

# 000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 TUTORBENCH: A BENCHMARK TO ASSESS TUTORING CAPABILITIES OF LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

## ABSTRACT

As students increasingly adopt large language models (LLMs) as learning aids, it is crucial to build models that are adept at handling the nuances of tutoring: they need to identify the core needs of students, be adaptive, provide personalized guidance, and be accurate. To this end, we introduce TUTORBENCH, a dataset and evaluation benchmark designed to rigorously evaluate the *core tutoring skills* of LLMs. The dataset comprises 1,490 samples curated by human experts, focused on high-school and AP-level curricula. The samples are drawn from three common tutoring tasks: (i) generating adaptive explanations tailored to a student’s confusion, (ii) providing actionable feedback on a student’s work, and (iii) promoting active learning through effective hint generation. To account for the inherent complexity of tutoring, samples are accompanied by sample-specific rubrics which are used to judge model responses during evaluation. TUTORBENCH uses a reliable and fine-grained automatic evaluation method that uses an LLM-judge and the sample-specific rubrics. We evaluate 16 frontier LLMs on TUTORBENCH and present a detailed analysis of their performance and behavior. Our results show that none of the frontier LLMs achieve a score of greater than 56%, showing a large room for improvement. We find that LLMs fall short in exhibiting the full range of tutoring skills needed to guide, diagnose, and support students effectively, with all the frontier models achieving less than a 60% pass rate on rubric criteria related to these skills. We also find that different model families exhibit varied strengths and limitations: the Claude models outperform others in supporting active learning, while they lag behind in the other two use cases. By releasing TUTORBENCH, we provide a comprehensive and unsaturated benchmark to guide the development of the next-generation of AI tutors.

## 1 INTRODUCTION

Large language models (LLMs) are rapidly transforming the way students learn and access educational support (Handa et al., 2025; Ammari et al., 2025; Geraghty & Goldstein, 2024; Scarlatos et al., 2025). Tools like ChatGPT, Gemini, and Claude have already serve as on-demand tutors for millions of learners worldwide. To encourage adoption, many providers offer variants specializing in tutoring (OpenAI, 2025b) and extend free access to students (Google, 2025). It is clear tutoring stands as a priority for model builders because it holds the promise of delivering transformative, personalized education at low cost. However, the true impact of AI-based tutoring is not yet fully understood: while many studies show improvements in learning outcomes (Vanzo et al., 2025), there are also reports that students perform worse when the tools are taken away (Bastani et al., 2025).

Given their rapid adoption, it is crucial to study and benchmark LLM behavior when they serve as tutors. Many existing LLM benchmarks focus on evaluating LLMs on advanced domain knowledge and reasoning (Phan et al., 2025; Glazer et al., 2024; He et al., 2024). While these dimensions are crucial for understanding the limits of model performance, they overlook the more nuanced, human-centered capabilities necessary for tutoring, such as clear and thorough explanation, adaptivity to a learner’s needs, and providing the right guidance. Some recent efforts have begun to explore how to benchmark LLMs in tutoring contexts (Maurya et al., 2025; Gupta et al., 2025). However, they are limited in scope as they focus on text-only examples from just one or two subjects like mathematics or computer science. This differs significantly from real-world usage, where students frequently upload pictures of their work and seek help on several subjects. SOTA LLMs such as GPT-4 achieve nearly perfect scores on both of these benchmarks according to their own reported results, indicating

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073



**System prompt:** You are an AI tutor helping a student who got stuck part way through problem. Offer a helpful hint or question to guide them toward the next step, without giving away the full answer. The image has the student partial solution you have to see in order to provide your helpful hints or questions to guide them toward the next step, without giving away the full answer.



**Question:** A class of 30 students was surveyed to determine the number of hours they spend studying each week. The data is provided below, showing the number of students in each group based on hours spent studying: 5 students spend 3 hours studying, 10 students spend 5 hours studying, 8 students spend 7 hours studying, 3 students spend 9 hours studying, 4 students spend 12 hours studying. Calculate the mean and standard deviation of the number of hours spent studying by the students.



Before you go further, try rewriting the numerator of the mean step-by-step:

1. List every (hours  $\times$  frequency) product for EACH of the five hour groups (3, 5, 7, 9, 12).
2. Add those five products together.
3. Divide that total by the full class size (30 students).

Does your current expression include the 9-hour and 12-hour groups? Once you have the correct mean, use it in  $\sigma^2 = \sum (x - \mu)^2 / N$  (remember to square the difference first, then multiply by the frequency, and finally sum over all groups).

$$\begin{aligned} \text{Mean} &= \frac{(3 \times 5) + (5 \times 10) + (7 \times 8) + (9 \times 3) + (12 \times 4)}{30} \\ &= \frac{15 + 50 + 56 + 27 + 48}{30} \\ &= \frac{196}{30} \approx 6.53 \end{aligned}$$

standard deviation :

I know that it relate to the mean I calculated but I don't remember the formula I thought it could be  $\frac{\sum (x - \mu)^2}{N}$ , but I'm not sure.

Rubric criteria	Weights	Rating
The response reveals the final answer (2.74) directly	-5	✗
The response should use clear formatting or structure (e.g., headers) to organize explanation into sections.	1	✓
⋮		
The response must explain that taking the square root of the variance gives the standard deviation in original units.	5	✗

Figure 1: An example from the TUTORBENCH dataset. Each sample includes a **system prompt** defining the tutoring goal (top left), a **question** along with the student’s partial work (top right), an **AI model’s response** (bottom left), and a set of **rubric criteria** for evaluation (bottom right). In this instance, the AI is prompted to provide a hint for a statistics problem.

the need for more nuanced benchmarks that expose shortcomings in model behavior. Moreover, they lack reliable auto-judge methods to evaluate LLMs’ tutoring capabilities.

To address the above gaps, we present **TUTORBENCH**<sup>1</sup>, a benchmark to evaluate the tutoring capabilities of LLMs. TUTORBENCH consists of 1490 conversations between a student-persona and a tutor-persona. The conversations cover 6 STEM subjects (Biology, Physics, Chemistry, Statistics, Calculus, and Computer Science) and focus on high school and AP curricula. The conversations are designed to reflect real-world usage by students and focus on three tutoring use cases: adaptive explanation generation, feedback and assessment, and active learning support. These use cases were chosen to test an LLM’s ability to calibrate its responses to individual students and support their learning journey rather than simply generating universal and standard solutions to questions. We provide more details about the use cases in Section 2.1. TUTORBENCH is also **multimodal**, with both text-only and image-based conversations. 828 samples in the dataset contain images of handwritten or printed work by students, reflecting real-world usage by students.

We design TUTORBENCH to be a **rubrics-based** evaluation to account for the open-ended nature of tutoring. Each sample in the dataset is accompanied by a set of sample-specific evaluation rubrics written by human experts that are self-contained, mutually exclusive, and verifiable. Together, these rubrics capture the requirements of a desirable response to the corresponding test sample. Model responses are graded by an LLM judge with a pass/fail rating on each rubric criterion.

We conduct data collection of TUTORBENCH with the help of human experts. The questions and the rubrics to each question were written by experts of the corresponding subject who hold a Bachelor’s or higher degree and have either tutoring or professional experience in the corresponding subject. In order to guarantee the difficulty level of TUTORBENCH, we then prompt 5 state-of-the-art LLMs (Gemini 2.5 Pro Google DeepMind (2025), Claude 3.7 Sonnet Anthropic (2025a), o3 OpenAI (2025a), DeepSeek-R1 DeepSeek (2025) and Llama 4 Maverick Meta AI (2025)) to respond to each collected test sample. We retain only conversations that result in a score of less than 50% for at least three of the five models. This results in a challenging benchmark, with the best performing model attaining a score of 55.65%.

<sup>1</sup>A sample subset of TUTORBENCH can be found at <https://huggingface.co/datasets/tutorbench/tutorbench>. The full dataset will be released soon.

108 Using the samples in TUTORBENCH, we evaluate and analyze 16 frontier LLMs. We report the  
 109 overall scores of the models in Section 3.1. Among all the models tested, Gemini 2.5 Pro and GPT-5  
 110 achieve the best performance, with 55.65% and 55.33% final scores respectively (100% being the  
 111 maximum final score). The Claude model series falls behind the Gemini 2.5 Pro and GPT model  
 112 series with non-negligible performance gaps, with Claude Opus 4.1 (Thinking) achieving a final  
 113 score of 50.78%. We observe that models achieve a score of only 47.16% on the use case of *adaptive*  
 114 *explanation generation*, indicating that frontier LLMs still struggle to generate effective personalized  
 115 responses. Further, they also achieve less than 55% on the other two use cases. Interestingly, models  
 116 of the Claude series outperform others in the “*active learning support*” use case, while trailing in  
 117 the overall performance. More detailed analysis and observations with respect to the use cases and  
 118 evaluation dimensions can be found in Section 3.2 and Section 3.3.

