

# Benchmarking Large Language Models for Diagnosing Students' Cognitive Skills from Handwritten Math Work

Anonymous ACL submission

## Abstract

Diagnosing students' cognitive skills from their handwritten math work is essential for personalized learning, as such work captures intermediate reasoning beyond final answers. Yet interpreting it remains labor-intensive for teachers, and automating it with LLMs is far from straightforward. Unlike mathematical problem-solving, where LLMs generate solutions themselves, cognitive diagnosis requires models to infer latent reasoning from student work. To systematically investigate this underexplored task, we construct MATHCOG, a benchmark dataset of 3,036 diagnostic verdicts across 639 student responses to 110 math problems, annotated by teachers using TIMSS-grounded cognitive skill checklists with evidential strength labels (*Evident/Vague*). Using MATHCOG, we evaluate 17 open and proprietary LLMs and find that (1) all models underperform ( $F1 < 0.521$ ) regardless of capability, and (2) performance degrades sharply under vague evidence. Error analysis reveals systematic failure patterns: models misattribute vague evidence as evident, over-infer from minimal cues, and hallucinate nonexistent evidence. These findings highlight fundamental limitations in LLMs' ability to reason under implicit evidential conditions, with direct implications for evidence-aware, teacher-in-the-loop designs in LLM-based cognitive diagnosis.

## 1 Introduction

Diagnosing students' cognitive skills from their problem-solving work is a long-standing goal in mathematics education, as it can reveal where and how students' reasoning breaks down beyond final correctness (Livingston, 2009; Jin et al., 2024). Cognitive skills refer to the mental procedures and knowledge structures that underlie competent performance (Anderson, 1982). In mathematics education, these skills are operationalized through frameworks such as TIMSS (Mullis, 2017), which

categorizes them into students' abilities to know, apply, and reason with mathematical concepts and procedures.

Handwritten math work provides rich insight into students' mathematical thinking, capturing intermediate reasoning steps, strategy choices, and partial understanding, which makes it particularly valuable for diagnosing cognitive skills. (Kheong, 1994; Jin et al., 2024). Despite its value, such handwritten work often provides incomplete, implicit, or unevenly expressed evidence of students' reasoning (Henderson et al., 2004; Rahbarnia et al., 2014), making reliable cognitive skill diagnosis inherently challenging. For example, students may omit some steps that could serve as clear evidence of certain skills (e.g., doing mental math and skipping writing out calculation steps); despite the absence of explicit evidence, one can infer that the students know how to calculate if they solved the later steps correctly. Despite recent advances in large language models (LLMs), including multimodal perception (Zhang et al., 2024b,a) and reasoning capabilities (Wang et al., 2023; Hao et al., 2024), their ability to diagnose and infer students' cognitive skills from handwritten work remains largely unexamined. While prior work has examined LLMs' own mathematical problem-solving abilities (Didolkar et al., 2024; Fang et al., 2024; Ahn et al., 2024), our work extends beyond problem solving to examine their ability to diagnose human problem-solving processes.

In this work, we systematically investigate how well existing LLMs diagnose students' cognitive skills in mathematics. Specifically, we focus on the degree of evidential strength in students' handwritten responses. We address the following research questions:

**RQ1.** How do different LLMs (varying in image input, reasoning, model size, and few-shot prompting) perform in diagnosing students'

083	cognitive skills?	with prior investigations into LLMs’ performance	133
084	<b>RQ2.</b> How does the evidential strength of student re-	in math-related contexts, and relevant benchmarks,	134
085	sponses affect LLM’s cognitive skill diagnosis	highlighting their current gaps.	135
086	performance?		
087	To answer these questions, we (i) constructed	<b>2.1 Cognitive Diagnosis in Mathematics</b>	136
088	MATHCOG, an expert-crafted benchmark dataset	<b>Education</b>	137
089	designed to evaluate cognitive skill diagnosis. In		
090	collaboration with 5 education experts and 15 mid-	Several frameworks have been proposed to char-	138
091	dle school teachers, we curated 12 middle school	acterize cognitive skills in learning and assess-	139
092	math topics and 110 problems, each with 50+ stu-	ment, including Cognitive Diagnostic Models	140
093	dent responses diagnosed by teachers based on a	(CDMs) (Leighton and Gierl, 2007), the PISA	141
094	problem-specific diagnostic checklist grounded in	Assessment and Analytical Framework (OECD,	142
095	the TIMSS cognitive framework (Mullis, 2017).	2023), and the TIMSS assessment frame-	143
096	Teachers provided binary judgements ( <i>Yes/No</i> ) for	work (Mullis, 2017). Among these, the TIMSS	144
097	each skill and annotated the presence of support-	is a comprehensive and math-specific framework	145
098	ing evidence as <i>Evident</i> or <i>Vague</i> , yielding 3,036	for cognitive diagnosis, comprising content and	146
099	diagnostic items across 639 student responses in	cognitive domains. The cognitive domain, which	147
100	total. Unlike existing benchmarks that rely on syn-	evaluates knowledge application, is divided into	148
101	thetic or automatically generated data, all items in	three key areas: Knowing ( <i>recalling</i> definitions,	149
102	MATHCOG are grounded in real student responses	<i>recognizing</i> mathematical entities, <i>classifying or</i>	150
103	and expert-agreed diagnostic judgments, which are	<i>ordering</i> quantities, <i>computing</i> ), Applying ( <i>deter-</i>	151
104	inherently difficult to collect at scale in educational	<i>mining</i> strategies, <i>representing</i> problems, <i>imple-</i>	152
105	settings.	<i>menting</i> solution procedures), and Reasoning ( <i>an-</i>	153
106	Using MATHCOG, we evaluated 17 closed- and	<i>alyzing, justifying</i> ). This research aims to explore	154
107	open-source LLMs spanning multiple model fam-	whether LLMs can substitute for the mapping to	155
108	ilies, sizes, and capabilities. Results indicate that	enable scalable and explainable cognitive diagno-	156
109	current LLMs struggle to reliably diagnose stu-	sis. To our knowledge, this is the first systematic	157
110	dents’ cognitive skills (all F1 scores < 0.521),	investigation of LLMs in the context of the TIMSS	158
111	with performance degrading sharply under vague-	framework.	159
112	evidence conditions and frequently misattributing		
113	weak evidence as strong. Qualitative error analysis	<b>2.2 LLM Capabilities in Mathematical Tasks</b>	160
114	further reveals the model’s failure patterns across		
115	model types. In particular, models often <i>misidentify</i>	Recent advances have enabled LLMs to achieve	161
116	or <i>hallucinate</i> supporting evidence, or <i>over-infer</i>	remarkable performance in mathematical problem-	162
117	students’ cognitive skills from incomplete work,	solving (Guo et al., 2025; OpenAI, 2024), yet they	163
118	showing that errors cascade throughout the diag-	show relative weaknesses in diagnosing student	164
119	nostic pipeline from recognizing student responses	abilities (Macina et al., 2025; Weitekamp et al.,	165
120	to interpreting them based on the rubric.	2025). These diagnostic challenges are further	166
121	Our contributions are threefold. (1) We intro-	complicated by current models’ struggles with mul-	167
122	duce MATHCOG, an expert-crafted benchmark	timodal inputs such as visual elements and hand-	168
123	dataset for cognitive skill diagnosis. (2) Using	written content (Zhang et al., 2024b; Baral et al.,	169
124	MATHCOG, we evaluate 17 large language models	2025; Liu et al., 2024). While Ma et al. (2025)	170
125	and analyze the effects of multimodality, reasoning,	demonstrated LLMs’ potential for cross-domain	171
126	model size, and evidential strength on diagnosis	cognitive diagnosis, existing research has limited	172
127	performance. (3) We analyze diagnostic error pat-	systematic exploration of LLMs’ diagnostic capa-	173
128	terns under vague-evidence conditions, identifying	bilities for cognitive skills in math specifically. We	174
129	recurring failure modes across model types.	address this critical gap in the field by providing the	175
130	<b>2 Related Work</b>	first comprehensive evaluation of LLMs’ ability to	176
131	We review existing frameworks for characterizing	diagnose cognitive skills in mathematical contexts	177
132	cognitive skills in learning and assessment, along	using an established educational framework.	178

