 4 Multi-turn Evaluations Using a Fine-Tuned Student Model

Beyond sufficient high-quality data, effective fine-tuning also requires a clear measure of good performance. A persistent issue in evaluating LLMs pedagogical performance is the lack of industry-wide standards for multi-turn evaluations [70, 71]. This contrasts with single-turn evaluations, which have become ubiquitous in the generative AI industry, including within the educational domain [31, 33, 56]. Here we introduce (i) a set of proxies for high-quality pedagogy in multi-turn dialogues, and (ii) a novel multi-turn evaluation protocol that combines a fine-tuned student model and a fine-tuned tutor model to generate a large number of synthetic dialogues. Each of these dialogues are then evaluated using traditional, single-shot evaluation methods and the results are then aggregated. These stochastic methods represent fast, scalable and reproducible measures of LLMs multi-turn performance.

## 4.1 Defining Proxies for High-Quality Pedagogy

A major challenge in education research is the lack of universally accepted pedagogical best practices [33, 72–75]. While certain behaviors are commonly cited as markers of effective teaching (e.g., not giving away the answer, asking questions, balancing talk time), good pedagogy is highly context-dependent, making simple heuristics insufficient to capture the full spectrum of quality instruction. Nevertheless, advancing LLM training requires well-defined evaluation criteria, and thus it is necessary to select at least *some* practical proxies for high-quality pedagogy.

In this preliminary case study, we focus on a handful of straightforward benchmarks, deferring more complex evaluations to future work. Our choice of benchmarks was informed by (i) prior academic studies (see below) and (ii) extensive user feedback from our PolyPilot experiment (see Section 2.2). Specifically, we highlight the following benchmarks:

- **Student talk-time:** The percentage of words uttered by the student relative to the total number of words in the dialogue. In a large-scale randomized controlled trial on the Polygence platform, increasing student talk time was positively associated with outcomes such as academic confidence and participant satisfaction [76].

- **Average number of words per tutor turn:** A measure closely tied to talk time, used to detect the prevalence of "wall-of-text" responses, a well-documented issue in LLM outputs [77–79].

- **Mean questions per interrogative turn:** PolyPilot user feedback highlighted that LLMs often display unnatural questioning styles such as asking a series of questions in the same turn. This metric captures the average number of questions per interrogative turn and serves as a proxy for human-like questioning practices, such as favoring open-ended inquiries. An interrogative turn is defined as a statement that has at least one question and questions are detected via question marks.

- **Number of turns before wrap-up:** LLMs are typically trained to resolve queries as quickly as possible, which can limit their ability to sustain extended and meaningful educational

10

dialogue. This metric is defined as the number of student and tutor turns before the tutor indicates that the discussion is over (e.g. telling the student to do the assigned work and report back when they are done).

- **Uncovering student background and learning context:** LLMs frequently suffer from the "first mile problem"—jumping into explanations without first eliciting information about the student's academic background or learning goals [80] (also see Appendix B). To evaluate this proxy, we track what percentage of all known information about a student is uncovered by a tutor during a dialogue. To do so, we extract and tabulate the information learned about the student in the original human-to-human conversation, which we treat as the theoretical maximum of 100%. We note that even the most skilled teachers would score less than 100% on this benchmark as every conversation is unique and uncovers a slightly different set of facts about the student. Nevertheless, this provides a directional measure of improvement in uncovering relevant information about the student (i.e., higher scores indicate better performance).

- **Checking coding skills for coding projects:** A targeted version of the previous benchmark, this tests whether the model probes a student's coding proficiency (a binary yes/no decision) before initiating technical projects. Our data indicates that nearly all human tutors begin by carefully assessing their students' coding backgrounds and even with thoughtful calibration of project difficulty, about 20% of coding projects still encounter challenges related to scope or student skill.

These benchmarks are intentionally simple, and in principle, any LLM could "ace" them through prompt-engineered rules (e.g., "always ask about the student's coding background before starting a computational project"). However, they serve only as *proxies* for educational quality. Our central thesis is that a model with strong pedagogical capabilities would perform well on these metrics, but optimizing solely for them does not guarantee high-quality teaching.

## 4.2   Fine-Tuned Student Model

Simulating students using generative AI is an active area of research, with prompt engineering being the prominent method of creating different student personas [43, 81–84]. We generally expect students to be easier to simulate than teacher models [34, 43], due to the fact that in most educational dialogues students are expected to *follow instructions* (similar to LLMs) as opposed to *leading with instructions*. However, prompt-engineered student simulators often lack the authenticity and variety that is needed to capture a wide range of different learning personas [34].

In our approach, we use student data to train a student model through parameter-efficient fine-tuning (PEFT) [85, 86]. This approach was inspired by the empirical realization that even human testers with relevant domain-expertise (e.g. the authors of this paper) are low-fidelity impersonators of students. At the same time, with access to extensive dialogue data about how students talk, behave and learn in actual sessions it is possible to train representative student models that mimic human learners with high fidelity, including the conversational benchmarks outlined in Section 4.1.

The steps of training a fine-tuned student model are as follows:

1. We select a large, random sample of projects. We assign non-descript ID's to track each project during the training process.

2. We perform multiple epochs of supervised fine-tuning on the student turns in these transcripts with a system prompt that contains the project ID, stopping the training before overfitting
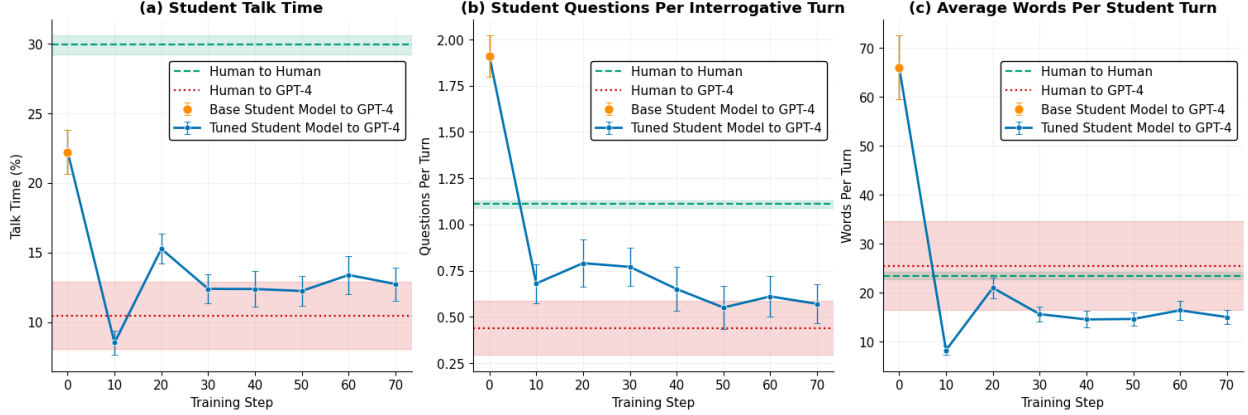
11

Figure 5: Comparing three core conversational statistics (talk time, questions per turn, words per turn) across four different types of dialogues: human to human, human to GPT-4, base student model (Gemini 2.0 Flash) to GPT-4, and tuned student model (Gemini 2.0 Flash tuned on Polygence student data) to GPT-4. The human-to-GPT-4 conversational data was obtained from our PolyPilot experiment (Section 2.2). We observe that simulated conversations between two prompt-engineered base models (large orange dot) produce significantly different ($p < 0.001$) conversational statistics from both dialogues involving only humans (green dashed line) and a human student and AI (red dotted line). Fine-tuning a model on student data (connected blue dots) progressively aligns its conversational statistics with those of actual humans conversing with AI (red dotted line). Humans interacting with AI also produces different conversational statistics than human-to-human conversations, highlighting that the prompt engineered base tutor model (GPT-4) impersonates a human tutor with limited fidelity. These results provide further evidence that simulated dialogues involving a fine-tuned student model approximate human-AI conversations better than conversations generated from two prompt-engineered AI models alone. Error bars and intervals represent 95% confidence intervals calculated with the Student's t-distribution (light green and light red intervals show the error bars for human-to-human and human-to-GPT-4 conversations respectively).

begins. This training on a large amount of student data gives the model a general sense of how students typically communicate and behave, while also ingraining individual student behavior into the model parameters.