119

## 120 2 BENCHMARK DESIGN

121

122 We design TUTORBENCH to cover a broad spectrum of common, real-world tutoring use cases,  
 123 including both text-only and multimodal questions on 6 STEM subjects. Use cases in TUTOR-  
 124 BENCH focus on assessing human-centered capabilities necessary for tutoring, such as explanation,  
 125 guidance, and adaptivity to a learner’s needs. We also design a reliable rubric-based eval to auto-  
 126 matically assess LLM capabilities on such subjective matters. Details regarding the data collection  
 127 process are provided in the Appendix (Section A.1).

128

### 129 2.1 TUTORING USE CASES

130

131 As large language models (LLMs) become increasingly integrated into educational workflows, it is  
 132 important to rigorously evaluate their effectiveness in different tutoring scenarios. TUTORBENCH  
 133 captures three core use cases that reflect essential tutoring behaviors: adaptive explanation, assess-  
 134 ment and feedback, and active learning support. We elaborate on these 3 use cases below.

135

136 **Use Case 1: Adaptive Explanation Generation.** Personalized instructions can help unlock a  
 137 student’s potential. One of the most impactful qualities of a good human tutor is their ability to  
 138 adapt explanations based on a student’s current understanding and knowledge gaps. The key skill  
 139 required here is to identify core misconceptions and shortcomings in a student’s understanding and  
 140 address them in a way that is easy to understand.

141 To evaluate this capability, we design a multi-turn interaction setup consisting of: an initial question  
 142 written in a *student-persona*, an initial explanation written in a *tutor-persona*, and a follow-up ques-  
 143 tion written in a student-persona. During evaluation, an LLM is presented with this triplet, along  
 144 with a system prompt/instruction to generate an explanation that directly addresses the knowledge  
 145 gap or misunderstanding exhibited in the follow-up question by the student. An example conversa-  
 146 tion is shown in Fig. 2. Model responses are evaluated on whether they can recognize the specific  
 147 context and background implied by the student’s follow-up, and whether they can produce a helpful,  
 148 focused, and to-the-point response tailored to the student’s current understanding. These aspects are  
 149 captured in the set of sample-specific rubric criteria that accompany every sample in the dataset.

150

151 **Use Case 2: Assessment and Feedback.** Students (or supervisors) often use LLMs as a means  
 152 to assess their (or their students’) work and get instant feedback. This reflects one of the most  
 153 promising applications of LLMs: providing real-time evaluation and correction that helps learners  
 154 iterate quickly. To capture this scenario, TUTORBENCH includes examples where the model is  
 155 shown a student solution and is asked to assess it. The model is expected to analyze the work,  
 156 identify any mistakes, classify the nature of each error (e.g., arithmetic, factual, conceptual, etc),  
 157 and generate feedback. An example of this use case is shown in Fig. 2.

158

159 **Use Case 3: Active Learning Support.** An essential aspect of high-quality tutoring is to pro-  
 160 mote active learning, wherein students are encouraged to engage directly with the problem-solving  
 161 process rather than passively receiving answers. Effective tutors guide students toward the solution  
 162 through scaffolding, such as hints, analogies, or intermediate steps, which support learning while  
 163 preserving student agency.

162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181

Consider the galvanic cell with the following half-reactions and standard reduction potentials:

$\text{Ag}^+ (\text{aq}) + \text{e}^- \rightarrow \text{Ag} (\text{s}) E^\circ = +0.80 \text{ V}$   
 $\text{Cu}^{2+} (\text{aq}) + 2\text{e}^- \rightarrow \text{Cu} (\text{s}) E^\circ = +0.34 \text{ V}$

The half-cells are connected, and a spontaneous reaction is observed. What is the overall cell potential ( $E^\circ_{\text{cell}}$ )?

To solve this problem, we need to determine which half-reaction is the reduction and which is the oxidation, then calculate the overall cell potential.

Step 1: Identify which reaction occurs at each electrode ... [full text not shown here for brevity]

Ok, so I understand that the answer is right, but I'm a little confused about why copper is reduced, and why we don't multiply the silver reaction by 2, so we can balance the 2 electrons in the copper reaction.

Which of the following questions is most relevant to understanding the Calvin cycle?

A) How does chlorophyll capture light?  
B) How is ATP used in the formation of 3-carbon carbohydrates?  
C) How is  $\text{NADP}^+$  reduced to  $\text{NADPH}$ ?  
D) How is ATP produced in chemiosmosis?

Hmm, let me think about this... the Calvin cycle is part of photosynthesis, I remember that much.

I think it's A) How does chlorophyll capture light?

The Calvin cycle is where plants make sugar, right? I remember my teacher saying it uses the ATP to build glucose molecules. I think the Calvin cycle uses ATP to form those 3-carbon compounds... I think they're called G3P or something?

Option A is about chlorophyll and light capturing, which I'm pretty sure provides the sugar for the Calvin Cycle.

Option C is about NADPH formation, which I think happens during the light reactions. The Calvin cycle uses NADPH, but doesn't create it.

Option D is about ATP production through chemiosmosis, which is definitely part of the light reactions where the proton gradient makes ATP.

So yeah, I'm going with A because the Calvin cycle is where the plant uses ATP to build sugar molecules, and it requires chlorophyll to fuel it.

Adaptive explanation generation

Assessment and feedback

Figure 2: Examples of two core use cases for an AI tutor. **Left:** The *adaptive explanation generation* scenario, showcasing a multi-turn dialogue in electrochemistry. The system must provide a targeted clarification in response to a student’s specific follow-up question regarding standard reduction potentials. **Right:** The *assessment and feedback* scenario, where a student provides a reasoned but incorrect answer to a biology question. The system’s task is to analyze the student’s reasoning, identify the misconception about the Calvin cycle, and provide corrective feedback.

To test this behavior, we present LLMs with incorrect or partially correct student responses and ask them to generate helpful hints that enable the student to take the next step. These hints must strike a careful balance, providing just enough information to resolve confusion without directly giving away the final answer. This task simulates scenarios where a student is stuck midway and requires targeted intervention to proceed. An example of this use case is presented in Fig. 1.

## 2.2 MULTIMODALITY

One key feature of TUTORBENCH is its multimodal design, reflecting the authentic ways students engage with LLMs for tutoring. For each use case illustrated above, we create both text-only and multimodal samples. The images in the dataset consist of hand-written text/diagrams, typed content, or screenshots, mirroring real settings where students often communicate not only through written text but also by sharing images of their work or problem-solving steps. This multimodal component is essential for evaluating LLMs as effective AI tutors, since it tests their ability to interpret and provide feedback across different forms of student input. An example of a multimodal sample is shown in Fig. 1 where the student solution is in the form of an image of their hand-written work.

## 2.3 LLM JUDGE WITH SAMPLE-SPECIFIC RUBRICS

Given the open-ended nature of the three use cases described above, we adopt a rubric-based evaluation using an LLM-judge based on Claude Sonnet 4 (Anthropic, 2025b) to evaluate model responses. Similar to recent research on rubrics-based model training and evaluation Arora et al. (2025), TUTORBENCH creates separate rubric criteria for each example (ranging from 3 to 39 per example) that help capture fine-grained qualities of a given response. This results in a total of 15,220 rubrics across the whole dataset. These criteria enable efficient and reliable evaluation using an LLM-judge, bypassing a time-consuming human judge for every new model release.

Rubric criteria are designed to be self-contained, mutually exclusive, and collectively comprehensive. This design ensures that they can be applied consistently and without overlap, while also covering the full space of desirable response qualities. Granularity of the criteria makes them partic-

216 ularly well-suited for grading with LLM-based judges. An illustration of the sample-specific rubrics  
 217 is shown in Fig. 1 (bottom right). More examples can be found in the sample dataset, a URL to  
 218 which is provided above.

219 We additionally introduce a weighting scheme to reflect the relative importance of different criteria.  
 220 Rubrics designated as *critical* are assigned a weight of 5, while *non-critical* rubrics carry a weight  
 221 of 1. In some cases, *critical* rubrics may be assigned a negative weight of -5 to penalize undesirable  
 222 behaviors. For instance, in active learning settings, a rubric may require that the final answer not  
 223 be revealed directly, and a violation of this principle would result in a strong negative score. To  
 224 quantify overall model performance, we compute a weighted average of the binary scores assigned  
 225 to each criterion. The final score is normalized to the range [0,1].