179	<b>2.3 Benchmarks for Evaluating Mathematical</b>	<b>3.1 Diagnostic Checklist</b>	224
180	<b>Assessment Tasks</b>		
181	In educational settings, several benchmarks have	To systematically diagnose cognitive skills re-	225
182	been proposed to evaluate LLMs’ performance in	flected in student responses, we developed a di-	226
183	mathematical assessment tasks, including grading,	agnostic checklist commonly applicable to isomor-	227
184	skill recognition, and error analysis. For exam-	phic problems. Each checklist consists of binary	228
185	ple, <a href="#">Lucy et al. (2024)</a> introduced MathFish, which	question items mapped to one of the 15 cogni-	229
186	aligns 9,900 problems with 385 K–12 standards to	tive skills defined in the TIMSS 2019 assessment	230
187	assess models’ ability to identify targeted mathe-	framework. The checklist items were adapted from	231
188	matical skills and concepts. Automated scoring of	TIMSS skill descriptions and refined through ex-	232
189	open-ended constructed responses in math word	pert review by five mathematics curriculum and	233
190	problems has also been studied ( <a href="#">Hellman et al.,</a>	evaluation experts with PhD degrees in education	234
191	<a href="#">2023</a> ). Other benchmarks focus on misconception	and practical experience in making math assess-	235
192	detection in multiple-choice questions ( <a href="#">Kaggle and</a>	ment guidelines. Experts gave feedback on the clar-	236
193	<a href="#">Eedi, 2020</a> ) or handwriting recognition in student	ity, granularity, and validity of the checklist items.	237
194	work ( <a href="#">Baral et al., 2025</a> ). While prior research has	The experts also commented that cognitive skills	238
195	focused on error detection and scoring, diagnosing	could be reliably assessed from the given problem	239
196	complex cognitive skills from students’ open-ended	types. They noted that the problems predominantly	240
197	handwritten responses remains underexplored and	require numerical computation and application of	241
198	under-resourced. To address this gap, we introduce	known procedures, which makes them well-suited	242
199	a novel dataset for cognitive skill diagnosis that	for assessing <i>knowing</i> and <i>applying</i> skills. In con-	243
200	enables systematic evaluation of LLMs’ ability to	trast, higher-order <i>reasoning</i> skills (e.g., justifying,	244
201	interpret implicit reasoning and diagnose students’	analyzing, generalizing) are less consistently ob-	245
202	cognitive skills beyond final answers.	servable from students’ written responses in this	246
203		dataset. Based on this observation, we intentionally	247
204	<b>3 Dataset: MATHCOG</b>	scoped our checklist to the <i>knowing</i> and <i>applying</i>	248
205	To evaluate LLMs’ performance on cognitive skill	domains to ensure reliable and evidence-grounded	249
206	diagnosis, we created a benchmark dataset com-	diagnosis. The checklists were refined through two	250
207	prising secondary-school math problems, handwri-	iterative rounds, with each version independently	251
208	ten student responses, diagnostic checklists, and	reviewed by two experts.	252
209	teacher-generated verdicts (Table 1). The math		
210	problem and student response data are from AI-	<b>3.2 Teacher-generated Verdict</b>	253
211	Hub <sup>1</sup> , which provides OCR transcriptions of hand-	We recruited 15 middle school math teachers to	254
212	written work. Each topic includes multiple iso-	evaluate 796 student responses based on prede-	255
213	morphic problems that share the same problem-	defined diagnostic checklists. Teachers had an av-	256
214	-solving procedures but differ in numerical val-	erage of 6.1 years of teaching experience (SD =	257
215	ues (e.g., “Solve $x^2 + 2x - 3 = 0$ ” and “Solve	4.3; min = 2.5, max = 20). Each check item was	258
216	$x^2 + x - 2 = 0$ ”) ( <a href="#">Reed et al., 1990</a> ; <a href="#">Morrison</a>	assessed along two dimensions: evidential strength	259
217	<a href="#">et al., 2015</a> ). From this data, we focused on topics	( <i>Evident/Vague</i> ) and correctness ( <i>Yes/No</i> ). “ <i>Evident</i> ”	260
218	from grades 7–9, where problems are sufficiently	meant there was clear evidence to support the judg-	261
219	complex to elicit students’ cognitive processes. We	ment, whereas “ <i>Vague</i> ” signified insufficient evi-	262
220	further excluded topics with fewer than 50 student	dence. A “ <i>Yes</i> ” response indicated that the student	263
221	responses to ensure sufficient data for reliable eval-	fully demonstrated the cognitive action specified,	264
222	uation. Through this filtering process, the resulting	while a “ <i>No</i> ” indicated otherwise, including par-	265
223	dataset contains 15 topics, 137 problems, and 796	tially demonstrated responses. To account for the	266
	student responses.	subjectivity in cognitive diagnosis, each student’s	267
		response was evaluated by three teachers with over-	268
		lapping assignments, allowing us to measure inter-	269
		rater agreement (see Table 6 in Appendix. Fol-	270
		lowing the threshold established in <a href="#">Graham et al.</a>	271
		(2012), we excluded topics with agreement below	272
		70%, resulting in a final dataset of 12 topics, 110	273

<sup>1</sup>This research used datasets from ‘The Open AI Dataset Project (AI-Hub, S. Korea)’. All data information can be accessed through [AIHub](#).

Problem	Student Response	Diagnostic Checklist	Verdict
There is a trapezoid with a lower side length of 8 cm and a height of 4 cm. If the area of this trapezoid is not less than 28 cm <sup>2</sup> , find how much more cm the length of the upper side of the trapezoid must be.	$(8+x) \times 4 \times \frac{1}{2} \leq 28$ $6+2x \geq 28$ $2x \geq 12$ $x \geq 6$	<b>Recall:</b> Does the student remember the formula for the area of a shape correctly?	Evident Yes
		<b>Compute:</b> Has the student calculated the linear and constant terms correctly?	Vague No
		<b>Determine:</b> Did the student know the need to set up an equation and then solve it to find the range of solutions that meet the conditions?	Evident Yes
		<b>Represent:</b> Has the given situation been expressed correctly?	Evident No
		<b>Implement:</b> When simplifying an expression, does the student keep the expression correct by performing the same operation on both sides?	Vague Yes

Table 1: A sample from MATHCOG. Each data point is composed of a math problem, student response, relevant diagnostic checklist, and verdict for each check item.

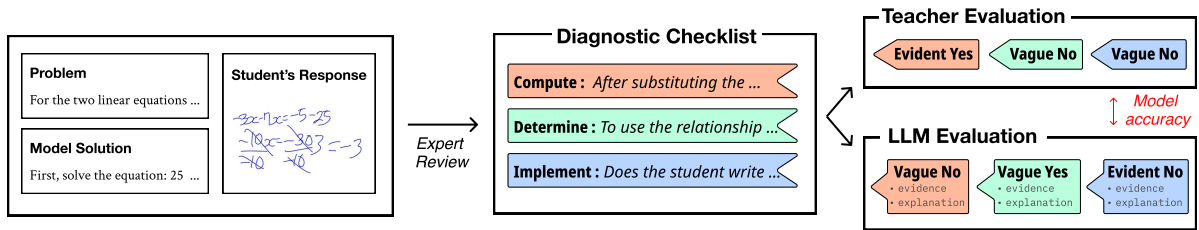


Figure 1: Diagnostic pipeline of MATHCOG. Each instance includes a math problem and student response from an open dataset with a TIMSS-aligned diagnostic checklist developed with mathematics curriculum experts. Math teachers evaluate each checklist item for correctness and evidential strength, and LLM predictions are compared with these teacher verdicts.

274 problems, and 639 student responses. As a result, 275 the benchmark focuses on diagnostic cases with 276 stable expert agreement, providing a reliable basis 277 for evaluating model performance. These topics 278 cover three-fourths of the content domains defined 279 in TIMSS 2019 (see Table 5 in Appendix), and the 280 dataset remains reasonably sized given the substan- 281 tial human effort required for expert annotation.

## 282 4 Experimental Setting

283 Using MATHCOG as a benchmark dataset, we evalu- 284 ated a diverse set of LLMs with varying input 285 modalities, reasoning capabilities, model sizes, and 286 prompting strategies to address RQ1, and analyzed 287 how the evidential strength of student responses 288 affects diagnosis performance to address RQ2.