3. We randomly choose a small subset of students ($n = 10$) for actual simulations and selectively activate them by using exactly the same prompt –including the associated project-specific ID– that was used during the fine-tuning phase.

4. In addition, we use an LLM to extract a comprehensive set of details that the student reveals about themselves during the student-tutor conversation and include it in the system prompt. This further reinforces the fine-tuned behavior of the student.

5. We run single-shot evaluations on the student model to confirm that it responds to a set of selected queries as expected.

Upon completing training, we carry out qualitative and quantitative assessments to verify that the model's behavior closely reflects that of the original students. Qualitative checks performed by Polygence personnel confirm that these fine-tuned and prompt-reinforced student models respond to questions about their academic background, project interests, etc. in a similar style as their live counterparts.

A summary of our quantitative checks is provided in Fig. 5. Specifically, in Fig. 5 we study the conversational statistics that emerge when this fine-tuned student model is paired with an off-the-shelf model (GPT-4) and compare the results in three other scenarios: (1) human-to-human interactions, (2) human-to-AI interactions, and (3) AI-to-AI interactions without fine-tuning. We specifically picked GPT-4 for this analysis because it allows us to use the conversational human-to-AI data from our PolyPilot experiment (see Section 2.2) for benchmarking.

Crucially, we find that simulated conversations between two AI base models (**Base Student Model to GPT-4**, large orange dot) produce markedly different conversational statistics from both dialogues involving only humans (**Human to Human**, green dashed line) and those between humans and an AI tutor (**Human to GPT-4**, red dotted line). A one-way Analysis of Variance (ANOVA) followed by Tukey's HSD post-hoc tests confirmed that these differences are highly statistically significant.

When interacting with GPT-4, the student base model talks less overall than a human student with a human tutor (mean of 22.2% vs. 29.9%, a difference of $-7.7$ percentage points, $p < 0.001$). At the same time, it talks substantially more than a human student interacting with GPT-4 (22.2% vs. 10.5%, a difference of $+11.8$ p.p., $p < 0.001$). Furthermore, the base model also includes significantly more questions per turn (mean of 1.91) than both the Human-to-AI scenario ($+1.47$ questions, $p < 0.001$) and the Human-to-Human scenario ($+0.80$ questions, $p < 0.001$). Finally, the base model's turns are substantially longer, averaging 66.1 words, which is significantly more than both the Human-to-AI scenario ($+40.5$ words, $p < 0.001$) and the Human-to-Human scenario ($+42.6$ words, $p < 0.001$). Notably, while the base model is a clear outlier on this metric, the average turn lengths for the Human-to-AI (25.5 words) and Human-to-Human (23.4 words) conditions were not statistically different from one another ($p = 0.589$).

In contrast, after training on our student data, the fine-tuned student model displays conversational statistics that begin to mimic those of humans interacting with AI. These results indicate that simulated dialogues involving a fine-tuned student model provide a more accurate representation of human-AI conversations than those generated from two prompt-engineered AI models alone.

## 4.3 Multi-Turn Evaluations

Our student simulator allows us to run multi-turn evaluations at scale on any LLMs. First, we pick a single fine-tuned student model by selecting a training checkpoint that performs well on our evaluations and that is well aligned with human-to-AI conversational statistics (Fig. 5). We then initiate a conversation between the student model and a tutor model that we are looking to evaluate by feeding the output of the tutor model into the student model and vice versa. This continues until either the end of the conversation is detected by a dedicated LLM judge (Gemini 2.5 Flash), or the number of turns hit a predetermined limit (this ensures that conversations do not continue indefinitely). Once the conversation is finished, it is analyzed for conversational statistics and passed on to a reasoning LLM judge for a single-shot evaluation (see e.g. 'Uncovering student background & learning context' in Section 4.1). We typically repeat this 10 times for each of the 10 student models, thereby ensuring that each evaluation point corresponds to 100 simulated multi-turn conversations. We chose these numbers to limit variability and to ensure reproducibility between consecutive evaluation experiments.

# 5 Benchmarking LLMs Against Human Performance

Before fine-tuning LLMs with authentic learning data, we first sought to evaluate the performance gap between state-of-the-art LLMs and humans in terms of conversational and pedagogical capabilities. To establish human performance for the six selected benchmarks outlined in the previous section, we analyze 80,000 hours of data from Polygence. We also evaluate off-the-shelf LLMs for the same benchmarks. Our methodology for performing these evaluations using a fine-tuned student model is discussed in Section 4.3.
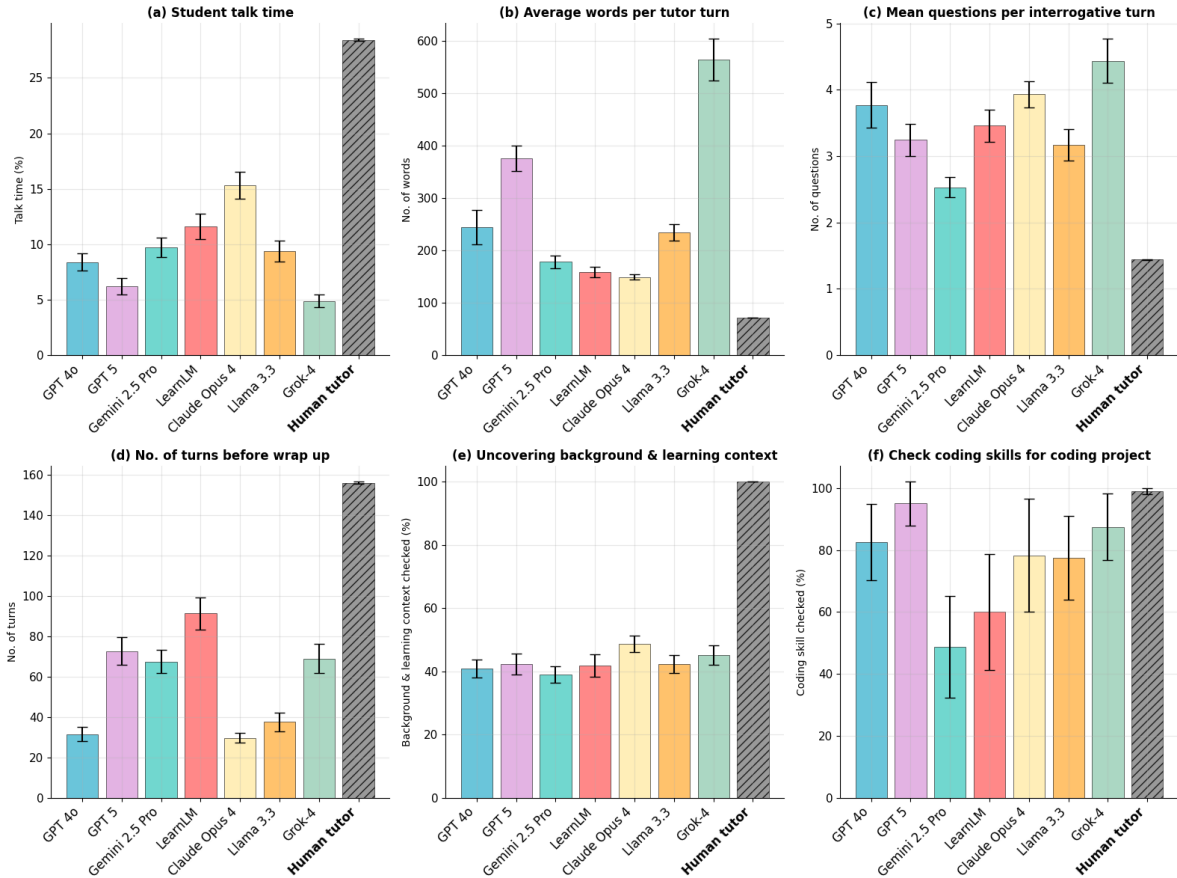


Figure 6: Benchmarking results for six conversational and pedagogical metrics comparing human tutors (hatched gray) with state-of-the-art LLMs from OpenAI, Google, Anthropic, Meta, and XAI. Benchmarks include (a) student talk time, (b) average words per tutor turn, (c) mean questions per interrogative turn, (d) number of turns before wrap-up, (e) uncovering student background and learning context, and (f) checking coding skills for coding projects. Human tutors consistently outperform LLMs across these benchmarks.