226 We conduct thorough experiments to show the alignment of the LLM-judge vs the human-expert  
 227 judge. To study the quality of the LLM-judge and compare it with human experts, we collect 3  
 228 human ratings per rubric criterion on a subset of 250 samples. Our experiments show that the LLM-  
 229 judge aligns better with the human experts than the median human expert, and achieves an F1 score  
 230 of 0.81 with respect to the majority vote. Details of experiments are described in Section 3.7.

#### 232 2.4 RUBRIC CRITERIA TAGS

234 To further enable fine-grained analysis of model performance, we annotate each rubric criterion with  
 235 four tags. Firstly, we annotate each rubric with one of the following *evaluation dimensions* to capture  
 236 the primary axis along which the model is being judged: instruction following, style and tone, truth-  
 237 fulness, visual reasoning, visual perception, conciseness and relevance, student level calibration,  
 238 and emotional component. These dimensions provide a structured lens for understanding different  
 239 facets of tutoring quality. We provide definitions of these dimensions in Appendix A.3. Secondly,  
 240 we tag the rubric criteria with specific *tutoring skills* being assessed: asking guiding questions, iden-  
 241 tifying core misconceptions, recognizing correct or incorrect student steps, including examples or  
 242 analogies, providing alternative solutions, stating definitions or theorems (or other knowledge), or  
 243 providing step-by-step help. Together, these two tags allow us to analyze not just whether the model  
 244 performs well, but also what type of tutoring capability is being demonstrated or lacking.

245 In addition, we annotate each rubric criterion with two complementary tags that enrich the dataset.  
 246 We annotate each rubric as being *explicit* or *implicit*, to indicate whether the criterion is related to  
 247 an explicit request in the prompt or implicitly inferred from the tutoring context. We also annotate  
 248 criteria as *objective* or *subjective* to denote whether the criterion judges a response according to  
 249 an objective standard (e.g., factual accuracy, correctness of a step) or a more subjective one (e.g.,  
 250 appropriateness of tone, empathy). We provide the distributions of the various tags in Appendix A.5.

### 252 3 EVALUATION AND ANALYSIS

#### 254 3.1 OVERALL MODEL PERFORMANCE

256 Model performance on TUTORBENCH is evaluated using the weighted average rubric rating  $\text{ARR}_w$ ,  
 257 defined on the  $j$ -th example as

$$259 \text{ARR}_w^j = \frac{\sum_{i=1}^{N_j} w_i^j \cdot \mathbb{1}_{r_i^j}}{\sum_{i=1}^{N_j} w_i^j \cdot \mathbb{1}_{w_i^j > 0}}, \quad (1)$$

262 where  $N_j$  is the number of criteria in the  $j$ -th example,  $w_i^j \in \{-5, 1, 5\}$  is the weight of the  $i$ -th  
 263 rubric for the  $j$ -th sample, and  $r_i^j \in \{0, 1\}$  is the fail/pass rating of the model response on the  $i$ -th  
 264 rubric of the  $j$ -th sample. The final score for a model is the average  $\text{ARR}_w$  across all the examples.

266 The evaluation results are presented in Table 1. We observe that Gemini 2.5 Pro (Google Deep-  
 267 Mind, 2025) and GPT-5 (GPT-5, 2025) obtain the best overall performance, with very similar scores  
 268 on both text-only and multimodal tests. Furthermore, none of the models surpass 56% overall perfor-  
 269 mance, highlighting the complex nature of TUTORBENCH. Finally, it is noteworthy that the recently  
 released gpt-oss-120b model performs close to the best models on the text-only subset.

Rank	Model	Text-Only (%)	Multimodal (%)	Overall (%)	CI
1	Gemini 2.5 Pro	<b>57.05</b>	<b>54.53</b>	<b>55.65</b>	$\pm 1.11$
2	GPT-5	<b>57.03</b>	53.97	55.33	$\pm 1.02$
3	o3 Pro	56.07	53.45	54.62	$\pm 1.02$
4	o3 Medium Effort	54.11	51.68	52.76	$\pm 1.00$
5	o3 High Effort	52.91	51.43	52.09	$\pm 1.01$
6	Claude Opus 4.1 (Thinking)	51.65	50.08	50.78	$\pm 1.05$
7	Claude Opus 4 (Thinking)	50.40	49.14	49.71	$\pm 1.02$
8	Claude Opus 4.1	49.51	45.72	47.40	$\pm 1.06$
9	Claude 3.7 Sonnet (Thinking)	45.67	47.07	46.45	$\pm 1.03$
10	Claude Opus 4	47.79	43.59	45.46	$\pm 1.06$
11	Llama 4 Maverick	39.54	40.73	40.20	$\pm 1.00$
12	GPT-4o	39.10	33.74	36.12	$\pm 0.96$
	gpt-oss-120b	56.01	N/A	N/A	$\pm 1.49$
	gpt-oss-20b	49.01	N/A	N/A	$\pm 1.53$
	DeepSeek-R1	48.38	N/A	N/A	$\pm 1.50$

Table 1: Evaluation of state-of-the-art LLMs on TUTORBENCH. Gemini 2.5 Pro achieves the highest score (55.65%), highlighting the complex nature of the samples in TUTORBENCH and the need for further advancement in models for tutoring applications. We report scores on the text-only samples for gpt-oss-120b, gpt-oss-20b, and DeepSeek-R1 because they do not support multimodal input.

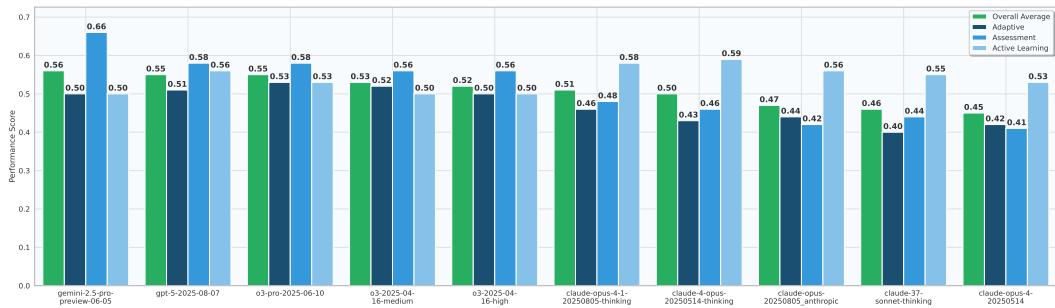


Figure 3: Model performance across three use cases: **Adaptive**, **Assessment**, and **Active Learning**. We observe a distinct difference between the performance of the Claude family of models compared to the other models, with the Claude models performing significantly better in providing active learning support, but still lagging behind other models overall.

### 3.2 PERFORMANCE BY USE CASE

Performance breakdown by use case is shown in Fig. 3. On average, models achieve a score of 47.16% on adaptive explanation generation, 51.56% on assessment and feedback, and 54.07% on active learning support. We also observe that the Claude family of models performs much better in the active learning support use case compared to the other models, although their overall performance is poorer.

### 3.3 PERFORMANCE ALONG EVALUATION DIMENSIONS

As explained in Section 2.4, we annotate each rubric criterion with one **evaluation dimensions**. This allows us to aggregate the pass/fail ratings on subsets with each of the tags. We present the mean pass/fail rate along these dimensions for the top 10 models (according to Table 1) in Fig. 4. While performance along these dimensions roughly follows the same trend as the overall performance, Gemini 2.5 Pro shows a much better performance than other models in recognizing student emotions such as confusion, frustration, curiosity, and generates responses with the right tone and style (e.g., by using headings, bullets, and LaTeX). On the other hand, GPT-5 and o3 Pro perform best on factual correctness (truthfulness), student-level calibration, and instruction following.