289 **Prompting.** We designed our prompting pro- 290 tocol rigorously to ensure controlled and repro- 291 ducible evaluation. LLMs were instructed to evalu- 292 ate each diagnostic check item, given a math 293 problem, its solution, a student’s response, and 294 a diagnostic checklist. Student’s response im- 295 age inputs were provided as OCR transcriptions, 296 with mathematical formulas and visual cues (e.g., 297 strikethroughs) represented in LaTeX. For multi- 298 modal settings, we additionally supplied images

299 of student responses. Since the original inputs 300 were in Korean, we machine-translated them into 301 English to prevent possible performance degrada- 302 tion due to language (Achiam et al., 2023); we 303 used Google Translate API for batch translation 304 and manually verified them. We employed Chain- 305 of-Thought prompting (Liu et al., 2024) to guide 306 LLMs to systematically address each check item 307 by first restating its content, identifying relevant 308 evidence, providing an explanation, and delivering 309 a final verdict. The verdict followed one of the 310 four categories used by teachers in MATHCOG. To 311 further examine the robustness of our prompt de- 312 sign, we additionally evaluated five representative 313 models under few-shot prompting using two exem- 314 plars covering all four verdict types (Table 8). The 315 few-shot experiment results and full system and 316 user prompts are provided in the Appendix.

317 **Models.** We evaluated 17 open and proprietary 318 LLMs from multiple model families, selected to 319 analyze the effects of multimodality, reasoning ca- 320 pability, and model size. For reasoning effort level, 321 SOTA models with controllable reasoning budgets 322 were each evaluated under both low and high ef- 323 fort conditions. For multimodality, selected mod- 324 els were tested under both text-only and image-

augmented conditions. For model size, models were grouped into large, medium, and small categories.

**Metrics.** We evaluated LLM outputs against teacher-provided ground-truth diagnostic labels. Overall diagnosis performance was assessed using **macro F1 score** and **accuracy** over the four verdict categories. We report macro F1 to mitigate the effect of label imbalance across verdicts (see Table 7 in Appendix). To address **RQ2**, we analyzed how the *evidential strength* of student responses affects diagnosis performance. To formalize this analysis, we treat each diagnosis instance as a (*student response, cognitive skill*) pair, indexed by  $i$ . For each instance  $i$ , the ground-truth labels consist of a skill judgment  $y_i \in \{\text{Yes, No}\}$  and evidential strength  $e_i \in \{\text{Evident, Vague}\}$ , and the corresponding model predictions  $\hat{y}_i$  and  $\hat{e}_i$ . To characterize how models handle evidential strength, we define two complementary metrics that capture distinct failure modes. Prior work has shown that model-generated explanations can be plausible yet fail to faithfully reflect their actual reasoning basis (Jacovi and Goldberg, 2020; Turpin et al., 2023; Lyu et al., 2024). In educational diagnosis, an analogous failure occurs when a model falsely attributes evidential support to its verdicts, either by overstating vague evidence, or by asserting evidence for incorrect diagnoses. We define two metrics to quantify each form of this failure. First, we define **evidence over-attribution (OverAttr)** as the frequency with which a model assigns “Evident” to cases where the ground-truth evidential strength is “Vague”:

$$\text{OverAttr} = \frac{|\{i \mid e_i = \text{Vague} \wedge \hat{e}_i = \text{Evident}\}|}{|\{i \mid e_i = \text{Vague}\}|} \quad (1)$$

This metric reflects the model’s tendency to overstate evidential support under vague evidence cases, independent of diagnosis correctness. Second, we define **evidence false-attribution (FalseAttr)** as the frequency with which a model assigns Evident to incorrect diagnoses:

$$\text{FalseAttr} = \frac{|\{i \mid \hat{y}_i \neq y_i \wedge \hat{e}_i = \text{Evident}\}|}{|\{i \mid \hat{y}_i \neq y_i\}|} \quad (2)$$

This metric captures a particularly misleading failure mode, where erroneous diagnoses are accompanied by strong evidential claims.

## 5 Results

This section reports results addressing our research questions.

### 5.1 RQ1. How do different LLMs perform in diagnosing students’ cognitive skills?

Overall, all evaluated LLMs showed limited performance in cognitive skill diagnosis, with most macro F1 scores below 0.5 (Figure 7, Table 2). The only model to exceed F1 of 0.5 was Gemini-3.1-Pro, which achieved F1 scores of 0.509 and 0.521 under low and high reasoning effort conditions, respectively, with the high-effort setting also yielding the highest accuracy (0.778) across all evaluated conditions. While overall accuracy was relatively moderate ( $M = .683$ ,  $SD = .054$ ), it obscures important diagnostic failures. Models frequently over-attributed evidential strength (OverAttr;  $M = .544$ ,  $SD = .112$ ) and falsely attributed evidence (FalseAttr;  $M = .579$ ,  $SD = .118$ ), often asserting strong evidence even when diagnoses were incorrect. Notably, even the best-performing models in terms of F1 score and accuracy showed high OverAttr and FalseAttr values, suggesting that performance gains often co-occur with unreliable evidential judgement.

Analysis of skill-specific performance reveals further insights into the limitations of current LLMs. **When averaged across all models, no skill category achieved an F1 score above 0.5**, indicating that fine-grained diagnosis remains challenging even at the individual skill level (Figure 4, Figure 9). While Gemini-3.1-Pro occasionally exceeded this threshold on individual skills such as *Implement* and *Recognize*, this was not consistent across models. Despite this overall limitation, we observed a performance gap between the *Knowing* and *Applying* cognitive skills. Contrary to our expectation that *Knowing* skills (e.g., *Recall*, *Recognize*) would be easier due to their surface-level nature, models performed better on *Applying* skills in both F1 score ( $t = 6.09$ ,  $p < .001$ ) and accuracy ( $t = 4.50$ ,  $p < .001$ ). However, this performance gain came with increased evidential errors. *Applying* skills exhibited higher evidence over-attribution ( $t = 5.42$ ,  $p < .001$ ) and false-attribution ( $t = 5.85$ ,  $p < .001$ ) than *Knowing* skills, indicating a tendency to overstate evidential support when diagnosing higher-level cognitive skills. Notably, *Recall* stood out as a particularly difficult skill, showing the lowest 1-OverAttr and

Family	R-Low				R-High				+Image				Large				Medium				Small			
	F1	Acc	OA	FA	F1	Acc	OA	FA	F1	Acc	OA	FA	F1	Acc	OA	FA	F1	Acc	OA	FA	F1	Acc	OA	FA
<b>GPT</b>	<i>GPT-5.2</i> .428 .664 .410 .482				<i>GPT-5.2</i> .434 .667 .326 .471				<i>GPT-4o</i> .448 .743 .656 .681				<i>GPT-4o</i> .412 .672 .648 .658				-	-	-	-	<i>GPT-4o-mini</i> .352 .622 .767 .784			
<b>Claude</b>	<i>Claude-4.6-Sonnet</i> .432 .646 .493 .570				<i>Claude-4.6-Sonnet</i> .459 .662 .419 .586				<i>Claude-3.5-Sonnet</i> .413 .691 .626 .559				<i>Claude3.5-Sonnet</i> .416 .656 .551 .574				-	-	-	-	-	-	-	-
<b>Gemini</b>	<i>Gemini-3.1-pro</i> .509 .775 .480 .628				<i>Gemini-3.1-pro</i> .521 .778 .454 .628				<i>Gemini-1.5-Flash</i> .429 .702 .471 .471				<i>Gemini-1.5-Pro</i> .440 .706 .520 .635				<i>Gemini-1.5-Flash</i> .432 .679 .502 .465				<i>Gemini-1.5-Flash-8b</i> .348 .559 .537 .525			
<b>DeepSeek</b>	-	-	-	-	<i>DeepSeek-R1</i> .442 .773 .767 .862				-	-	-	-	<i>DeepSeek-V3</i> .417 .714 .612 .595				-	-	-	-	-	-	-	-
<b>Llama</b>	<i>Llama-4-Maverick</i> .468 .758 .551 .623				<i>Llama-4-Maverick</i> .478 .763 .555 .622				-	-	-	-	<i>Llama-3.1-405B</i> .385 .624 .467 .307				<i>Llama-3.3-70B</i> .397 .635 .436 .409				<i>Llama-3.1-8B</i> .323 .654 .705 .626			
<b>Qwen</b>	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	<i>Qwen-2.5-72B</i> .416 .711 .568 .500				<i>Qwen-2.5-7B</i> .339 .705 .758 .798			

Table 2: Performance of LLMs on cognitive skill diagnosis, grouped by model family. **OA** and **FA** denote OverAttr and FalseAttr (lower is better). **R-Low** and **R-High** reflect two reasoning effort levels of the latest frontier models (e.g., thinking budget or inference compute); **+Image** denotes models prompted with student response images in addition to OCR text; and **Large**, **Medium**, and **Small** correspond to model size tiers based on parameter count. Specific models are indicated in italics above each value group. **Blue**: best; **Red**: worst value per column.