Fig. 6 shows our benchmarking results for humans (hatched gray) and a range of different state-of-the-art LLMs from OpenAI, Google, Anthropic, Meta, and XAI. In Fig. 6(a), we find that in human-to-human tutor-led dialogues students speak close to 30% of the time. Off-the-shelf models are typically more verbose than humans, leaving only 5-15% of talk time to students. A hallmark of dynamic dialogues is when students and tutors take frequent turns, rather than one party monopolizing the conversation. Fig. 6(b) shows that on average, human tutors speak only 72 words before passing the turn to the student. In contrast, most off-the-shelf models average 150-300

words per turn. Another important technique in human tutoring is asking open-ended questions and then pausing and letting students answer. Statistically speaking, human tutors typically ask 1-2 questions per interrogative turn (on average 1.5 in our dataset). In Fig. 6(c) we observe the inherent tendency of LLMs to ask a high number of questions (3-4 per interrogative turn). Another shortcoming of LLMs is the rapid drive towards 'resolving' conversations. Fig. 6(d) shows that most LLMs tend to end student conversations in 30-80 turns, which contrasts with human tutors, whose average session length is closer to 150-160 turns (note that session lengths are widely distributed, with a mean of 56 minutes and a standard deviation of 16 minutes despite the nominal 1-hour length).

Fig. 6(e) shows our most complex evaluation, which quantifies the extent to which off-the-shelf frontier models are able to uncover relevant information about the student's background and learning context. We find that most models score in the 40-45% interval, often missing key context (such as coding skills, goals, needs, motivation, etc.) about the student.

Fig. 6(f) shows results for a more specific version of the previous benchmark – the tendency of models to check the student's the coding background in cases where their project requires coding (a binary yes/no decision). Human data shows that tutors check students' coding background virtually any time the project involves coding. In contrast, we find that most off-the-shelf models perform in the 50-80% range, with only a few thinking models doing better.

Finally, we note that these benchmarking experiments have generally proven reproducible across runs and across different versions of our student model. For example, we consistently find that among the models we analyzed, Claude produces the highest student talk time, LearnLM maintains the conversation the longest, and GPT-5 (and previously o1) performs the best on checking the coding skills of the student.

# 6 Fine-Tuning on Authentic Learning Data

In this section, we demonstrate that performing parameter-efficient fine-tuning [85, 86] on state-of-the-art frontier models using high-quality post-training data (see Section 3) improves their performance on the benchmarks outlined in Section 4.1. Specifically, we fine-tune Google's Gemini 2.5 Flash and OpenAI's GPT 4o-08-06, both of which are state-of-the-art models and are the most recent non-thinking model versions[1] from their respective providers with API access to fine-tuning.

## 6.1 Benchmarking through automatic evaluations

Fig. 9 shows a representative set of results of our fine-tuning experiments across the six different benchmarks under consideration. The dashed green line shows results from the human training data for reference (note that these results are slightly different from Fig. 6 due to the effects of the cleaning process outlined in Section 3.3). We see that fine-tuning improves model performance on every benchmark, bringing the resulting values closer to the human training data benchmarks. Specifically, we find that student talk time increases, the average number of words per tutor turn goes down, the mean number of question per interrogative turn settles in the 1-2 question region[2], and the number of conversation turns increases. We also find that models show marked improvements on the complex evaluation of uncovering student background and learning context, although

---

[1]For Gemini 2.5 Flash we set the thinking budget to zero tokens.

[2]Note that GPT-4o-2024-0806 ask comparatively fewer questions (around 1.3) per interrogative turn than other off-the-shelf models. Nonetheless, fine-tuning changes how it is formulating its questions. We found that later models from OpenAI – such as GPT-4o-2024-11-20, o1, and GPT-5– regress to asking 3-4 questions per interrogative turn.
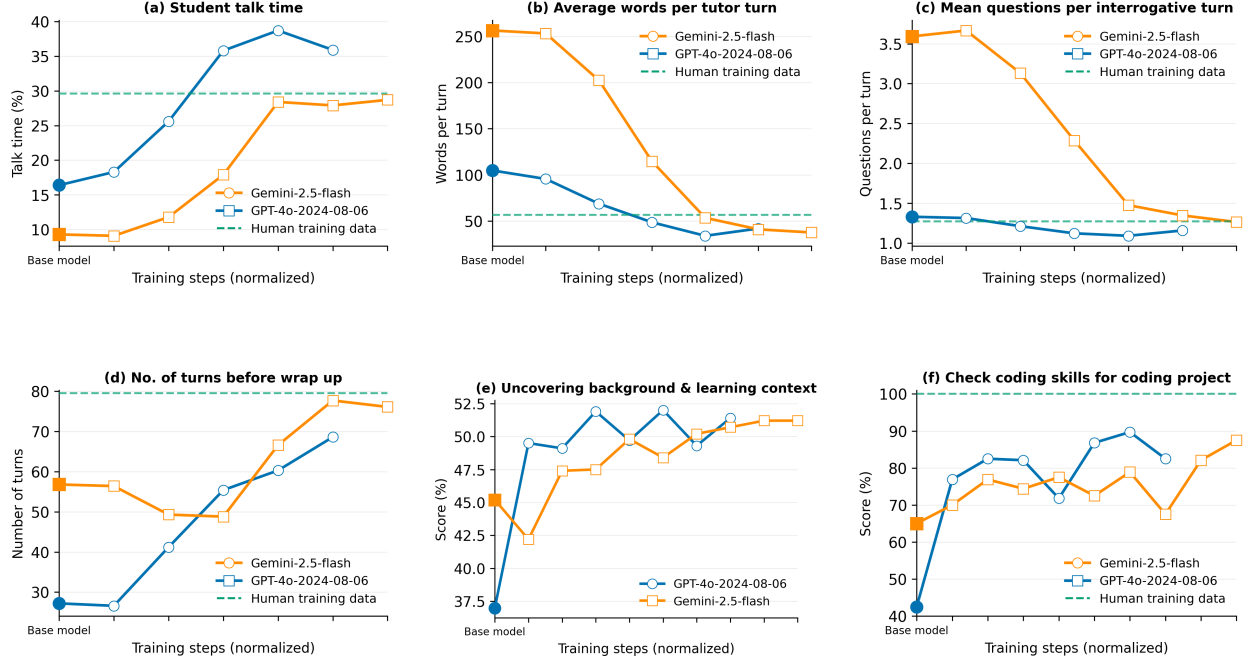
Figure 7: Fine-tuning results on six conversational and pedagogical benchmarks for Gemini 2.5 Flash (orange) and GPT-4o-2024-08-06 (blue). The dashed green line indicates performance from human training data. Benchmarks include (a) student talk time, (b) average words per tutor turn, (c) mean questions per interrogative turn, (d) number of turns before wrap-up, (e) uncovering student background and learning context, and (f) checking coding skills for coding projects. Fine-tuning consistently improves model performance, with simpler conversational metrics (a–d) converging more quickly than the more complex pedagogical benchmarks (e–f).

the absolute values are still far from 100%. We also find a corresponding improvement on the coding skill check benchmark.

The $x$-axes in Figs. 9(a)-(d) represent a smaller number of training steps than Figs. 9(e)-(f). This aligns with our empirical observation that training models to uncover student background, learning context, and coding skills requires significantly more time than training them on simpler conversational improvements. Finally, the near-monotonic improvement on all benchmarks as training progresses allows the selection of checkpoints, which show clear improvement on all benchmarks simultaneously (see Appendix D for an example).