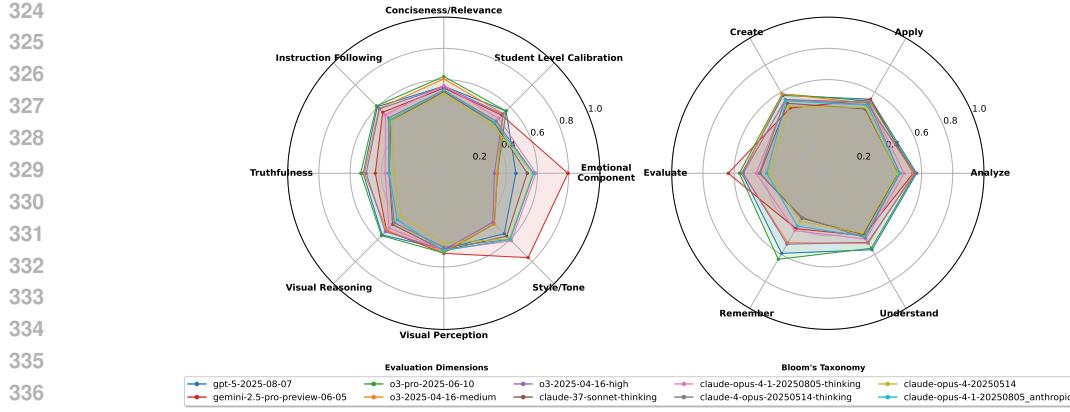


Figure 4: Model performance breakdown along evaluation dimensions and Bloom’s taxonomy categories. While the top-performing models GPT-5 and Gemini 2.5 Pro are close overall, their performance differs widely when measured along the above dimensions.

### 3.4 PERFORMANCE ALONG DIFFICULTY LEVEL

This study categorizes each sample using Bloom’s taxonomy, which organizes cognitive tasks into skill levels ranging from remembering and understanding to advanced ones such as analyzing, evaluating, and creating, in order to better assess model performance across different levels of cognitive demand. Three state-of-the-art LLMs (Gemini 2.5 Pro, Claude Sonnet 4, and GPT-5) were used to annotate each sample, with a category label assigned if at least two models agreed, covering over 97% of the samples. Final scores for each category were calculated by averaging  $ARR_w$  across all annotated samples. Results, shown in Figure 4, reveal that performance does not align with the taxonomy’s order of difficulty. For example, Gemini 2.5 Pro scores lower on “remember” than on “evaluate”, suggesting that while models may demonstrate advanced reasoning, they can struggle with recalling and explicitly stating content (when not explicitly prompted to), which is crucial for tutoring tasks that require effective context-driven communication.

### 3.5 PERFORMANCE ALONG TUTORING SKILLS

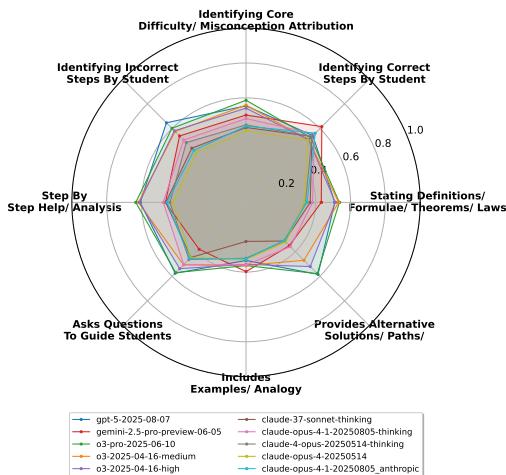
Tutoring is a complex task requiring several nuanced skills, such as calibrating to the student, recalling and applying knowledge, or breaking down problems into simpler steps. To evaluate models along with such finer skills, we classify the rubric criteria written by human-experts into 8 high-level skills: identifying core difficulty, identifying students’ correct steps, identifying students’ incorrect steps, recalling and stating knowledge, providing alternative solutions, including examples, asking questions to guide students, and providing step-by-step help. These skills were identified by manually inspecting 600 rubric criteria and by using an LLM to tag all the rubric criteria with one of the skills, if applicable.

We present the average pass/fail rate of LLM responses on these tags in Figure 5. Overall, we observe that models perform their best in identifying correct/incorrect steps by students (average scores of 53.7% and 53%), but they struggle to include alternate solutions, and examples/analogy in their responses (achieving average scores of 41.9% and 32.8%). Model performance varies across model families, with Gemini 2.5 Pro excelling at identifying correct/ incorrect steps and providing examples, while GPT-5 is stronger at spotting errors and misconceptions.

### 3.6 AN EVALUATION OF GPT-5 STUDY MODE

Along with the models in Table 1, we also evaluated OpenAI’s recently released *study mode* OpenAI (2025b). According to their website, study mode is “powered by custom system instructions written in collaboration with teachers, scientists, and pedagogy experts” to promote deeper learning. Since it is only available via the web interface, we collected responses from the website.

378 Unlike API models, study mode does not allow direct system prompts and instead engages users in  
 379 back-and-forth dialogue calibrated to their skill level. This makes apples-to-apples comparison with  
 380 TUTORBENCH difficult, since TUTORBENCH evaluates only the final response in a conversation  
 381 with a pre-specified history. For this reason, we exclude study mode from the main leaderboard and  
 382 report results separately.



400 Figure 5: Model performance breakdown along  
 401 tutoring skills: models struggle to include alter-  
 402 native solutions, examples, and analogies in their  
 403 responses. However, they perform relatively bet-  
 404 ter in identifying mistakes, correct steps, and core  
 405 misconceptions.

406 LLM-judge to grade model responses against the rubric criteria. We experimented with multiple  
 407 state-of-the-art models and chose Claude-4-Sonnet as our judge model, as it achieved the best align-  
 408 ment with human ratings.

409 To compare the LLM-judge with a  
 410 human expert, we first collect 3 sepa-  
 411 rate human-expert ratings across 250  
 412 examples, with a total of 2475 rubric  
 413 criteria. To capture the inter-human  
 414 ratings agreement, we select all the  
 415 rubric criteria rated by each human-  
 416 expert, and compare them to the  
 417 other two ratings. The mean inter-  
 418 human agreement across the 69 ex-  
 419 perts was 0.75, while agreement be-  
 420 tween the LLM-judge and human rat-  
 421 ings was 0.78. The distribution of  
 422 agreements per individual human ex-  
 423 pert who contributed to our annota-  
 424 tions is shown in Fig. 6, along with  
 425 the LLM-judge agreement.

426 We also study the alignment of the  
 427 LLM-judge with the majority vote on  
 428 rubrics. Firstly, we filter out the *criti-  
 429 cal* rubric criteria (with weights +5 or  
 430 -5) from the 250 examples, resulting  
 431 in a total of 1900 criteria with 3 rat-  
 ings each. We then use the majority vote (pass/fail) as the label for each rubric. The LLM judge

On TUTORBENCH, study mode scores  $46.94 \pm 1.06\%$  overall (48.19% text-only, 45.93% multimodal), indicating poor performance against rubric criteria. Comparing with API-based GPT-5, we note two key observations: (1) study mode often provides only partial, intermediate responses, typically ending with a counter-question (e.g., “Before we move on, do you want me to walk you through ...?”). These incomplete answers naturally score lower. (2) Even when complete, study mode responses perform worse than GPT-5’s, as they are more concise and address fewer rubric criteria.

Thus, study mode’s lower scores stem from both TUTORBENCH’s unsuitability for partial responses and from study mode’s tendency to neglect several rubric criteria even in complete answers. We provide examples and more details in the appendix (Section A.4).

### 3.7 QUALITY OF LLM-JUDGE USED IN AUTOMATED EVALUATIONS

To achieve scalable evaluation, we rely on an LLM-judge to grade model responses against the rubric criteria. We experimented with multiple state-of-the-art models and chose Claude-4-Sonnet as our judge model, as it achieved the best alignment with human ratings.

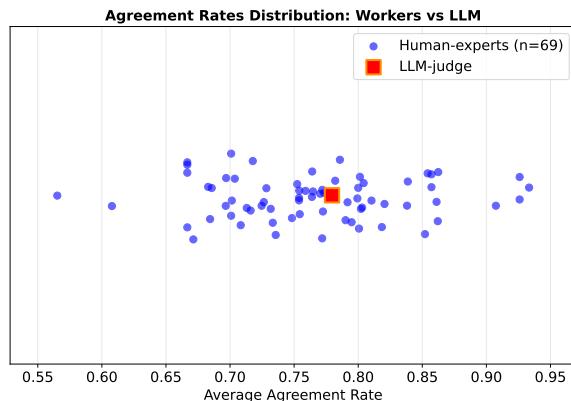


Figure 6: We measure the quality of the LLM-judge used in our evaluations by comparing its rating on model responses across 250 samples from TUTORBENCH with the majority vote obtained using 3 human-expert ratings. We observe that the LLM-judge ranks better than the median human-expert, thus demonstrating a strong alignment of the LLM-judge with human ratings.