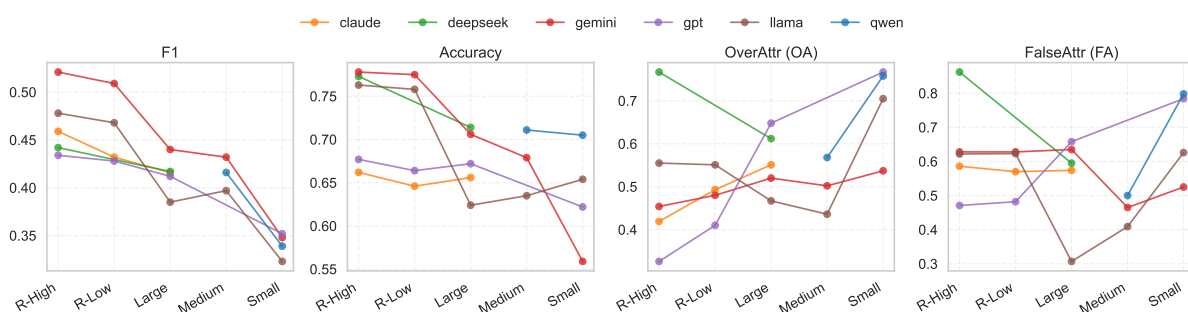


Figure 2: Performance of six LLM families across five evaluation conditions. R-High and R-Low denote high and low reasoning effort levels of the latest frontier models; Large, Medium, and Small refer to model size tiers. Missing points indicate no available model for that condition.

1-FalseAttr values among all skills (Figure 4), suggesting that models were most prone to misattributing evidence when diagnosing this higher-level cognitive skill.

We further examined the impact of **multimodality**, **reasoning effort level**, and **model size** as factors that may influence diagnostic performance (Table 2). **Multimodal input provided modest but inconsistent benefits (Figure 3)**. Models supporting image input consistently outperformed their text-only counterparts in accuracy, indicating that access to handwritten layouts and visual cues helps mitigate errors introduced by OCR transcription. However, improvements in F1 score were inconsistent, suggesting that visual input can sometimes introduce noise and does not uniformly improve fine-grained skill diagnosis. **Higher reasoning effort generally yielded better diagnostic performance**. Among sota models evaluated under both reasoning effort conditions, high effort consistently outperformed low effort in F1 and

accuracy across model families (Figure 2). Evidence over-attribution (OverAttr) was also somewhat lower under higher reasoning effort, while false-attribution (FalseAttr) remained comparable across the two conditions. These results suggest that increased reasoning compute offers modest but consistent benefits for cognitive diagnosis, though the gains were not large enough to substantially close the performance gap observed across all evaluated models. **Model size showed a strong positive relationship with F1 performance** (Spearman’s  $\rho = .716, p = .009$ ), indicating that larger models more robustly interpret student responses (Figure 2). However, size was not significantly associated with accuracy ( $\rho = .399, p = .199$ ), OverAttr ( $\rho = -.417, p = .177$ ), or FalseAttr ( $\rho = -.246, p = .441$ ), though smaller models tended to show higher rates of overconfident evidential claims. These results suggest that scaling up model size improves diagnostic judgment but does not guarantee more reliable evidential reasoning.



Figure 3: Effect of image input on diagnostic performance across three multimodal models (Claude-3.5-Sonnet, Gemini-1.5-Flash, GPT-4o).

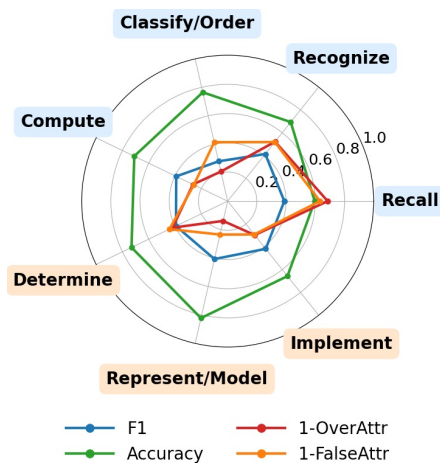


Figure 4: Radar chart showing model performance across seven cognitive skills. The metrics include F1, Accuracy, 1-OverAttr, and 1-FalseAttr (higher is better). Skills are grouped by cognitive domain: Knowing (blue) and Applying (orange).

Taken together, these findings suggest that model capability and scale alone are insufficient to ensure reliable cognitive skill diagnosis.

## 5.2 RQ2. How does the evidential strength of student responses affect LLM’s cognitive skill diagnosis performance?

Our results reveal a **strong dependency of LLM diagnostic performance on the evidential strength** of student responses. Across all models, both F1 score ( $t = 64.60, p < .001$ ) and accuracy ( $t = 17.80, p < .001$ ) were substantially higher when student responses contained *evident* evidence for a cognitive skill than when evidence was *vague* (Fig. 5). The magnitude of this gap was substantial (mean  $\Delta F1 = .51$ ; mean  $\Delta Acc = .41$ ), indicating that current LLMs rely heavily on explicit evidence and exhibit limited robustness when reasoning un-

der weak or implicit evidential conditions.

Beyond degraded performance, **models also exhibited over-attribution bias under vague evidence**. Specifically, they frequently labeled responses as *Evident* even when the ground-truth evidential strength was *Vague*, resulting in consistently high **OverAttr** values. More concerning, models also showed a strong tendency toward **FalseAttr**, in which they asserted *evident* evidence despite arriving at incorrect cognitive diagnoses. In such cases, models produced unsupported explanations that appeared plausible yet were not grounded in students’ actual responses. Notably, these two tendencies were strongly correlated across models (Spearman’s  $\rho = .837, p < .001$ ), indicating that models prone to over-attributing evidential strength are also more likely to produce incorrect diagnoses with asserted evidence. This behavior is particularly problematic for educational use, as it can mislead teachers and students by presenting incorrect diagnoses with seemingly well-justified rationales (Kim et al., 2025).

To better understand how such errors arise, we conducted a qualitative error analysis of model outputs under incorrectly diagnosed vague evidence cases. We examined all components generated in the models’ CoT responses, including the identified evidence, explanations, and final verdicts, focusing on how the models misinterpreted students’ answers, generated unsupported reasoning, and ultimately reached incorrect diagnoses. Three authors collaboratively analyzed 372 cases (13.2%) sampled from a total of 2,816 error instances, covering 17 model responses, and iteratively derived five recurring error types with substantial inter-rater reliability (Fleiss’  $\kappa = 0.74$ ). The remaining cases were then independently labeled using the agreed error

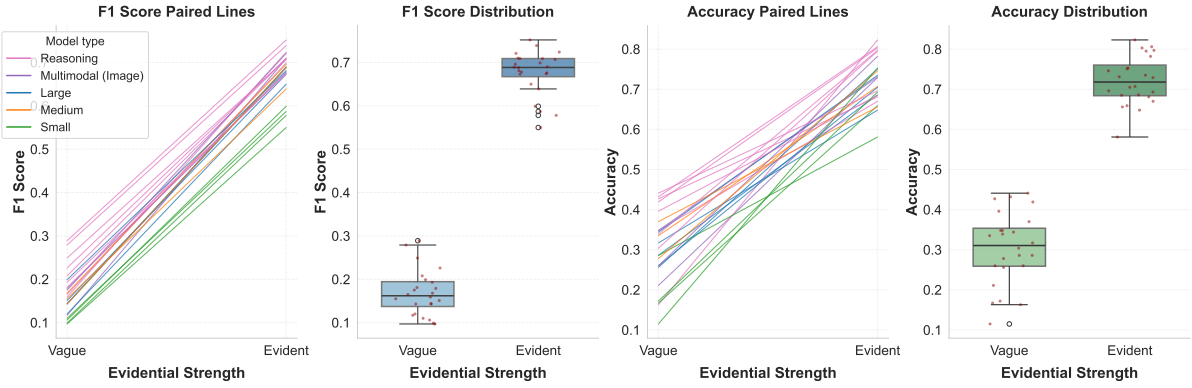


Figure 5: Effect of evidential strength on LLM diagnostic performance across all evaluated models. Paired line plots (left) show per-model changes in F1 and accuracy from vague to evident evidence conditions, with each line representing one model colored by model type. Box plots (right) summarize the overall distribution of F1 and accuracy under each condition.