## 6.2 Human Evaluations

In addition to the large-scale automated evaluations outlined above, we also performed limited human evaluations of the resulting fine-tuned models to confirm the reported performance improvements. In these manual experiments, human testers (internal team members) impersonated students and entered into dialogues with the fine-tuned tutor model. This allowed us to directly observe reduced verbosity and enhanced turn taking, a marked shift towards fewer but more open-ended questions, a more natural flow of conversations, and clear improvement in context-setting and in-depth understanding of student background before diving into further tutoring activity. A rigorous, large-scale human evaluation of our fine-tuned models will be addressed in an upcoming

report.

# 7 Conclusion

In this technical report, we analyzed the fundamental limitations of prompt-engineering LLMs for education and outlined the importance of post-training frontier models on authentic learning data. We introduced a novel framework for performing multi-turn evaluations on LLMs using a fine-tuned student model trained on authentic student data. We used our framework to quantify the gap between frontier models and human tutors in terms of six conversational and pedagogical benchmarks. We also showed that fine-tuning on authentic tutor data improved the performance of frontier models on all of our benchmarks.

# 8 Outlook

This preliminary report focused on the supervised fine-tuning of frontier models, which is only the first step in post-training of LLMs. Encouraged by these early results, we are now focused on realizing the full benefit of post-training through reinforcement learning from human feedback (RLHF) [26], which is a natural next step given the availability of authentic learning dialogues involving humans.

Our early efforts focused on establishing the simple, measurable benchmarks outlined in this report, but more sophisticated evaluations are needed to capture the full richness of human pedagogy. Specifically, we will prioritize creating benchmarks that capture the nuances of longitudinal interactions between students and tutors. These types of interactions are only possible over extended periods of time and are critical in fostering rapport and driving tangible learning outcomes.

We also note that evaluations in this report were focused on scalable and automated processes and (limited) human feedback from non-learners. Our aim is to scale and quantify *student feedback* on actual model performance by incorporating post-trained models into the student journey on the Polygence platform.

# 9 Acknowledgements

## Author Information

Dr. Demszky's contribution to this publication was as a consultant and was not part of her Stanford University duties and responsibilities. Janos Perczel is the co-founder and CEO of Polygence, and Jin Chow is the co-founder and COO of Polygence.

## Author Contributions

J.P. conceived the project, developed the methodology, and executed the technical work. D.D. contributed technical advice and ideas. J.C. assisted with data cleaning and model testing. J.P., J.C. and D.D. wrote the paper.

## References

[1] Benjamin S. Bloom. The 2 Sigma Problem: The Search for Methods of Group Instruction as Effective as One-to-One Tutoring. *Educational Researcher*, 13(6):4–16, 1984. URL `https://journals.sagepub.com/doi/10.3102/0013189X013006004`.

[2] World Bank, UNESCO, UNICEF, USAID, FCDO, and Bill & Melinda Gates Foundation. The State of Global Learning Poverty: 2022 Update. `https://www.worldbank.org/en/topic/education/publication/state-of-global-learning-poverty`, 2022.

[3] Greg Kestin, Kelly Miller, Anna Klales, Timothy Milbourne, and Gregorio Ponti. AI tutoring outperforms in-class active learning: an RCT introducing a novel research-based design in an authentic educational setting. *Scientific Reports*, 15(1):17458, June 2025. ISSN 2045-2322. doi: 10.1038/s41598-025-97652-6. URL `https://doi.org/10.1038/s41598-025-97652-6`.

[4] Rose E. Wang, Ana T. Ribeiro, Carly D. Robinson, Susanna Loeb, and Dora Demszky. Tutor CoPilot: A Human-AI Approach for Scaling Real-Time Expertise, 2025. URL `https://arxiv.org/abs/2410.03017`.

[5] Martín De Simone, Federico Tiberti, Wuraola Mosurola, Federico Manolioco, Maria Barron, and Eliott Dikoru. From Chalkboards to Chatbots: Transforming Learning in Nigeria, One Prompt at a Time, 2025. URL `https://blogs.worldbank.org/en/education/From-chalkboards-to-chatbots-Transforming-learning-in-Nigeria`.

[6] Owen Henkel, Hannah Horne-Robinson, Nessie Kozhakhmetova, and Amanda Lee. Effective and Scalable Math Support: Experimental Evidence on the Impact of an AI-Math Tutor in Ghana. In *International Conference on Artificial Intelligence in Education*, pages 373–381. Springer, 2024.

[7] OpenAI Global Affairs. Building an AI-Ready Workforce: A Look at College Student Chat-GPT Adoption in the US, 2025. URL `https://cdn.openai.com/global-affairs/openai-edu-ai-ready-workforce.pdf`.

[8] CivicScience. ChatGPT Is Still Leading the AI Wars but Google Gemini Is Gaining Ground, 2025. URL `https://civicscience.com/chatgpt-is-still-leading-the-ai-wars-but-google-gemini-is-gaining-ground/`.

[9] Chunpeng Zhai, Santoso Wibowo, and Lily D. Li. The Effects of Over-Reliance on AI Dialogue Systems on Students' Cognitive Abilities: A Systematic Review. *Smart Learning Environments*, 11, 2024. doi: 10.1186/s40561-024-00316-7. URL `https://slejournal.springeropen.com/articles/10.1186/s40561-024-00316-7`.

[10] Hamsa Bastani, Osbert Bastani, Alp Sungu, Haosen Ge, Özge Kabakci, and Rei Mariman. Generative AI Without Guardrails Can Harm Learning: Evidence from High School Mathematics. *Proceedings of the National Academy of Sciences*, 122(26):e2422633122, 2025. doi: 10.1073/pnas.2422633122.

[11] Nataliya Kosmyna, Eugene Hauptmann, Ye Tong Yuan, Jessica Situ, Xian-Hao Liao, Ashly Vivian Beresnitzky, Iris Braunstein, and Pattie Maes. Your Brain on ChatGPT: Accumulation of Cognitive Debt when Using an AI Assistant for Essay Writing Task, 2025. URL `https://arxiv.org/abs/2506.08872`.

[12] Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. Finetuned Language Models Are Zero-Shot Learners, 2022. URL `https://arxiv.org/abs/2109.01652`.

[13] Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction Tuning with GPT-4, 2023. URL `https://arxiv.org/abs/2304.03277`.

[14] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan. Training a Helpful and Harmless Assistant with Reinforcement Learning from Human Feedback, 2022. URL `https://arxiv.org/abs/2204.05862`.

[15] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022. URL `https://arxiv.org/abs/2203.02155`.

[16] Gemini Team, Rohan Anil, et al. Gemini: A family of highly capable multimodal models, 2025. URL `https://arxiv.org/abs/2312.11805`.

[17] Anthropic. The Claude 3 Model Family: Opus, Sonnet, Haiku, 2024. URL `https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model_Card_Claude_3.pdf`.

[18] AI@Meta. Llama 3 Model Card, 2024. URL `https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md`.

[19] OpenAI, Josh Achiam, et al. GPT-4 Technical Report, 2024. URL `https://arxiv.org/abs/2303.08774`.

[20] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023. URL `https://arxiv.org/abs/2310.06825`.

[21] Robert A. Bjork and Elizabeth Ligon Bjork. Desirable Difficulties in Theory and Practice. *Journal of Applied Research in Memory and Cognition*, 9(4):475–479, 2020. doi: 10.1016/j.jarmac.2020.09.003.

[22] Jason M. Lodge, Gregor Kennedy, Lori Lockyer, Amael Arguel, and Mariya Pachman. Understanding Difficulties and Resulting Confusion in Learning: An Integrative Review. *Frontiers in Education*, 3:1–10, 2018. doi: 10.3389/feduc.2018.00049. URL `https://www.frontiersin.org/journals/education/articles/10.3389/feduc.2018.00049`.