432 achieves an F1-score of 0.82. The highest F1-score achieved by a single human-expert was 0.91,  
 433 computed across a set of 321 criteria that were rated by them. This demonstrates that the LLM-judge  
 434 has a strong performance and aligns well with human-expert ratings.  
 435

## 436 437 4 RELATED WORK

438  
 439 **Evaluation Frameworks for AI-Tutor Models.** Evaluating tutoring quality has emerged as an  
 440 important area of study. Maurya et al. (2025) propose MRBench, an evaluation benchmark with  
 441 192 math tutoring conversations. They propose a uniform taxonomy for evaluating tutoring capa-  
 442 bilities based on learning principles. In a similar effort to evaluate the impact of using LLM-based  
 443 tutors, Lyu et al. (2024) conduct a field study of 50 students by using an LLM-based coding as-  
 444 sistant. While the students who used the assistant achieved statistically significant improvement in  
 445 their performance, they expressed concerns about the limited role the assistant played in developing  
 446 critical thinking skills. Similarly, Liu et al. (2024) document the integration of AI-based tools into  
 447 Harvard’s CS50 course.  
 448

449 **LLMs as Tutors.** Recent research has also explored designing improved tutors using LLMs.  
 450 Pal Chowdhury et al. (2024) develop a structured, rule-based tutoring framework with LLMs. Their  
 451 system embeds guardrails and predefined pedagogical strategies, allowing LLMs to generate tutor  
 452 responses while maintaining instructional control. Compared to free-form GPT-4 tutors, this hybrid  
 453 approach reduced redundant or unhelpful dialog turns and better adhered to pedagogical principles.  
 454 Similarly, Wang et al. (2024) propose Bridge, which uses cognitive task analysis to capture expert  
 455 teachers’ decision-making processes, such as diagnosing student errors and selecting remediation  
 456 strategies, before feeding them into LLMs.  
 457

458 In contrast to these efforts, our dataset targets multimodal tutoring across six high school STEM  
 459 subjects, moving beyond single-domain math word problems. Our work also introduces explicit  
 460 rubric-based evaluation criteria tied to each sample, enabling efficient and scalable evaluation of  
 461 multiple models. Our dataset is also substantially larger and covers a diverse set of use cases, topics,  
 462 and evaluation dimensions. Overall, TUTORBENCH provides a broader testbed for measuring and  
 463 advancing LLMs’ tutoring capabilities.  
 464

## 465 5 LIMITATIONS OF TUTORBENCH

466 While our dataset provides a structured foundation for evaluating LLM tutoring, it has important lim-  
 467 itations. Its scope of instructional use cases is constrained: we focus on three representative tutoring  
 468 scenarios but omit other valuable tasks such as generating practice problems, designing exercises,  
 469 or introducing new concepts. The dataset also evaluates only final responses to pre-formulated con-  
 470 versations, limiting assessment of adaptability in dynamic, multi-turn exchanges. Our dataset only  
 471 incorporates images of students’ work as input and does not test visual content generation. Our  
 472 STEM focus provides depth in reasoning and problem-solving but excludes humanities domains  
 473 that rely on narrative and interpretive skills. Finally, we isolate model outputs from the broader  
 474 context of UI/UX design.  
 475

## 476 6 CONCLUSION

477 We introduce TUTORBENCH, a dataset and evaluation benchmark to assess tutoring capabilities of  
 478 LLMs. We focus on three specific user scenarios: (i) adapting explanations to a student with a  
 479 given background, (ii) providing assessment of students’ work, and (iii) providing hints for active  
 480 learning. TUTORBENCH is a rubrics-based evaluation and each sample in the dataset is accompani-  
 481 ed by specific rubric criteria. These rubric criteria are weighted based on importance and are also tagged  
 482 with several attributes that allow for finer analysis. Another important aspect of TUTORBENCH is  
 483 the support for multimodal input: the dataset includes samples with images containing students’  
 484 work to reflect real-world usage. We also provide a comprehensive evaluation of a select set of  
 485 state-of-the-art LLMs and identify growth areas. Overall, TUTORBENCH serves as one of the first  
 486 datasets to comprehensively evaluate LLMs on tutoring-based applications.  
 487

486     **Reproducibility Statement** To facilitate reproducibility of our results, we provide a sample ver-  
 487     sion of our dataset, which is publicly available at <https://huggingface.co/datasets/tutorbench/tutorbench>. The sample dataset consists of 30 samples, 10 from each use-case  
 488     (wit 5 text-only and 5 multimodal). We will soon release the full dataset along with the evaluation  
 489     code to enable exact reproduction of our experiments and to support future research on improving  
 490     the tutoring capabilities of large language models.  
 491

492     **LLM usage acknowledgement:** We acknowledge that large language models (LLMs) were used  
 493     to assist in the preparation of this manuscript, specifically for improving clarity and polishing the  
 494     language of certain sentences. All substantive ideas, analyses, and conclusions remain the responsi-  
 495     bility of the authors.  
 496

## 497     REFERENCES

499     Tawfiq Ammari, Meilun Chen, S. M. Mehedi Zaman, and Kiran Garimella. How students (re-  
 500     ally) use chatgpt: Uncovering experiences among undergraduate students. *arXiv preprint*  
 501     arXiv:2505.24126, 2025. doi: 10.48550/arXiv.2505.24126.  
 502

503     Anthropic. Claude 3.7 sonnet and claude code, February 2025a. URL <https://www.anthropic.com/news/claude-3-7-sonnet>. Product announcement.  
 504

505     Anthropic. Claude sonnet 4: Hybrid reasoning model from anthropic. <https://www.anthropic.com/claude/sonnet>, May 2025b. Also detailed in the Claude 4 announce-  
 506     ment (May 22 2025) and Wikipedia entry (accessed 2025-09-06).  
 507

508     Rahul K Arora, Jason Wei, Rebecca Soskin Hicks, Preston Bowman, Joaquin Quiñonero-Candela,  
 509     Foivos Tsimpourlas, Michael Sharman, Meghan Shah, Andrea Vallone, Alex Beutel, et al. Health-  
 510     bench: Evaluating large language models towards improved human health. *arXiv preprint*  
 511     arXiv:2505.08775, 2025.  
 512

513     Hamsa Bastani, Osbert Bastani, Alp Sungu, Haosen Ge, Özge Kabakçı, and Rei Mariman.  
 514     Generative ai without guardrails can harm learning: Evidence from high school mathemat-  
 515     ics. *Proceedings of the National Academy of Sciences*, 122(26):e2422633122, 2025. doi:  
 516     10.1073/pnas.2422633122. URL <https://www.pnas.org/doi/abs/10.1073/pnas.2422633122>.  
 517

518     DeepSeek. Deepseek-r1 release, January 2025. URL <https://api-docs.deepseek.com/news/news250120>. Model announcement and licensing details.  
 519

520     Sean Geraghty and Nash Goldstein. 11th grader takes an ai tutoring deep dive: And a human tutor-  
 521     ing expert takes notes. Education Next (blog post), November 22 2024. URL <https://www.educationnext.org/11th-grader-takes-an-ai-tutoring-deep-dive/>.  
 522     Last updated November 22, 2024.  
 523

524     Elliot Glazer, Ege Erdil, Tamay Besiroglu, Diego Chicharro, Evan Chen, Alex Gunning, Caroline  
 525     Falkman Olsson, Jean-Stanislas Denain, Anson Ho, Emily de Oliveira Santos, Olli Järvinen,  
 526     Matthew Barnett, Robert Sandler, Matej Vrzala, Jaime Sevilla, Qiuyu Ren, Elizabeth Pratt, Lionel  
 527     Levine, Grant Barkley, Natalie Stewart, Bogdan Grechuk, Tetiana Grechuk, Shreepranav Varma  
 528     Enugandla, and Mark Wildon. Frontiernmath: A benchmark for evaluating advanced mathematical  
 529     reasoning in ai. *arXiv preprint arXiv:2411.04872*, November 2024. doi: 10.48550/arXiv.2411.  
 530     04872.  
 531

532     Google. Google gemini: Free pro plan for students. <https://gemini.google/students/>,  
 533     2025. Accessed September 25, 2025. Offer ends October 6, 2025.  
 534

535     Google DeepMind. Gemini 2.5: Our most intelligent ai model, March  
 536     2025. URL <https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/>. Announcing availability of  
 537     Gemini 2.5 Pro.  
 538

539     OpenAI GPT-5. GPT-5 System Card. Online; August 7, 2025, 2025. URL <https://openai.com/index/gpt-5-system-card/>. Accessed: 2025-08-27.

540 Adit Gupta, Jennifer Reddig, Tommaso Calo, Daniel Weitekamp, and Christopher J. MacLellan.  
 541 Beyond final answers: Evaluating large language models for math tutoring. *arXiv preprint*  
 542 *arXiv:2503.16460*, February 2025. doi: 10.48550/arXiv.2503.16460.