Error code	Description	Count (%)
E1. Evidence Misidentification	Incorrectly identifying or relying on insufficient evidence from the student response.	659 (23.40%)
E2. Rubric Misinterpretation	Misunderstanding the diagnostic rubric or evaluation criteria.	279 (9.91%)
E3. Over-Inference	Inferring unstated reasoning beyond the information provided in the student response.	942 (33.45%)
E4. Consistency Issue	Inconsistencies between evidence, explanation, and final verdict.	379 (13.46%)
E5. Hallucination	Introducing content or reasoning not present in the student response.	557 (19.78%)

Table 3: Error types and their distribution identified from qualitative analysis of LLM diagnostic failures under vague evidence conditions.

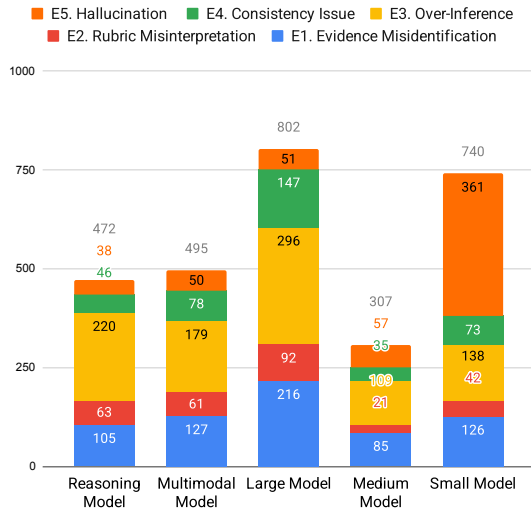


Figure 6: Error type distribution across model categories. Each bar represents the total number of diagnostic errors under vague evidence conditions, broken down by each error type.

taxonomy. Table 3 summarizes the distribution of error types across model categories. Errors related to *over-inference* were most prevalent (33.5%), followed by *evidence misidentification* (23.4%) and *hallucination* (19.8%). We further observed systematic differences across model types (Figure 6): smaller models exhibited a higher proportion of

*hallucination* errors (48.78%), whereas reasoning-oriented models committed *over-inferred* more frequently (46.61%) beyond the information explicitly provided in student responses.

## 6 Conclusion

We introduced MATHCOG, a benchmark dataset comprising 3,036 diagnostic verdicts across 639 student responses, annotated by teachers using a TIMSS-grounded cognitive framework. Evaluating 17 LLMs, we find that diagnostic performance remains limited across all model types and cognitive skill categories, with performance degrading sharply under vague evidence conditions. Error analysis reveals systematic failure patterns that vary across model types. These findings suggest that the core challenge lies in the inherent difficulty of interpreting implicit diagnostic evidence, rather than in model capability alone. Responsible deployment of LLMs in educational assessment will require task designs that elicit more explicit student reasoning and teacher-in-the-loop workflows that support critical interpretation of model outputs.

## 545 Limitations

546 This work has several limitations. First, the Math-  
547 Cog benchmark is constructed from a specific edu-  
548 cational dataset consisting of middle-school math-  
549 ematics responses, which may limit the generaliz-  
550 ability of our findings to other subjects, grade lev-  
551 els, or educational contexts. Second, the cognitive  
552 skills in this work are defined using TIMSS-based  
553 categories (*know*, *apply*, *reason*), which provide  
554 an operational taxonomy but may simplify the di-  
555 versity and complexity of students’ real problem-  
556 solving processes. Third, cognitive skill diagnosis  
557 is inherently subjective, and teacher annotations  
558 may reflect individual interpretations of student  
559 work. Although we mitigated this by requiring  
560 three independent teacher judgments per item and  
561 excluding topics with inter-rater agreement below  
562 70%, residual subjectivity in the ground-truth la-  
563 bels may still affect evaluation reliability. Finally,  
564 while we explored multiple prompting strategies in-  
565 cluding Chain-of-Thought and few-shot prompting,  
566 our evaluation is limited to a specific task formula-  
567 tion; alternative formulations or more specialized  
568 prompting approaches may lead to different model  
569 behaviors and diagnostic outcomes.

## 570 Ethical Impact

571 This work analyzes students’ mathematical  
572 problem-solving responses to study whether large  
573 language models can diagnose cognitive skills  
574 from partial evidence. The dataset consists of  
575 anonymized student work and does not contain per-  
576 sonally identifiable information. A potential risk  
577 is that LLM-based diagnostic systems may pro-  
578 duce incorrect or overconfident interpretations of  
579 reasoning if used without appropriate verification.  
580 By systematically identifying such failure patterns,  
581 our benchmark aims to support the development  
582 and evaluation of more reliable reasoning analysis  
583 systems.

## 584 References

585 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama  
586 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
587 Diogo Almeida, Janko Altenschmidt, Sam Altman,  
588 Shyamal Anadkat, and 1 others. 2023. Gpt-4 techni-  
589 cal report. *arXiv preprint arXiv:2303.08774*.

590 Janice Ahn, Rishu Verma, Renze Lou, Di Liu, Rui  
591 Zhang, and Wenpeng Yin. 2024. Large language  
592 models for mathematical reasoning: Progresses and  
593 challenges. *arXiv preprint arXiv:2402.00157*.

John R Anderson. 1982. Acquisition of cognitive skill. *Psychological review*, 89(4):369. 594  
595

Sami Baral, Li Lucy, Ryan Knight, Alice Ng, Luca Sol-  
dainei, Neil T Heffernan, and Kyle Lo. 2025. Drawe-  
dumath: Evaluating vision language models with  
expert-annotated students’ hand-drawn math images. *arXiv preprint arXiv:2501.14877*. 596  
597  
598  
599  
600

Aniket Didolkar, Anirudh Goyal, Nan Rosemary Ke,  
Siyuan Guo, Michal Valko, Timothy Lillicrap, Danilo  
Jimenez Rezende, Yoshua Bengio, Michael C Mozer,  
and Sanjeev Arora. 2024. Metacognitive capabilities  
of llms: An exploration in mathematical problem  
solving. *Advances in Neural Information Processing  
Systems*, 37:19783–19812. 601  
602  
603  
604  
605  
606  
607

Meng Fang, Xiangpeng Wan, Fei Lu, Fei Xing, and  
Kai Zou. 2024. Mathodyssey: Benchmarking math-  
ematical problem-solving skills in large language  
models using odyssey math data. *arXiv preprint  
arXiv:2406.18321*. 608  
609  
610  
611  
612

Matthew Graham, Anthony Milanowski, and Jackson  
Miller. 2012. Measuring and promoting inter-rater  
agreement of teacher and principal performance rat-  
ings. *Online Submission*. 613  
614  
615  
616

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao  
Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shi-  
rong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025.  
Deepseek-r1: Incentivizing reasoning capability in  
llms via reinforcement learning. *arXiv preprint  
arXiv:2501.12948*. 617  
618  
619  
620  
621  
622

Shibo Hao, Yi Gu, Haotian Luo, Tianyang Liu, Xiyan  
Shao, Xinyuan Wang, Shuhua Xie, Haodi Ma,  
Adithya Samavedhi, Qiyue Gao, and 1 others. 2024.  
Llm reasoners: New evaluation, library, and analysis  
of step-by-step reasoning with large language models.  
*arXiv preprint arXiv:2404.05221*. 623  
624  
625  
626  
627  
628

Scott Hellman, Alejandro Andrade, and Kyle Haber-  
mehl. 2023. Scalable and explainable automated  
scoring for open-ended constructed response math  
word problems. In *Proceedings of the 18th Workshop  
on Innovative Use of NLP for Building Educational  
Applications (BEA 2023)*, pages 137–147. 629  
630  
631  
632  
633  
634

Charles Henderson, Edit Yerushalmi, Vince H Kuo, Pa-  
tricia Heller, and Kenneth Heller. 2004. Grading  
student problem solutions: The challenge of sending  
a consistent message. *American Journal of Physics*,  
72(2):164–169. 635  
636  
637  
638  
639

Alon Jacovi and Yoav Goldberg. 2020. Towards faith-  
fully interpretable nlp systems: How should we de-  
fine and evaluate faithfulness? In *Proceedings of the  
58th annual meeting of the association for computa-  
tional linguistics*, pages 4198–4205. 640  
641  
642  
643  
644

Hyoungwook Jin, Yoonsu Kim, Yeon Su Park, Bekzat  
Tilekbay, Jinho Son, and Juho Kim. 2024. Using  
large language models to diagnose math problem-  
solving skills at scale. In *Proceedings of the Eleventh  
ACM Conference on Learning@ Scale*, pages 471–  
475. 645  
646  
647  
648  
649  
650



## A Appendix

### A.1 Prompts

The blue text represents the programmatically filled arguments, and the orange text represents LLM-generated output.