[23] Elizabeth Ligon Bjork and Robert A. Bjork. Making Things Hard on Yourself, But in a Good Way: Creating Desirable Difficulties to Enhance Learning. In M. A. Gernsbacher, R. W. Pew, L. M. Hough, and J. R. Pomerantz, editors, *Psychology and the Real World: Essays Illustrating Fundamental Contributions to Society*, pages 56–64. Worth Publishers, New York, NY, USA, 2011.

[24] Benedict Du Boulay. *Intelligent Tutoring Systems That Adapt to Learner Motivation*, pages 103–128. Nova Science Publishers Inc, 10 2018. ISBN 978-1-53614-086-6.

[25] Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep Reinforcement Learning from Human Preferences. *arXiv preprint* , 2017. URL `https://arxiv.org/abs/1706.03741`.

[26] Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-Tuning Language Models from Human Preferences. 2019. URL `https://arxiv.org/abs/1909.08593`.

[27] Stanford University IT. AI Demystified: What Is Prompt Engineering?, 2024. URL `https://uit.stanford.edu/service/techtraining/ai-demystified/prompt-engineering`.

[28] Anthropic. Introducing Claude for Education, 2024. URL `https://www.anthropic.com/news/introducing-claude-for-education`.

[29] AI Education News. Breaking: OpenAI Releases Study Mode, 2024. URL `https://aieducation.substack.com/p/breaking-openai-releases-study-mode`.

[30] AI Education News. Breaking: Google Introduces Guided Mode, 2024. URL `https://aieducation.substack.com/p/breaking-google-introduces-guided`.

[31] LearnLM Team, Abhinit Modi, Aditya Srikanth Veerubhotla, Aliya Rysbek, Andrea Huber, Ankit Anand, Avishkar Bhoopchand, Brett Wiltshire, Daniel Gillick, Daniel Kasenberg, Eleni Sgouritsa, Gal Elidan, Hengrui Liu, Holger Winnemoeller, Irina Jurenka, James Cohan, Jennifer She, Julia Wilkowski, Kaiz Alarakyia, Kevin R. McKee, Komal Singh, Lisa Wang, Markus Kunesch, Miruna Pîslar, Niv Efron, Parsa Mahmoudieh, Pierre-Alexandre Kamienny, Sara Wiltberger, Shakir Mohamed, Shashank Agarwal, Shubham Milind Phal, Sun Jae Lee, Theofilos Strinopoulos, Wei-Jen Ko, Yael Gold-Zamir, Yael Haramaty, and Yannis Assael. Evaluating Gemini in an Arena for Learning. 2025. URL `https://doi.org/10.48550/arXiv.2505.24477`.

[32] Jin Wang and Wenxiang Fan. The Effect of ChatGPT on Students' Learning Performance, Learning Perception, and Higher-Order Thinking: Insights from a Meta-Analysis. *Humanities and Social Sciences Communications*, 12(1):621, 2025. doi: 10.1057/s41599-025-04787-y. URL `https://www.nature.com/articles/s41599-025-04787-y`.

[33] Irina Jurenka, Markus Kunesch, Kevin R. McKee, et al. Towards Responsible Development of Generative AI for Education: An Evaluation-Driven Approach, 2024. URL `https://goo.gle/LearnLM`.

[34] Julia M. Markel, Steven G. Opferman, James A. Landay, and Chris Piech. GPTeach: Interactive TA Training with GPT-Based Students. In *Proceedings of the Tenth ACM Conference on Learning @ Scale (L@S '23)*, page 11, Copenhagen, Denmark, July 2023. ACM. doi: 10.1145/3573051.3593393.

[35] Zihan Wang, Jiaze Chen, Zhicheng Liu, Markus Mak, Yidi Du, Geonsik Moon, Luoqi Xu, Aaron Tua, Kunshuo Peng, Jiayi Lu, Mingfei Xia, Boqian Zou, Chenyang Ran, Guang Tian, Shoutai Zhu, Yeheng Duan, Zhenghui Kang, Zhenxing Lin, Shangshu Li, Qiang Luo, Qingshen Long, Zhiyong Chen, Yihan Xiao, Yurong Wu, Daoguang Zan, Yuyi Fu, Mingxuan Wang, and Ming Ding. AetherCode: Evaluating LLMs' Ability to Win in Premier Programming Competitions, 2025. URL https://arxiv.org/abs/2508.16402v1.

[36] Hanzhao (Maggie) Lin and Heng-Tze Cheng. Gemini Achieves Gold-Medal Level Performance at the International Collegiate Programming Contest World Finals, 2025. URL https://deepmind.google/discover/blog/gemini-achieves-gold-level-performance-at-the-international-collegiate-programming-contest-world-finals/.

[37] Tong Xie, Yuwei Wan, Wei Huang, Zhenyu Yin, Yixuan Liu, Shaozhou Wang, Qingyuan Linghu, Chunyu Kit, Clara Grazian, Wenjie Zhang, Imran Razzak, and Bram Hoex. DARWIN Series: Domain Specific Large Language Models for Natural Science. *arXiv preprint arXiv:2308.13565*, 2023. URL https://arxiv.org/abs/2308.13565v1.

[38] Chuxuan Hu, Yuxuan Zhu, Antony Kellermann, Caleb Biddulph, Suppakit Waiwitlikhit, Jason Benn, and Daniel Kang. Breaking Barriers: Do Reinforcement Post Training Gains Transfer To Unseen Domains?, 2025. URL https://arxiv.org/abs/2506.19733.

[39] Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi-Rui Tam, Keith Stevens, Abdullah Barhoum, Nguyen Minh Duc, Oliver Stanley, Richárd Nagyfi, E. S. Shahul, Sameer Suri, David Glushkov, Arnav Dantuluri, Andrew Maguire, Christoph Schuhmann, Huu Nguyen, and Alexander Mattick. OpenAssistant Conversations – Democratizing Large Language Model Alignment, 2023. URL https://arxiv.org/abs/2304.07327.

[40] Richard S. Sutton. The Bitter Lesson, 2019. URL http://www.incompleteideas.net/IncIdeas/BitterLesson.html.

[41] Justin Vasselli, Christopher Vasselli, Adam Nohejl, and Taro Watanabe. NAISTeacher: A Prompt and Rerank Approach to Generating Teacher Utterances in Educational Dialogues. In *Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)*, pages 772–784, 2023.

[42] Jakub Macina, Nico Daheim, Sankalan Pal Chowdhury, Tanmay Sinha, Manu Kapur, Iryna Gurevych, and Mrinmaya Sachan. Mathdial: A dialogue tutoring dataset with rich pedagogical properties grounded in math reasoning problems, 2023. URL https://arxiv.org/abs/2305.14536.

[43] Anaïs Tack and Chris Piech. The AI Teacher Test: Measuring the Pedagogical Ability of Blender and GPT-3 in Educational Dialogues, 2022. URL https://arxiv.org/abs/2205.07540.

[44] Yann Hicke, Abhishek Masand, Wentao Guo, and Tushaar Gangavarapu. Assessing the efficacy of large language models in generating accurate teacher responses, 2023. URL https://arxiv.org/abs/2307.04274.

[45] Rania Abdelghani, Yen-Hsiang Wang, Xingdi Yuan, Tong Wang, Pauline Lucas, Hélène Sauzéon, and Pierre-Yves Oudeyer. Gpt-3-driven pedagogical agents to train children's curious question-asking skills. *International Journal of Artificial Intelligence in Education*,

34(2):483–518, June 2023. ISSN 1560-4306. doi: 10.1007/s40593-023-00340-7. URL http://dx.doi.org/10.1007/s40593-023-00340-7.

[46] Jakub Macina, Nico Daheim, Lingzhi Wang, Tanmay Sinha, Manu Kapur, Iryna Gurevych, and Mrinmaya Sachan. Opportunities and challenges in neural dialog tutoring, 2023. URL https://arxiv.org/abs/2301.09919.