543  
 544 Kunal Handa, Drew Bent, Alex Tamkin, Miles McCain, Esin Durmus, Michael Stern,  
 545 Mike Schiraldi, Saffron Huang, Stuart Ritchie, Steven Syverud, Kamyia Jagadish, Mar-  
 546 garet Vo, Matt Bell, and Deep Ganguli. Anthropic education report: How uni-  
 547 versity students use claude, 2025. URL <https://www.anthropic.com/news/anthropic-education-report-how-university-students-use-claude>.

548  
 549 Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi  
 550 Hu, Xu Han, Yujie Huang, Yuxiang Zhang, Jie Liu, Lei Qi, Zhiyuan Liu, and Maosong  
 551 Sun. Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilin-  
 552 gual multimodal scientific problems. *arXiv preprint arXiv:2402.14008*, February 2024. doi:  
 553 10.48550/arXiv.2402.14008.

554  
 555 Rongxin Liu, Carter Zenke, Charlie Liu, Andrew Holmes, Patrick Thornton, and David J. Malan.  
 556 Teaching cs50 with ai: Leveraging generative artificial intelligence in computer science ed-  
 557 ucation. In *Proceedings of the 55th ACM Technical Symposium on Computer Science Edu-  
 558 cation V. 1*, SIGCSE 2024, pp. 750–756, New York, NY, USA, 2024. Association for Com-  
 559 puting Machinery. ISBN 9798400704239. doi: 10.1145/3626252.3630938. URL <https://doi.org/10.1145/3626252.3630938>.

560  
 561 Wenhao Lyu, Yimeng Wang, Tingting (Rachel) Chung, Yifan Sun, and Yixuan Zhang. Evaluating  
 562 the effectiveness of llms in introductory computer science education: A semester-long field study.  
 563 In *Proceedings of the Eleventh ACM Conference on Learning @ Scale, L@S '24*, pp. 63–74,  
 564 New York, NY, USA, 2024. Association for Computing Machinery. ISBN 9798400706332. doi:  
 565 10.1145/3657604.3662036. URL <https://doi.org/10.1145/3657604.3662036>.

566  
 567 Kaushal Kumar Maurya, Kv Aditya Srivatsa, Ksenia Petukhova, and Ekaterina Kochmar. Unifying  
 568 AI tutor evaluation: An evaluation taxonomy for pedagogical ability assessment of LLM-powered  
 569 AI tutors. In Luis Chiruzzo, Alan Ritter, and Lu Wang (eds.), *Proceedings of the 2025 Conference  
 570 of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human  
 571 Language Technologies (Volume 1: Long Papers)*, pp. 1234–1251, Albuquerque, New Mexico,  
 572 April 2025. Association for Computational Linguistics. doi: 10.18653/v1/2025.naacl-long.57.  
 573 URL <https://aclanthology.org/2025.naacl-long.57/>.

574  
 575 Meta AI. The llama 4 herd: The beginning of a new era of natively mul-  
 576 timodal intelligence, April 2025. URL <https://ai.meta.com/blog/llama-4-multimodal-intelligence/>. Announcing Llama 4 Scout and Llama 4  
 577 Maverick.

578  
 579 OpenAI. Introducing openai o3 and o4-mini, April 2025a. URL <https://openai.com/index/introducing-o3-and-o4-mini/>. Model announcement; later updates noted on  
 580 page.

581  
 582 OpenAI. Introducing study mode. <https://openai.com/index/chatgpt-study-mode/>, July 2025b. Accessed September 25, 2025. Feature launched July  
 583 29, 2025.

584  
 585 Sankalan Pal Chowdhury, Vilém Zouhar, and Mrinmaya Sachan. Autotutor meets large language  
 586 models: A language model tutor with rich pedagogy and guardrails. In *Proceedings of the  
 587 Eleventh ACM Conference on Learning @ Scale (L@S '24)*, pp. 5–15, Atlanta, GA, USA, 2024.  
 588 Association for Computing Machinery. doi: 10.1145/3657604.3662041.

589  
 590 Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li, Josephina Hu, Hugh Zhang, et al. Humanity’s  
 591 last exam. *arXiv preprint arXiv:2501.14249*, January 2025. doi: 10.48550/arXiv.2501.14249.

592  
 593 Alexander Scarlatos, Naiming Liu, Jaewook Lee, Richard Baraniuk, and Andrew Lan. Train-  
 594 ing llm-based tutors to improve student learning outcomes in dialogues. *arXiv preprint*  
*arXiv:2503.06424v1*, March 2025. doi: 10.48550/arXiv.2503.06424.

594 Alessandro Vanzo, Sankalan Pal Chowdhury, and Mrinmaya Sachan. GPT-4 as a homework tutor  
 595 can improve student engagement and learning outcomes. In Wanxiang Che, Joyce Nabende, Eka-  
 596 terina Shutova, and Mohammad Taher Pilehvar (eds.), *Proceedings of the 63rd Annual Meeting*  
 597 *of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 31119–31136,  
 598 Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-89176-  
 599 251-0. doi: 10.18653/v1/2025.acl-long.1502. URL <https://aclanthology.org/2025.acl-long.1502/>.  
 600

601  
 602 Rose Wang, Qingyang Zhang, Carly Robinson, Susanna Loeb, and Dorottya Demszky. Bridg-  
 603 ing the novice-expert gap via models of decision-making: A case study on remediating math  
 604 mistakes. In *Proceedings of the 2024 Conference of the North American Chapter of the As-  
 605 sociation for Computational Linguistics: Human Language Technologies (Volume 1: Long  
 606 Papers)*, Mexico City, Mexico, June 2024. Association for Computational Linguistics. URL  
 607 <https://aclanthology.org/2024.naacl-long.120/>.  
 608

## 609 A APPENDIX

### 610 A.1 DATA COLLECTION PROCESS

611 TUTORBENCH consists of high-school-level questions from STEM subjects. The questions and the  
 612 rubrics to each question were written by experts of the corresponding subject who have a Bachelor’s  
 613 or higher degree and have either tutoring or professional experience in the corresponding subject.  
 614 Each example starts with a question, which is then followed by use-case-specific content. For the  
 615 adaptive explanation use case, it is followed by (i) a teacher’s explanation of the answer and (ii) a  
 616 student’s follow-up question asking for clarification on a specific part of that explanation. For the  
 617 feedback and assessment use case, the initial question is followed by an incorrect solution written  
 618 in a student persona. For the active learning support use case, the initial question is followed by a  
 619 partial solution written in a student persona.  
 620

621 Subsequently, the human experts write a list of rubric criteria to evaluate model responses. The  
 622 rubric criteria are formulated based on a “golden tutoring response” written by them. Each rubric  
 623 criterion is annotated with several tags as explained in Section 2.4, and assigned a weight to indicate  
 624 their importance. The weighting system is discussed in Section 2.3.  
 625

626 Five state-of-the-art models, Gemini 2.5 Pro, GPT-o3, Claude 3.7 Sonnet , DeepSeek-R1, Llama 4  
 627 Maverick are then prompted with the full example, along with use-case-specific instructions. The  
 628 obtained responses are then graded against the rubric criteria. We keep only those samples for which  
 629 3 out of the 5 models achieve less than 50% score (weighted) across the rubric criteria.  
 630

### 631 A.2 DISTRIBUTION OF SAMPLES OVER SUBJECTS AND USE CASES

### 632 A.3 RUBRIC CRITERIA TAG DEFINITIONS

633 As explained in Section 2.4, we annotate each rubric criterion with several tags to enable fine-grained  
 634 analysis of model behavior. In Table 2, we provide definitions of the *evaluation dimensions* tag.  
 635

### 636 A.4 GPT-5 STUDY MODE RESPONSE EXAMPLES

637 In Section 3.6, we provided details about why the study mode responses from GPT-5 collected from  
 638 the website rank lower than GPT-5 itself. Here, we provide illustrative examples of cases where  
 639 either the study mode response is indeed poorer or is truncated due to the model seeking user input  
 640 before proceeding. In the latter case, the model response is incomplete (probably by design) and  
 641 hence scores poorly against the rubric criteria. We present three examples in Figs. 8 to 10 where  
 642 their responses are shown side-by-side along with the corresponding rubric criteria. These examples  
 643 highlight some of the shortcomings of study mode: the responses are sometimes terse, or fail to get  
 644 into the details of the problem, or sometimes expect a response from the student before proceeding.  
 645 By grounding the comparison in rubric-based criteria, we provide concrete evidence for why study  
 646 mode ranks lower in Table 1.  
 647