#### A.1.1 System Prompt

```
# **Task Description**
You are a middle school math teacher tasked
with evaluating students' mathematical thinking
skills based on their responses to math
problems. Your goal is to analyze a given
student's response and determine whether they
exhibit specific cognitive skills in solving
the problem. Your evaluation must be strict
and evidence-based, meaning that every
verdict must be backed by direct evidence from
the response. If no clear evidence exists, do
not assume correctness.

## **Evaluation Categories**
For each thinking skill in the checklist, you
must classify the student's performance into
one of the following categories:

- Evident Yes: The student's response
provides clear and explicit evidence that the
check item is met. A direct quote from the
response can confirm this.
- Vague Yes: The response suggests that the
check item might be satisfied, but no specific
part of the response directly proves it.
- Evident No: The response explicitly
contradicts or fails to meet the check item,
with clear evidence demonstrating the error or
omission.
- Vague No: The response does not appear to
satisfy the check item, but there is no direct
evidence confirming whether the student
considered it or not.

## **Input Format**
You will receive the following data:
- Problem: A math problem given to a
student.
- Answer: The correct step-by-step solution.
- Response: The student's response to the
problem.
- Check Items: A set of specific skills to
evaluate.

## **Output Format**
Return a valid JSON object structured as
follows:
```json
{
  "skills": [
    {
      "checkItem": "<Check Item's [Label]
and the Following Question>",
      "evidence": "<Directly Quoted Part
of Response>",
      "explanation": "<Explanation About
Why the Evidence Supports the
Verdict>",
```

```
"verdict": "Evident Yes" | "Vague
Yes" | "Evident No" | "Vague No"
```

```
    }
  ]
}
```

#### A.1.2 User Prompt

```
# **Task**
student responses to math problems, extract
direct evidence, and strictly classify thinking
skills according to the given categories.

Return a valid JSON object structured as
follows:
```json
{
  "skills": [
    {
      "checkItem": "<Check Item's [Label]
and the Following Question>",
      "evidence": "<Directly Quoted Part
of Response>",
      "explanation": "<Explanation About
Why the Evidence Supports the
Verdict>",
      "verdict": "Evident Yes" | "Vague
Yes" | "Evident No" | "Vague No"
    }
  ]
}
```

Problem
There is a two-digit natural number whose tens
digit is 1. If the number that changes the tens
and ones digits of this natural number is 9 less
than 5 times the first number, find the first
number.

Answer
If the number in the ones place is  $x$ , this
natural number is  $10 + x$ . The number where the
tens digit and the ones digit are swapped is
 $10x + 1$  because the tens digit is  $x$  and the ones
digit is 1. The changed number is 9 smaller than
5 times the first number, so
 $10x + 1 = 5(10 + x) - 9$ 
 $10x + 1 = 5x + 41$ 
 $5x = 40$ 
 $x = 8$ 
Therefore, the first number is 18.

Response
 $1 \ x \ 1 \ 10x + 1 = 5(10 + x) - 9$ 
 $10x + 1 = 50 + 5x - 9$ 
 $41 \ 5x = 40 \ x = 8$ 

Check Items
- [recognize_0] Is the student aware that the
relationship between two numbers switched can be
expressed as
addition/subtraction/multiplication?
- [recognize_1] Has the student figured out all
the elements necessary to solve the problem
(original numbers, swapped numbers, difference or
multiple relationships, etc.)?
- [compute_0] Has the student performed algebraic
calculations correctly, such as expanding
expressions and organizing similar terms?
- [determine_0] Has the student chosen the
appropriate approach to solving the problem
(e.g., establishing a linear equation)?
- [represent_0] Are two-digit numbers expressed
correctly in the form  $10x + y$  (or equivalent)?
- [implement_0] Does the student proceed
consistently with the process of expanding the
equation and moving terms to find  $x$  (or  $y$ )
```