[47] Katherine Stasaski, Kimberly Kao, and Marti A. Hearst. CIMA: A Large Open Access Dialogue Dataset for Tutoring. In *Proceedings of the Fifteenth Workshop on Innovative Use of NLP for Building Educational Applications*, pages 52–64, 2020.

[48] Andrew Caines, Helen Yannakoudakis, Helena Edmondson, Helen Allen, Pascual Pérez-Paredes, Bill Byrne, and Paula Buttery. The teacher-student chatroom corpus, 2020. URL https://arxiv.org/abs/2011.07109.

[49] Abhijit Suresh, Jennifer Jacobs, Margaret Perkoff, James H. Martin, and Tamara Sumner. Fine-tuning transformers with additional context to classify discursive moves in mathematics classrooms. In Ekaterina Kochmar, Jill Burstein, Andrea Horbach, Ronja Laarmann-Quante, Nitin Madnani, Anaïs Tack, Victoria Yaneva, Zheng Yuan, and Torsten Zesch, editors, *Proceedings of the 17th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2022)*, pages 71–81, Seattle, Washington, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.bea-1.11. URL https://aclanthology.org/2022.bea-1.11/.

[50] Dorottya Demszky and Heather Hill. The NCTE Transcripts: A Dataset of Elementary Math Classroom Transcripts, 2023. URL https://arxiv.org/abs/2211.11772.

[51] National Tutoring Observatory. National Tutoring Observatory, 2025. URL https://nationaltutoringobservatory.org/.

[52] Louis DiPietro. National Tutoring Observatory to Accelerate the Science of Teaching. *Cornell Chronicle.* URL https://news.cornell.edu/stories/2025/01/national-tutoring-observatory-accelerate-science-teaching.

[53] National Science Foundation. Capturing and Leveraging Data from Teacher-Student Interactions to Improve STEM Learning: An Incubator Project. NSF Award #2321499, 2023. URL https://www.nsf.gov/awardsearch/showAward?AWD_ID=2321499. Division of Research on Learning in Formal and Informal Settings (DRL).

[54] Polygence. Polygence: Experiential Research Mentorship Programs for High School Students. https://www.polygence.org/, 2025.

[55] Google DeepMind and Google Cloud. LearnLM: Integrating Pedagogical Capabilities into Gemini, 2025. URL https://cloud.google.com/solutions/learnlm?hl=en.

[56] Abhinit Modi, Aditya Srikanth Veerubhotla, Aliya Rysbek, Andrea Huber, Brett Wiltshire, Brian Veprek, Daniel Gillick, Daniel Kasenberg, Derek Ahmed, Irina Jurenka, James Cohan, Jennifer She, Julia Wilkowski, Kaiz Alarakyia, Kevin R. McKee, Lisa Wang, Markus Kunesch, Miruna Pîslar, Nikhil Joshi, Parsa Mahmoudieh, Paul Jhun, Sara Wiltberger, Shakir Mohamed, Shashank Agarwal, Shubham Milind Phal, Sun Jae Lee, Theofilos Strinopoulos, Wei-Jen Ko, Amy Wang, Ankit Anand, Avishkar Bhoopchand, Dan Wild, Divya Pandya, Filip Bar, Garth Graham, Holger Winnemoeller, Mahvish Nagda, Prateek Kolhar, Renee Schneider,

Shaojian Zhu, Stephanie Chan, Steve Yadlowsky, Viknesh Sounderajah, Yannis Assael, et al. LearnLM: Improving Gemini for Learning, 2025. URL `https://arxiv.org/abs/2412.16429`.

[57] Reuters. Scale AI's Bigger Rival Surge AI Seeks Up to $1 Billion Capital Raise, Sources Say, 2025. URL `https://www.reuters.com/business/scale-ais-bigger-rival-surge-ai-seeks-up-1-billion-capital-raise-sources-say-2025-07-01/`.

[58] Alex Tamkin, Miles McCain, Kunal Handa, Esin Durmus, Liane Lovitt, Ankur Rathi, Saffron Huang, Alfred Mountfield, Jerry Hong, Stuart Ritchie, Michael Stern, Brian Clarke, Landon Goldberg, Theodore R. Sumers, Jared Mueller, William McEachen, Wes Mitchell, Shan Carter, Jack Clark, Jared Kaplan, and Deep Ganguli. Clio: Privacy-Preserving Insights into Real-World AI Use, 2024. URL `https://arxiv.org/abs/2412.13678`.

[59] 567-labs. Kura: A Simple Reproduction of the CLIO Paper for Chat Data Analysis, 2025. URL `https://github.com/567-labs/kura`.

[60] Polygence. High School Research Student Kobe Builds a Wind Tunnel, 2020. URL `https://www.polygence.org/blog/high-school-research-aerodynamics`.

[61] Lydia Olivieri. Lydia Olivieri's Passion for Fashion History and Recreating an 18th Century Dress, 2023. URL `https://www.polygence.org/blog/lydia-olivieri-passion-for-fashion-history`.

[62] Polygence. How can AI be utilized to standardize ballet teaching?, 2025. URL `https://www.polygence.org/projects/research-project-how-can-ai-be-utilized-to-standardize-ballet-teaching-`.

[63] Polygence. High School Neuroscience Research Student Tori Records a 7-Episode Podcast Exploring Human Memories, 2020. URL `https://www.polygence.org/blog/high-school-research-neuroscience-alzheimers`.

[64] Arkapravo Sen. Real-Time Field Concussion Detection System Presentation at Polygence's Symposium, 2024. URL `https://www.youtube.com/watch?v=0aB8Nn5d5Go`.

[65] Giulia Di Lisio, Amaia Halty, Ana Berástegui, Antonio Milá Roa, and Alba Couso Losada. The longitudinal associations between teacher-student relationships and school outcomes in typical and vulnerable student populations: a systematic review. *Social Psychology of Education*, 28(1):144, July 2025. ISSN 1573-1928. doi: 10.1007/s11218-025-10107-8. URL `https://doi.org/10.1007/s11218-025-10107-8`.

[66] Sanne G. A. van Herpen, Femke Hilverda, and Manja Vollmann. A Longitudinal Study on the Impact of Student–Teacher and Student–Peer Relationships on Academic Performance: The Mediating Effects of Study Effort and Engagement. *European Journal of Higher Education*, 2024. doi: 10.1080/21568235.2024.2414760.

[67] Rimm-Kaufman, Sara. Improving students' relationships with teachers to provide essential supports for learning, 2025. URL `https://www.apa.org/education-career/k12/relationships`.

[68] Markus Krause and Nancy Chang. Achieving 10,000× Training Data Reduction with High-Fidelity Labels, 2025. URL `https://research.google/blog/achieving-10000x-training-data-reduction-with-high-fidelity-labels/`.

[69] ElevenLabs. ElevenLabs Scribe: Speech-to-Text / Audio-to-Text Model, 2025. URL `https://elevenlabs.io/audio-to-text`.

[70] Wai-Chung Kwan, Xingshan Zeng, Yuxin Jiang, Yufei Wang, Liangyou Li, Lifeng Shang, Xin Jiang, Qun Liu, and Kam-Fai Wong. MT-Eval: A Multi-Turn Capabilities Evaluation Benchmark for Large Language Models, 2024. URL `https://arxiv.org/abs/2401.16745`.

[71] Ved Sirdeshmukh, Kaustubh Deshpande, Johannes Mols, Lifeng Jin, Ed-Yeremai Cardona, Dean Lee, Jeremy Kritz, Willow Primack, Summer Yue, and Chen Xing. MultiChallenge: A Realistic Multi-Turn Conversation Evaluation Benchmark Challenging to Frontier LLMs, 2025. URL `https://arxiv.org/abs/2501.17399`.