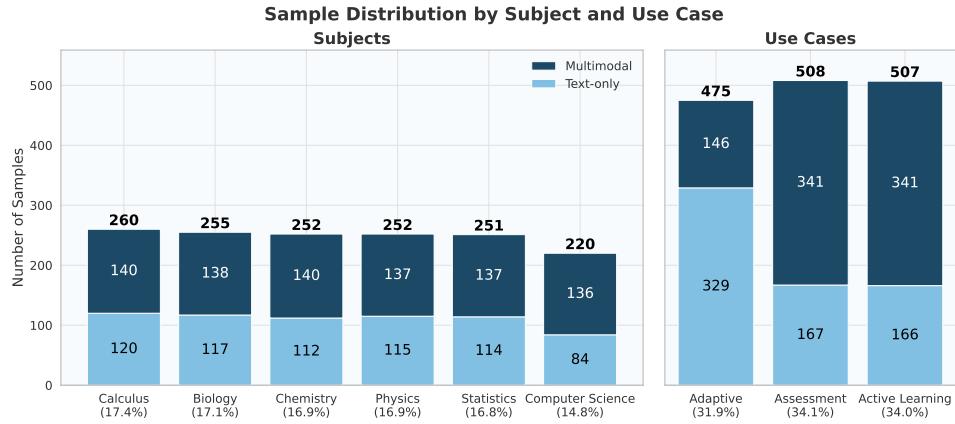


Figure 7: Distribution of samples by subject (left) and use case (right). The stacked bars indicate the number of Text-only (light blue) and Multimodal (dark blue) samples in each category. The distribution across the six subjects is relatively balanced, while the counts for the use cases vary more significantly. The labels ‘Adaptive’, ‘Assessment’, and ‘Active Learning’ correspond to adaptive explanation generation, assessment and feedback, and active learning support, respectively.

Tag	Definition
Instruction following	Criterion checks for adherence to system instructions and prompt instructions
Truthfulness	Criterion checks for factual accuracy, capturing the “precision” of the model, and does not penalize missed content, incorrect formatting
Conciseness and Relevance	Criterion checks for how direct, on-topic, and efficient the content is. Examples are criteria that check for what should not be in the response, unnecessary formatting
Style and Tone	Criterion checks for clarity, fluency, and appropriateness of tone
Visual Perception [Multimodal Only]	Criterion checks whether the model correctly identifies the content from the image required to solve the problem.
Visual Reasoning [Multimodal Only]	Criterion checks whether the model correctly reasons based on the provided image. If the reasoning does not depend on the image, this tag is not used.
Student-level Calibration	Criterion checks if the model response accounts for the explicit/ implicit knowledge level of the student.
Emotional Component Calibration	Criterion checks if the model pays attention to a student’s enthusiasm, frustrations, overwhelm, confusion, and other emotional cues mentioned in the prompt.

Table 2: Definitions of the evaluation dimensions used to tag rubric criteria.

### A.5 DISTRIBUTION OF RUBRIC CRITERIA TAGS

In this section, we present the distribution of the four tags that each rubric criterion was annotated with when applicable. In Fig. 11, we show the distribution of the objectivity and explicitness tags. We observe that the majority of the rubric criteria are both objective and explicit. In Fig. 12a, we show the distribution of the evaluation dimensions, and in Fig. 12b, we show the distribution of tutoring skills across the rubrics.

### A.6 SYSTEM PROMPTS

To generate model responses for evaluation, we use system prompts to lightly guide the LLMs to assume the role of a tutor. We refrain from using overly prescriptive system prompts to measure the “natural” tutoring abilities of LLMs. This helps us create an evaluation system that is reflective of

702		
703		
704		
705		
706		
707		
708		
709		
710		
711		
712		
713		
714		
715		
716		
717		
718		
719		
720		
721		
722		
723		
724		
725		
726		
727		
728		
729		
730		
731		
732		
733		
734		
735		
736		
737		
738		
739		
740		
741		
742		
743		
744		
745		
746		
747		
748		
749		
750		
751		
752		
753		
754		
755		

 (image)

Question: What is responsible for the seasonal change in the coat color of Arctic foxes?

In winter, Arctic foxes usually have a white coat. That's because lower temperatures in winter denature the pigment molecules in the fox's coat, causing it to be lighter in color. In summer, when the temperature rises again, the coat color changes back to brown.

Study mode  
(web)

GPT-5  
(API)

Think about how temperature affects pigment production in the fox's fur — what happens to those pigment molecules in the warmer months compared to the colder months? That contrast will help you explain the seasonal color change.

Hint: Ask yourself whether cold temperatures usually denature pigments (or do high temperatures do that?). Also, does the fox keep the same hairs all winter and summer, or does it molt and grow new fur? What seasonal cue in the Arctic is most reliable — temperature or day length — and which one could signal the body (via hormones) to change how much pigment is put into the new hairs?

	Study mode (web)	GPT-5 (API)
The response must acknowledge the student's current work by referring to their claim that lower temperatures in winter denature the pigment molecules in the fox's coat, causing it to be lighter in color.	0	0
The response must identify that the student is stuck due to a conceptual misunderstanding of what actually causes denaturation (i.e., high rather than low temperatures and extreme pH).	0	1
The response must hint that the seasonal coat-color changes observed in the Arctic fox cannot be caused by denaturation by encouraging the student to consider whether denaturation could occur at sub-zero arctic temperatures.	0	1
The response must hint that the seasonal coat-color changes observed in the Arctic fox cannot be caused by denaturation by encouraging the student to consider whether hair follicles whose pigment molecules are irreversibly denatured would ever be able to produce dark hairs again in the summer.	0	0
The response must provide a hint promoting the student to consider that denaturation is irreversible and therefore unlikely to explain the seasonal coat-color change.	0	0
The response must hint that a hormone affected by day length changes plays a role in regulating melanin production in Arctic foxes.	0	1
The response must hint that changes in day length can biologically influence fur color through melanin production.	0	1
The model must be able to reason that the student's response that the Arctic fox's pigment molecules are denatured in the winter is incorrect because denaturation cannot occur at the extreme low temperatures in which the Arctic fox lives.	0	1
The response should mention that the Arctic fox grows a new coat each season, which may help explain reversible changes in coat color.	0	1
The response must not state that the correct final answer is that the seasonal coat-color changes observed in the Arctic fox are due to photoperiodic melatonin regulation of melanin.	1	1

Figure 8: An example comparing study mode and GPT-5 on the use case of *active learning support*. We can observe that the rubrics require the response to correct a misconception in the student regarding the relationship between temperature and fur color in arctic foxes and guide the student to consider other factors (such as daylight duration and hormonal changes) affecting the fur color during winter. As observed, the response from GPT-5 (API) addresses the rubric criteria directly, while the response from study mode is brief and does not refer to any other factors at all.

real-world usage, where students do not use elaborate instructions. Below, we outline the system prompt used for each use case.

#### Adaptive explanation generation:

- **Text-only samples:** You are an AI tutor helping a high school student understand a concept. Answer their question clearly and adjust your explanation based on what the student says they're confused about.

756		
757		
758		
759		
760		
761		
762		
763		
764		
765		
766		
767		
768		
769		
770		
771		
772		
773		
774		
775		
776		
777		
778		
779		
780		
781		
782		
783		
784		
785		
786		
787		
788		
789		
790		
791		
792		
793		
794		
795		
796		
797		
798		
799		
800		
801		
802		
803		
804		
805		
806		
807		
808		
809		

**Study mode (web)**

A right prism has a base that is an isosceles trapezoid with legs that are \$5\$ inches in length and with one parallel side that is \$6\$ inches longer than the other parallel side. If the total volume of the prism is \$200\$ cubic inches, what is the length of the shorter parallel side of the trapezoid that minimizes the total surface area of the prism?

Initial solution (content omitted here for brevity)

I'm really confused. I don't understand the surface area calculation. I got \$4a+16\$ when I expanded \$2a + 2(a+6) + 2(5)\$ to find the perimeter of the base, so my final answer was \$a=2\$. What did I do wrong?

**GPT-5 (API)**

Great question—you're right to focus on the surface area setup. Here's where the mix-ups happened and how to fix them.

1) Which perimeter to use  
For a right prism, the total surface area is:  
SA = 2B + P·H  
- B = area of one base  
- P = perimeter of one base (not doubled)  
- H = prism height

So you should use P = a + (a+6) + 5 + 5 = 2a + 16.  
Writing 2a + 2(a+6) + 2(5) actually gives 2P, which double-counts the perimeter.

2) Small algebra slip  
Even if you were computing 2P, expanding 2(a+6) should give 2a + 12, not 2a + 6. So 2a + 2(a+6) + 2(5) = 2a + (2a + 12) + 10 = 4a + 22, not 4a + 16.