```

877 according to the planned solution strategy?
878
879 **Output:**
880 {
881   "skills": [
882     {
883       "checkItem": "[recognize_0] Is the
884       student aware that the relationship
885       between two numbers switched can be
886       expressed as
887       addition/subtraction/multiplication?",
888       "evidence": "10 x+1=5(10+x)-9",
889       "explanation": "The student sets up
890       the equation 10x+1=5(10+x)-9, which
891       shows they understand the relationship
892       between the original and swapped
893       numbers involves subtraction and
894       multiplication.",
895       "verdict": "Evident Yes"
896     }, {
897       "checkItem": "[recognize_1] Has the
898       student figured out all the elements
899       necessary to solve the problem
900       (original numbers, swapped numbers,
901       difference or multiple relationships,
902       etc.)?",
903       "evidence": "10 x+1=5(10+x)-9",
904       "explanation": "The student uses the
905       correct equation to represent the
906       relationship between the original and
907       swapped numbers, indicating they have
908       identified the necessary elements.",
909       "verdict": "Evident Yes"
910     },
911     ...
912   ]
913 }

```

## A.2 Data Collection

### A.2.1 Instructions Given to Education Experts and Teachers

The data collection procedures were approved by the Institutional Review Board (IRB). Prior to participation, informed consent was obtained from all participants, including both the education experts and the middle school teachers. The guideline materials used in this study are publicly available on the Open Science Framework (OSF). The materials can be accessed at: [https://osf.io/67kz9/overview?view\\_only=430e88aa45cf493db06f573dfb3eaae2](https://osf.io/67kz9/overview?view_only=430e88aa45cf493db06f573dfb3eaae2). These materials document the procedures and criteria used during the data collection process.

### A.2.2 Recruitment and Payment

*Education Experts for Checklist Refinement:* We recruited five education experts who work at the Korea Institute for Curriculum and Evaluation (KICE)<sup>2</sup> via cold email invitations. Given their expertise in curriculum design and assessment, each expert received 300,000 KRW (approximately USD 200) as compensation for approximately three hours of participation.

<sup>2</sup><https://www.kice.re.kr/main.do?s=english>

*Middle School Teachers for Cognitive Diagnosis:* We recruited 15 middle school teachers through advertisements posted in online teacher communities in Korea. Each teacher received 150,000 KRW (approximately USD 100) as compensation for approximately three hours of participation.

## A.3 MathCog Dataset Details

This section provides supplementary details of the MATHCOG dataset, including the diagnostic workflow, cognitive skill taxonomy, inter-rater reliability of teachers' diagnoses, dataset composition, and representative failure cases of model diagnoses.

### A.3.1 TIMSS Cognitive Skill Framework

The MATHCOG dataset is grounded in the TIMSS cognitive framework, which categorizes mathematical thinking into the domains of *knowing*, *applying*, and *reasoning*. Table 4 provides the detailed checklist items used to operationalize each cognitive skill in our diagnostic process.

### A.3.2 Problem Topics

The dataset covers a diverse set of middle school mathematics topics to ensure coverage of different cognitive processes and problem types. Table 5 lists the full set of problem topics included in MATHCOG.

### A.3.3 Inter-Rater Agreement

To ensure the reliability of cognitive skill annotations, we measured agreement among teachers on diagnostic labels. Table 6 reports the inter-rater agreement statistics for the annotated dataset.

### A.3.4 Dataset Distribution

We analyze the distribution of diagnostic labels across cognitive skills and evidence conditions to characterize the dataset composition. Table 7 presents detailed statistics of label distributions in MATHCOG.

## A.4 Experiment Result Details

### A.4.1 Few-Shot Prompting Results

Table 8 reports the diagnostic performance of five representative models under few-shot prompting, compared to their zero-shot counterparts. Few-shot prompting did not consistently improve F1 or accuracy, and in several cases led to decreased performance. However, few-shot examples did reduce tendencies toward evidence over-attribution

| Domains   | Cognitive Skills     | Description                                                                                                                                                                                                                      |
|-----------|----------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Knowing   | Recall               | Recall definitions, terminology, number properties, units of measurement, geometric properties, and notation (e.g., $a \times b = ab$ , $a + a + a = 3a$ ).                                                                      |
|           | Recognize            | Recognize numbers, expressions, quantities, and shapes. Recognize entities that are mathematically equivalent (e.g., equivalent familiar fractions, decimals, and percents; different orientations of simple geometric figures). |
|           | Classify/Order       | Classify numbers, expressions, quantities, and shapes by common properties.                                                                                                                                                      |
|           | Compute              | Carry out algorithmic procedures for $+$ , $-$ , $\times$ , $\div$ or a combination of these with whole numbers, fractions, decimals, and integers. Carry out straightforward algebraic procedures.                              |
|           | Retrieve             | Retrieve information from graphs, tables, texts, or other sources.                                                                                                                                                               |
|           | Measure              | Use measuring instruments; and choose appropriate units of measurement.                                                                                                                                                          |
| Applying  | Determine            | Determine efficient/appropriate operations, strategies, and tools for solving problems for which there are commonly used methods of solution.                                                                                    |
|           | Represent/Model      | Display data in tables or graphs; create equations, inequalities, geometric figures, or diagrams that model problem situations; and generate equivalent representations for a given mathematical entity or relationship.         |
|           | Implement            | Implement strategies and operations to solve problems involving familiar mathematical concepts and procedures.                                                                                                                   |
| Reasoning | Analyze              | Determine, describe, or use relationships among numbers, expressions, quantities, and shapes.                                                                                                                                    |
|           | Integrate/Synthesize | Link different elements of knowledge, related representations, and procedures to solve problems.                                                                                                                                 |
|           | Evaluate             | Evaluate alternative problem solving strategies and solutions.                                                                                                                                                                   |
|           | Draw conclusions     | Make valid inferences on the basis of information and evidence.                                                                                                                                                                  |
|           | Generalize           | Make statements that represent relationships in more general and more widely applicable terms.                                                                                                                                   |
|           | Justify              | Provide mathematical arguments to support a strategy or solution.                                                                                                                                                                |

Table 4: Fifteen cognitive skills and their descriptions defined in the TIMSS 2019 framework.

984 and false-attribution, suggesting that in-context ex- 1004  
985 amples may encourage more cautious, evidence- 1005  
986 grounded responses even when overall diagnostic 1006  
987 accuracy remains limited. 1007

#### 988 A.4.2 Model-Level Performance 1008

989 Figure 7 presents model-level performance across 1009  
990 different model families and types. Overall, per- 1010  
991 formance differences across models remain mod- 1011  
992 est, with no model consistently achieving high F1 1012  
993 scores. This suggests that cognitive skill diagnosis 1013  
994 remains challenging regardless of model architec- 1014  
995 ture or scale. 1015

#### 996 A.4.3 Model-wise Confusion Patterns 1016

997 Figure 8 shows model-wise confusion patterns 1017  
998 across verdict categories. 1018

#### 999 A.4.4 Skill-Level Performance 1000

1000 Figure 9 shows performance and evidential attri-  
1001 bution errors across cognitive skills. Performance  
1002 varies across skills, while evidence over-attribution  
1003 and false-attribution remain consistently high. This

1004 indicates that models struggle not only with skill  
1005 classification but also with accurately grounding  
1006 their predictions in appropriate evidence. 1007

1008 Consistent patterns are observed in the model-  
1009 wise results (Tables 9–12), which further show that  
1010 no model achieves uniformly strong performance  
1011 across all skills and that attribution-related errors  
1012 persist across models. 1013

#### 1014 A.5 Examples of Model Errors 1015

1016 Figure 10 provides illustrative examples of failure  
1017 cases in cognitive skill diagnosis. These examples  
1018 highlight common error patterns, such as misidenti-  
1019 fying relevant evidence, hallucinating unsupported  
1020 reasoning steps, and over-inferring skills from in-  
1021 complete student work. 1022

| Topics                                                                             | Content Domain | Difficulty |
|------------------------------------------------------------------------------------|----------------|------------|
| 1. Problems involving the digits of numbers                                        | Number         | 1, 3       |