[72] Mark Dynarski, Roberto Agodini, Sheila Heaviside, Timothy Novak, Nancy Carey, Larissa Campuzano, Barbara Means, Robert Murphy, William Penuel, Hal Javitz, et al. Effectiveness of Reading and Mathematics Software Products: Findings from the First Student Cohort. Technical Report NCEE 2007-4005, U.S. Department of Education, Institute of Education Sciences, National Center for Education Evaluation and Regional Assistance, Washington, DC, 2007. URL `https://files.eric.ed.gov/fulltext/ED496015.pdf`.

[73] Junlei Li and David Klahr. Cognitive Research and Elementary Science Instruction: From the Laboratory, to the Classroom, and Back. *Journal of Science Education and Technology*, 14(2): 217–238, 2005. doi: 10.1007/s10956-005-4425-8.

[74] David Klahr. What Do We Mean? On the Importance of Not Abandoning Scientific Rigor When Talking About Science Education. *Proceedings of the National Academy of Sciences*, 110(Supplement_3):14075–14080, 2013. doi: 10.1073/pnas.1304115110.

[75] Amy Ogan. Designing Culturally-Relevant Educational Technology at a Global Scale, 2023. URL `https://learnlab.org/learning-science-and-engineering-seminar/`.

[76] Dorottya Demszky and Jing Liu. Measuring Conversational Moves in 1:1 Tutoring. In *L@S '23: Proceedings of the Tenth ACM Conference on Learning at Scale*. ACM, July 2023. doi: 10.1145/3573051.3597418.

[77] Eleftheria Briakou, Zhongtao Liu, Colin Cherry, and Markus Freitag. On the Implications of Verbose LLM Outputs: A Case Study in Translation Evaluation, 2024. URL `https://arxiv.org/abs/2410.00863`.

[78] Sania Nayab, Giulio Rossolini, Marco Simoni, Andrea Saracino, Giorgio Buttazzo, Nicolamaria Manes, and Fabrizio Giacomelli. Concise Thoughts: Impact of Output Length on LLM Reasoning and Cost, 2024. URL `https://arxiv.org/abs/2407.19825`.

[79] Yushi Zhang et al. Verbosity Compensation Behavior of Large Language Models, 2024. URL `https://arxiv.org/abs/2411.07858`.

[80] Clayton Cohn, Surya Rayala, Namrata Srivastava, Joyce Horn Fonteles, Shruti Jain, Xinying Luo, Divya Mereddy, Naveeduddin Mohammed, and Gautam Biswas. A Theory of Adaptive Scaffolding for Large Language Model-Based Pedagogical Agents, 2025. URL `https://arxiv.org/abs/2508.01503`.

[81] Tao Wu, Jingyuan Chen, Lin Wang, Mengze Li, Yumeng Zhu, Ang Li, Kun Kuang, and Wu Fei. Embracing Imperfection: Simulating Students with Diverse Cognitive Levels Using LLM-based Agents, 2025. URL `https://arxiv.org/abs/2505.19997`.

24

[82] Haoxuan Li, Jifan Yu, Xin Cong, Yang Dang, Daniel Zhang-li, Yisi Zhan, Huiqin Liu, and Zhiyuan Liu. Exploring LLM-based Student Simulation for Metacognitive Cultivation, 2025. URL https://arxiv.org/abs/2502.11678.

[83] Yash Mahajan, Naman Bansal, Eduardo Blanco, and Santu Karmaker. ALIGN-SIM: A Task-Free Test Bed for Evaluating and Interpreting Sentence Embeddings through Semantic Similarity Alignment. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 7393–7428, Miami, Florida, USA, November 2024. Association for Computational Linguistics.

[84] Bihao Hu, Jiayi Zhu, Yiying Pei, and Xiaoqing Gu. Exploring the potential of LLM to enhance teaching plans through teaching simulation. *npj Science of Learning*, 10(1):7, February 2025. ISSN 2056-7936. doi: 10.1038/s41539-025-00300-x. URL https://doi.org/10.1038/s41539-025-00300-x.

[85] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. URL https://arxiv.org/abs/2106.09685.

[86] John Schulman and Thinking Machines Lab. Lora without regret. *Thinking Machines Lab: Connectionism*, 2025. doi: 10.64434/tml.20250929. URL https://thinkingmachines.ai/blog/lora/.

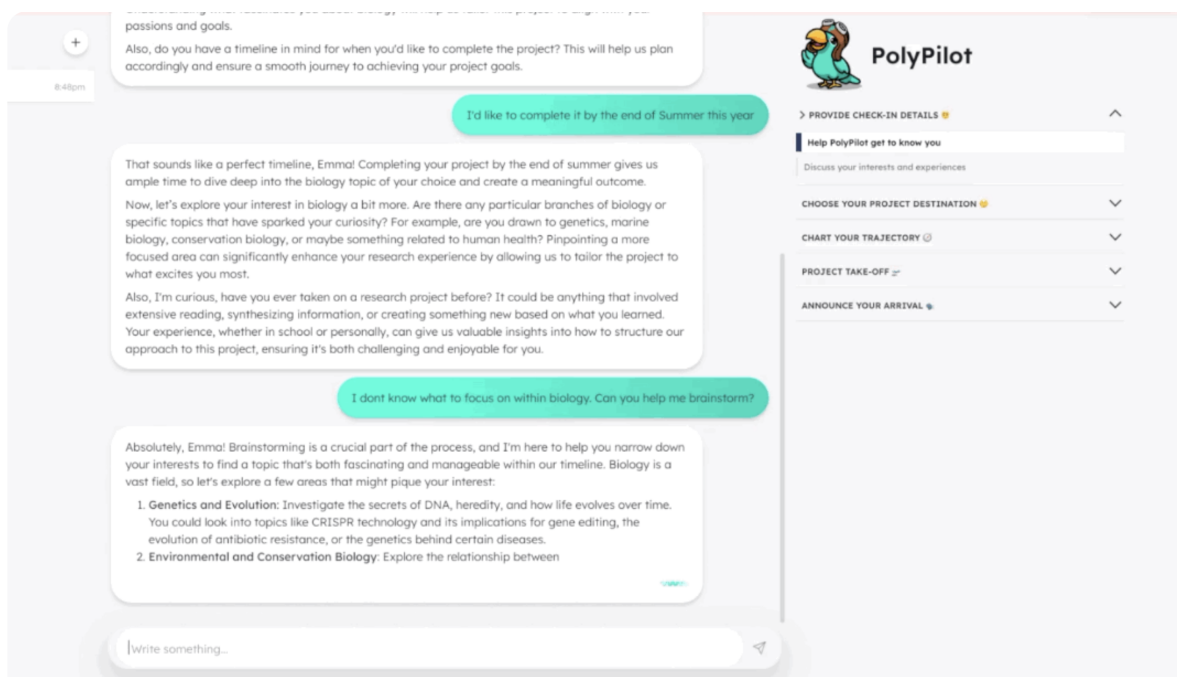# A  PolyPilot Experiment: Building a Prompt-Engineered Tutor



Figure 8: Screenshot of the PolyPilot interface. The right-hand side shows the different pre-defined stages in the student-journey

PolyPilot was a high-conviction product bet in early 2024 to build an effective AI tutor for project-based learning ("research mentor") by prompt engineering GPT-4. PolyPilot was designed to be a ChatGPT-like interface, but adapted to the needs of students looking to complete a long-term project (see screenshot in Fig. 8). The student journey was broken up into five distinct stages:

- **Provide check-in details** – This stage focused on establishing context about the student's background, interests, and goals.

- **Choose your project destination** – This stage helped students scope out their project idea.

- **Chart your trajectory** – This stage guided the students through the background work needed to complete the project.

- **Write your first draft** – This stage focused on starting the writing process.

- **Announce your arrival** – This stage helped students finish writing their artifact and showcasing it to the public.