3) Correct setup and result  
- Height of trapezoid base: ht = 4 (from 5^2 = ht^2 + 3^2).  
- Base area B = 1/2(a + a + 6)4 = 4a + 12.  
- Prism height from volume 200 = B·H = 200/(4a + 12) = 50/(a + 3).  
- Perimeter: P = 2a + 16.

Surface area:  
SA = 2B + P·H = 2(4a + 12) + (2a + 16)50/(a + 3)  
= 8a + 24 + 1000/(a + 3)

Differentiate and set to zero:  
d(SA)/da = 8 - 500/(a + 3)^2 = 0  
=> (a + 3)^2 / 50 = 25  
=> a + 3 = sqrt(25) = 5  
=> a = 5 - 3 = 2 inches.

Quick check: with a = 2 (your answer), SA = 240; with a = 4.91, SA = 227, which is smaller. So the minimum occurs at a = 4.91 inches.

Does the SA = 2B + P·H formula and the perimeter P = 2a + 16 make sense now?

Figure 9: An example comparing study mode and GPT-5 on the use case of *adaptive explanation generation*. We can observe that the GPT-5 response clearly identifies that the student is confused about the formula of the perimeter and offers clarification, unlike study mode. Further, study mode does not provide the correct answer, but expects the student to respond before proceeding further.

- **Multimodal samples:** You are an AI tutor helping a high school student understand a concept. Answer their question clearly and adjust your explanation based on what the student says they’re confused about.

### Assessment and feedback:

- **Text-only samples:** You are an AI tutor reviewing a student’s answer to a question. Evaluate whether it is correct, identify any mistakes, and explain your reasoning clearly. Provide an assessment of the student incorrect solution in the first response
- **Multimodal samples:** You are an AI tutor reviewing a student’s answer to a question. Evaluate whether it is correct, identify any mistakes, and explain your reasoning clearly. Provide an assessment of the student incorrect solution present in the image.

810		
811		
812		
813		
814		
815		
816		
817		
818		
819		
820		
821		
822		
823		
824		
825		
826		
827		
828		
829		
830		
831		
832		
833		
834		
835		
836		
837		
838		
839		
840		
841		
842		
843		
844		
845		
846		
847		
848	Study mode (web)	GPT-5 (API)
849		
850		
851		
852		
853		
854		
855		
856		
857		
858		
859		
860		
861		
862		
863		

Figure 10: An example comparing study mode and GPT-5 on the use case of *assessment and feedback*. We can observe that study mode provides an inaccurate assessment and concludes that the student’s solution was **solid**. However, GPT-5 identifies the right source of error and provides suggestions on how to avoid such errors in the future.

Note that for text-only samples, the initial solution appears in the conversation history as a model response to the initial question. Hence, the system prompt refers to it as the ‘first response’.

### Active learning support:

- **Text-only samples:** You are an AI tutor helping a student who got stuck partway through a problem. Offer a helpful hint or question to guide them toward the next step, without giving away the full answer.
- **Multimodal samples:** You are an AI tutor helping a student who got stuck partway through a problem. Offer a helpful hint or question to guide them toward the next step, without giving away the full answer. The image has the student partial solution you have to see in

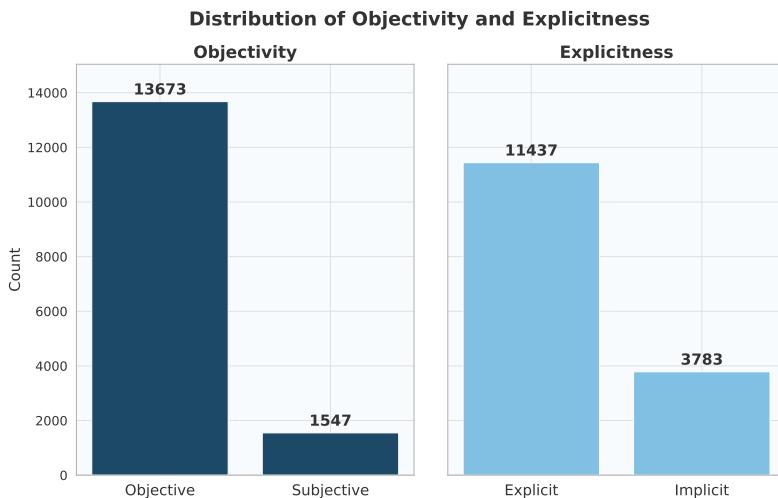


Figure 11: Distribution of objectivity and explicitness.

order to provide your helpful hints or questions to guide them toward the next step, without giving away the full answer

While manually collecting responses from the GPT-5 study model from its website, we use the following preambles:

#### Adaptive explanation generation:

- **Text-only samples** You will be presented with an initial question, a solution to the question, and a follow-up question seeking clarification. Help answer the follow-up question.
- **Multimodal samples:** You will be presented with an initial question, a solution to the question, and a follow-up question seeking clarification. The initial question will refer to the attached image. If an image is not provided, do not answer the question and say 'Image not provided'. Otherwise, answer the follow-up question.

#### Assessment and feedback:

- **Text-only samples:** You will be presented with an initial question and a solution to the question. Provide an assessment of the solution.
- **Multimodal samples:** You will be presented with an initial question and a solution to the question. Provide an assessment of the solution.

#### Active learning support:

- **Text-only samples:** Offer a helpful hint or question to guide me toward the next step.
- **Multimodal samples:** Offer a helpful hint or question to guide me toward the next step. My solution is shown in the attached image

918

919

920

921

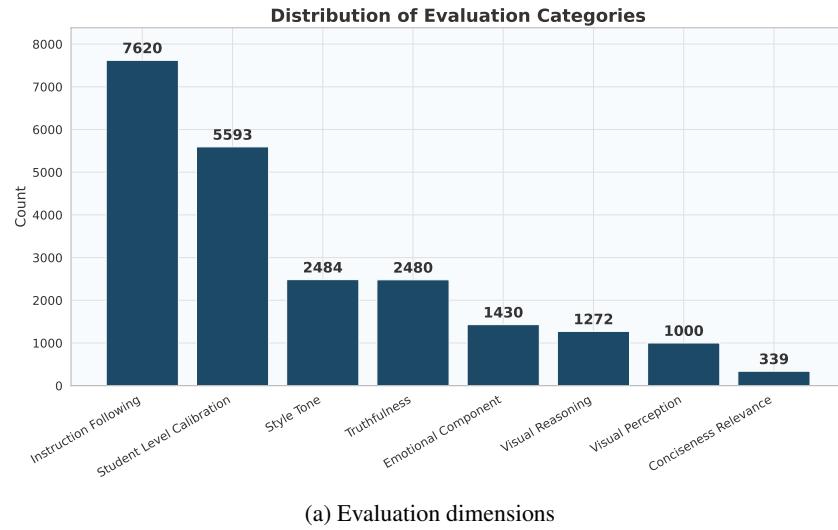
922

923

924

925

926



931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

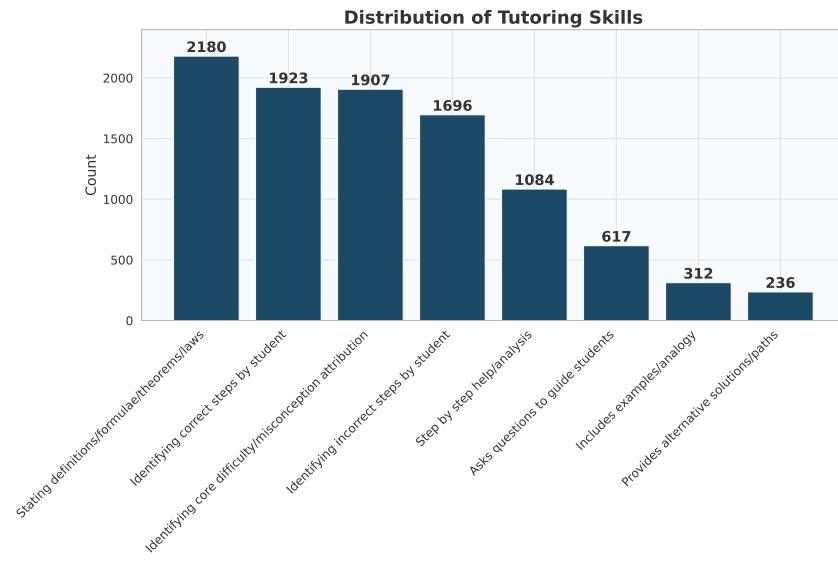


Figure 12: Distributions across (a) evaluation dimensions and (b) tutoring skills