For each stage we customized our prompts and upon detecting stage completion, the prompt for the next session was loaded on the backend. We iteratively refined the prompts to better align GPT-4 with good pedagogy and to handle edge cases. Initially, we experimented with just five prompts for the five distinct stages and the prompts quickly ballooned to 500-1000 words in length. We observed severe limitations in GPT-4's ability to follow all instructions. Later, we broke down each of the 4 stages into 3-6 distinct sub-stages and progressively refined the prompts to be only 200-300 words. These improved the performance, but the overall pedagogical capabilities of the product remained limited. To improve product quality and reliability, we introduced human checkpoints where an actual tutor checked the progress of the students and decided whether the student would need to do more work before proceeding to the next stage. We found that these human checkpoints were helpful, but didn't fully address the underlying issues. The product was used by $n = 71$ students generating valuable data about student-AI interactions.

## B  Observations on Anthropic's *Learning Mode*, OpenAI's *Study Mode*, and Google's *Guided Learning*

In recent months, multiple model developers have released educational LLMs that have been prompt-engineered to improve their pedagogical behavior. To get a directional sense of their behavior, we manually tested Anthropic's *Learning Mode* [28], OpenAI's *Study Mode* [29], and Google's *Guided Learning* [30] (which integrates LearnLM [31]). Topics ranged from solving simple quadratic equations, to learning Latin, to studying the advanced physics of black holes. Goals ranged from learning new topics, to refining existing knowledge, to writing research papers. Among others, we found the following patterns:

- **Missing learning context**: All three models spend practically no time establishing the learning context, such as the pre-existing level of understanding, the learning goal or the motivation of the student. As a result, models dive quickly into discussions of specific topics that may be inappropriate for the level of the student or drive towards outcomes that do not reflect the learning goal of the student. In cases where the model asked questions about the student's existing level of understanding, we found limited evidence that this was taken into account during the conversation.

- **Multiple-choice-style questioning:** All three models appear to struggle with asking open-ended questions, instead defaulting to multiple-choice style questions (typically 3) that dramatically restrict the space of choices. We theorize that this is a consequence of the underlying off-the-shelf model behavior that is learned from data containing bullet-point-style, overly-structured responses. This typically narrows the topics and/or course of action quickly without appropriately exploring the ideas and goals of the student.

- **One question per statement:** We found that Gemini and GPT typically ask exactly one question per statement (occasionally two), usually placed at the end of the statement, which we theorize is a simple consequence of a prompt instruction. While this behavior is an improvement over the large number of questions off-the-shelf models typically ask (see Section 5), this simple, static rule represents a coarse approximation of high-quality Socratic questioning. (We found that Claude asks a large number of questions per statement similar to its off-the-shelf version.)

- **Wall of text:** All three models are verbose (though less so than their off-the-shelf versions), which appears to get worse as the conversation progresses. We theorize that the underlying training data ingrains verbosity into models and the effectiveness of steering models towards brevity with prompts diminishes as a progressively larger part of the context window is taken up by the conversation.

- **Inability to deal with confusion:** Not giving away answers has been a prominent (and potentially over-simplified) focal point of improving LLMs for education. GPT appears to still give away answers with minimal prompting. Gemini and Claude immediately and repeatedly rephrase and reexplain the question or problem when faced with confusion – making little effort to understand the source or level of the confusion.

We note that while each of the specific behaviors highlighted above could be improved through targeted prompting, these issues are merely indicators of a larger issue with prompt engineering. As noted previously (and highlighted in the LearnLM team's initial report [33]), prompt engineering alone is unlikely to fully encode the vast and complex space of effective tutoring, as it would require an exhaustive, context-specific rule-based description of all of good pedagogy – an impossible task even with unlimited context length.

# C   Overlap of Polygence Data with Top Student Activities Reported by OpenAI

We map the tutoring activities in our dataset to the top 27 categories of ChatGPT usage by 18-24 year old students, as reported by OpenAI [7]. We find that approximately 78% of our data overlaps with the top 10 ChatGPT use cases by students. Table 2 below shows the specific percentage of students utilizing each of the top 10 activities on OpenAI's platform. We separately show whether that activity is covered by our data.

# D   Selecting a Checkpoint With Improved Performance Across All Benchmarks

In Section 6 we showed that models show general improvement across all selected benchmarks when trained on authentic learning data. In Fig. 9 we choose a specific checkpoint and show that the fine-tuned version outperforms the base model on all six of the benchmarks simultaneously.

| Top 10 ChatGPT Usage Categories by Students | % of Students | Data Overlap |
|---|---|---|
| Starting papers/projects | 49% | ✓ |
| Summarize texts | 48% | ✗ |
| Brainstorm creative projects | 45% | ✓ |
| Explore topics | 42% | ✓ |
| Edit writing | 42% | ✓ |
| Mathematical problem-solving | 39% | ✓ |
| Exam preparation | 36% | ✗ |
| Academic research | 34% | ✓ |
| Tutoring | 32% | ✓ |
| Essay drafting | 32% | ✓ |

Table 2: Overlap of Polygence data with the top 10 ChatGPT use cases by 18-24 year old students, as reported by OpenAI [7].
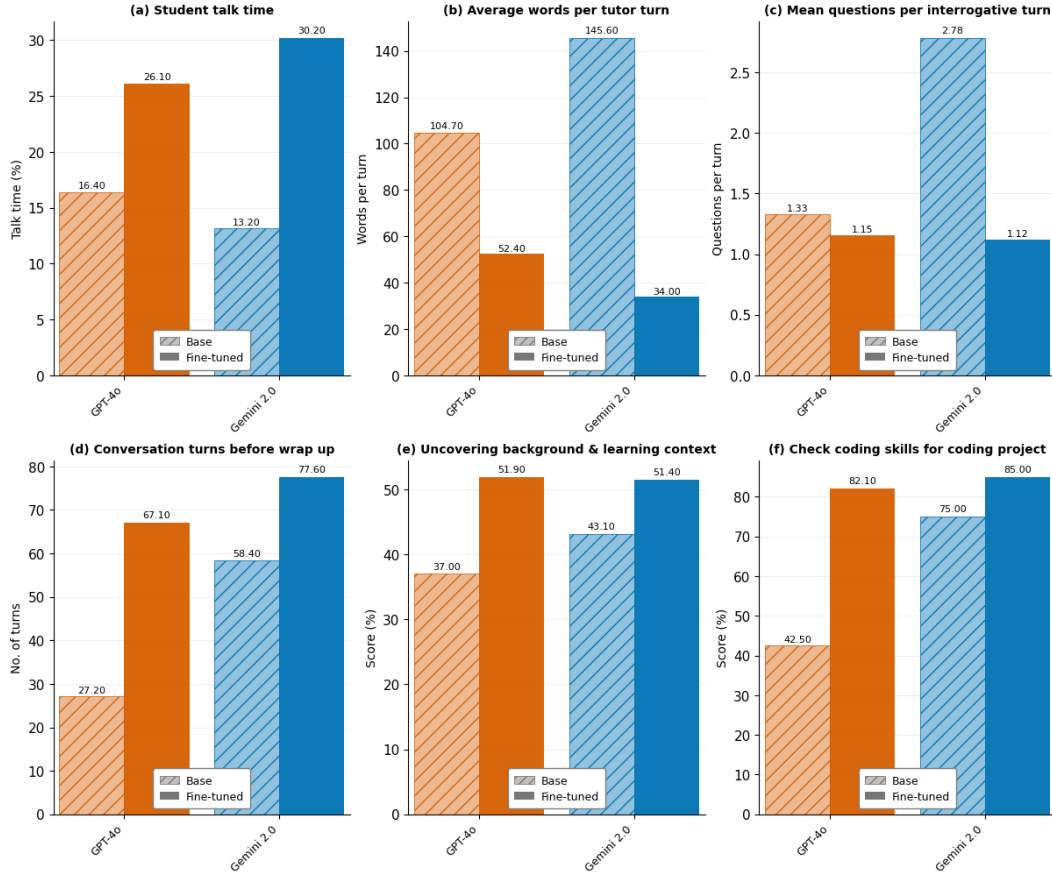


Figure 9: Performance of a specific checkpoint a fine-tuned Gemini 2.0 model on all six benchmarks from Section 4.1. We observe simultaneous improvement across all benchmarks.