| 2. Solving for unknowns under special conditions on the solution                   | Algebra        | 1, 2, 3    |
| 3. Finding unknowns when two equations have the same solution                      | Algebra        | 1, 2       |
| 4. Applying linear inequalities to geometric figures                               | Geometry       | 2, 3       |
| 5. Applying linear inequalities to pricing                                         | Algebra        | 1, 2       |
| 7. Using square roots to calculate lengths in geometric figures                    | Geometry       | 2, 3       |
| 8. Performing simple addition and subtraction of square roots                      | Number         | 1, 2       |
| 9. Applying multiplication formulas                                                | Number         | 2, 3       |
| 10. Determining coefficients or constants to complete the square                   | Number         | 1, 2, 3    |
| 11. Solving quadratic equations by factoring                                       | Algebra        | 1, 2       |
| 12. Rewriting expressions in perfect square form                                   | Algebra        | 1, 2       |
| 15. Finding the quadratic function given the vertex and another point on its graph | Algebra        | 1, 2       |

Table 5: The topics and difficulty levels of problems. This information comes from the original dataset’s metadata, which details the topic and difficulty level (1-3) of each problem. The isomorphic problems require the same mathematical concept to solve, but the difference in numbers makes one more tricky and complicated than the other.

| Topics      | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  | 10 | 11  | 12 | 13 | 14 | 15 |
|-------------|----|----|----|----|----|----|----|----|----|----|-----|----|----|----|----|
| % Agreement | 95 | 70 | 96 | 95 | 89 | 68 | 74 | 86 | 80 | 80 | 100 | 88 | 40 | 57 | 78 |

Table 6: The inter-rater percentage absolute agreement of each topic. The percentage indicates the ratio of unanimous verdicts in each teacher group.

|                   | Recall | Recognize | Classify/Order | Compute | Determine | Represent | Implement |
|-------------------|--------|-----------|----------------|---------|-----------|-----------|-----------|
| Evident Yes       | 413    | 434       | 92             | 522     | 441       | 122       | 307       |
| Evident No        | 62     | 56        | 8              | 125     | 54        | 33        | 140       |
| Vague Yes         | 20     | 16        | 2              | 26      | 25        | 0         | 18        |
| Vague No          | 35     | 28        | 0              | 20      | 16        | 3         | 18        |
| Student responses | 530    | 534       | 102            | 693     | 536       | 158       | 483       |

Table 7: Distribution of verdicts and number of diagnosed student responses in each cognitive skill. Note that the number of student responses can be larger than the sum of the four labels because some diagnostic checklists have two items for the same cognitive skill.

| Few-shots         | Precision   | Recall      | F1          | Accuracy    | OverAttr    | FalseAttr   |
|-------------------|-------------|-------------|-------------|-------------|-------------|-------------|
| claude-3-5-sonnet | .398        | .485        | .410        | .646        | <b>.515</b> | .647        |
| deepseek-V3       | .404        | <b>.493</b> | .404        | .624        | <b>.493</b> | <b>.497</b> |
| gemini-1.5-pro    | .429        | <b>.517</b> | .430        | .665        | <b>.427</b> | <b>.552</b> |
| gpt-4o            | .390        | <b>.485</b> | .396        | .637        | <b>.577</b> | <b>.587</b> |
| llama-3.1-405B    | <b>.442</b> | <b>.494</b> | <b>.429</b> | <b>.677</b> | .476        | .457        |

Table 8: Performance of five state-of-the-art models in few-shot experiments. Bold values denote improvements over the corresponding zero-shot results.

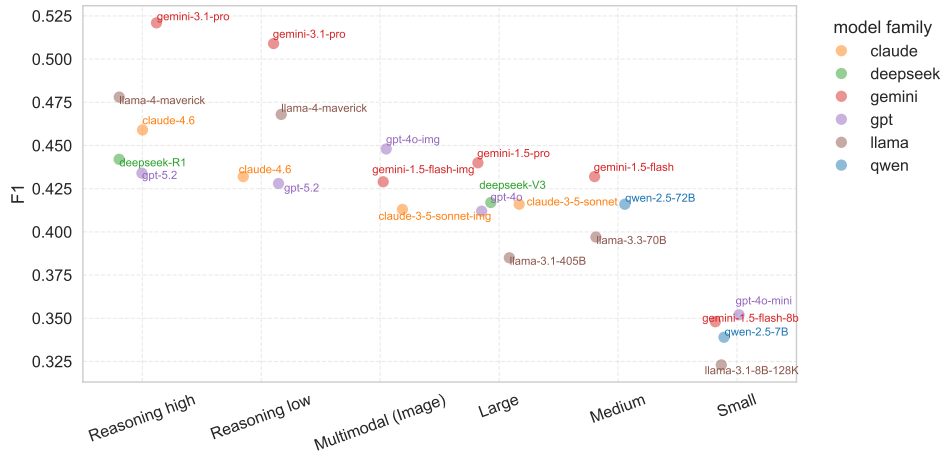


Figure 7: F1 score across by model type and family.

| Model                  | Recall | Recognize | Classify/Order | Compute | Determine | Represent | Implement |
|------------------------|--------|-----------|----------------|---------|-----------|-----------|-----------|
| GPT-5.2-Low            | 0.362  | 0.425     | 0.297          | 0.431   | 0.373     | 0.404     | 0.427     |
| GPT-5.2-High           | 0.346  | 0.420     | 0.302          | 0.411   | 0.353     | 0.456     | 0.448     |
| GPT-4o-Img             | 0.443  | 0.471     | 0.335          | 0.399   | 0.410     | 0.408     | 0.413     |
| GPT-4o                 | 0.400  | 0.368     | 0.333          | 0.392   | 0.357     | 0.395     | 0.440     |
| GPT-4o-mini            | 0.361  | 0.339     | 0.178          | 0.344   | 0.297     | 0.392     | 0.322     |
| Claude-Sonnet-4.6-Low  | 0.419  | 0.366     | 0.298          | 0.413   | 0.415     | 0.438     | 0.412     |
| Claude-Sonnet-4.6-High | 0.460  | 0.390     | 0.338          | 0.436   | 0.437     | 0.377     | 0.452     |
| Claude-3.5-Sonnet-Img  | 0.364  | 0.374     | 0.328          | 0.431   | 0.407     | 0.393     | 0.412     |
| Claude-3.5-Sonnet      | 0.358  | 0.382     | 0.325          | 0.424   | 0.438     | 0.473     | 0.396     |
| Gemini-3.1-Low         | 0.438  | 0.525     | 0.309          | 0.452   | 0.493     | 0.407     | 0.531     |
| Gemini-3.1-High        | 0.474  | 0.561     | 0.224          | 0.449   | 0.481     | 0.394     | 0.598     |
| Gemini-1.5-Flash-Img   | 0.362  | 0.415     | 0.226          | 0.387   | 0.405     | 0.417     | 0.462     |
| Gemini-1.5-pro         | 0.395  | 0.429     | 0.319          | 0.367   | 0.392     | 0.416     | 0.468     |
| Gemini-1.5-Flash       | 0.398  | 0.422     | 0.266          | 0.378   | 0.414     | 0.395     | 0.451     |
| Gemini-1.5-Flash-8b    | 0.276  | 0.330     | 0.185          | 0.337   | 0.334     | 0.456     | 0.349     |
| Deepseek-R1            | 0.469  | 0.483     | 0.357          | 0.376   | 0.384     | 0.416     | 0.395     |
| Deepseek-V3            | 0.402  | 0.439     | 0.224          | 0.372   | 0.390     | 0.392     | 0.380     |
| Llama-Maverick-Low     | 0.438  | 0.481     | 0.232          | 0.407   | 0.418     | 0.384     | 0.410     |
| Llama-Maverick-High    | 0.430  | 0.482     | 0.232          | 0.409   | 0.461     | 0.394     | 0.460     |
| Llama-3.1-405B         | 0.344  | 0.342     | 0.231          | 0.373   | 0.385     | 0.349     | 0.409     |
| Llama-3.3-70B          | 0.337  | 0.399     | 0.312          | 0.376   | 0.409     | 0.469     | 0.393     |
| Llama-3.1-8B-128K      | 0.322  | 0.266     | 0.334          | 0.316   | 0.271     | 0.285     | 0.291     |
| Qwen-2.5-72B           | 0.387  | 0.462     | 0.296          | 0.369   | 0.371     | 0.394     | 0.394     |
| Qwen-2.5-7B            | 0.321  | 0.337     | 0.231          | 0.343   | 0.338     | 0.463     | 0.308     |

Table 9: F1 by model (rows) and skill (columns).



Figure 8: Model-wise confusion patterns in cognitive skill diagnosis across verdict categories (Weak/Strong  $\times$  Yes/No and Vague/Evident  $\times$  Yes/No). Each heatmap shows the normalized distribution of predicted labels for a given model. Darker colors indicate higher proportions.

| Model                  | Recall | Recognize | Classify/Order | Compute | Determine | Represent/Model | Implement |
|------------------------|--------|-----------|----------------|---------|-----------|-----------------|-----------|
| GPT-5.2-Low            | 0.502  | 0.697     | 0.784          | 0.691   | 0.705     | 0.816           | 0.644     |
| GPT-5.2-High           | 0.549  | 0.717     | 0.814          | 0.698   | 0.690     | 0.829           | 0.648     |
| GPT-4o-Img             | 0.700  | 0.730     | 0.833          | 0.759   | 0.754     | 0.861           | 0.712     |
| GPT-4o                 | 0.589  | 0.610     | 0.804          | 0.716   | 0.683     | 0.829           | 0.675     |
| GPT-4o-mini            | 0.543  | 0.622     | 0.412          | 0.622   | 0.700     | 0.829           | 0.600     |
| Claude-Sonnet-4.6-Low  | 0.568  | 0.596     | 0.775          | 0.713   | 0.646     | 0.715           | 0.642     |
| Claude-Sonnet-4.6-High | 0.574  | 0.635     | 0.735          | 0.722   | 0.681     | 0.728           | 0.644     |
| Claude-3.5-Sonnet-Img  | 0.564  | 0.710     | 0.902          | 0.737   | 0.715     | 0.823           | 0.627     |
| Claude-3.5-Sonnet      | 0.534  | 0.674     | 0.804          | 0.688   | 0.692     | 0.810           | 0.602     |
| Gemini-3.1-Low         | 0.691  | 0.803     | 0.794          | 0.762   | 0.825     | 0.880           | 0.760     |
| Gemini-3.1-High        | 0.696  | 0.830     | 0.775          | 0.740   | 0.819     | 0.867           | 0.793     |
| Gemini-1.5-Flash-Img   | 0.615  | 0.670     | 0.784          | 0.722   | 0.754     | 0.810           | 0.696     |
| Gemini-1.5-pro         | 0.666  | 0.719     | 0.824          | 0.677   | 0.731     | 0.880           | 0.665     |
| Gemini-1.5-Flash       | 0.583  | 0.669     | 0.725          | 0.696   | 0.720     | 0.829           | 0.667     |
| Gemini-1.5-Flash-8b    | 0.366  | 0.566     | 0.343          | 0.610   | 0.659     | 0.797           | 0.545     |
| Deepseek-R1            | 0.743  | 0.805     | 0.873          | 0.753   | 0.804     | 0.873           | 0.712     |
| Deepseek-V3            | 0.649  | 0.730     | 0.775          | 0.734   | 0.720     | 0.848           | 0.673     |
| Llama-Maverick-Low     | 0.689  | 0.773     | 0.833          | 0.769   | 0.810     | 0.842           | 0.700     |
| Llama-Maverick-High    | 0.692  | 0.787     | 0.833          | 0.766   | 0.819     | 0.861           | 0.702     |
| Llama-3.1-405B         | 0.487  | 0.566     | 0.843          | 0.680   | 0.681     | 0.703           | 0.621     |
| Llama-3.3-70B          | 0.458  | 0.584     | 0.765          | 0.685   | 0.728     | 0.816           | 0.621     |
| Llama-3.1-8B-128K      | 0.592  | 0.678     | 0.873          | 0.658   | 0.707     | 0.791           | 0.542     |
| Qwen-2.5-72B           | 0.619  | 0.766     | 0.824          | 0.723   | 0.726     | 0.854           | 0.648     |
| Qwen-2.5-7B            | 0.634  | 0.700     | 0.637          | 0.729   | 0.787     | 0.829           | 0.636     |

Table 10: Accuracy by model (rows) and skill (columns)

| Model                  | Recall | Recognize | Classify/Order | Compute | Determine | Represent/Model | Implement |
|------------------------|--------|-----------|----------------|---------|-----------|-----------------|-----------|
| GPT-5.2-Low            | 0.218  | 0.341     | 0.500          | 0.587   | 0.366     | 1.000           | 0.556     |
| GPT-5.2-High           | 0.236  | 0.250     | 0.500          | 0.500   | 0.220     | 0.667           | 0.417     |
| GPT-4o-Img             | 0.436  | 0.477     | 1.000          | 0.848   | 0.732     | 1.000           | 0.833     |
| GPT-4o                 | 0.364  | 0.614     | 1.000          | 0.826   | 0.756     | 1.000           | 0.722     |
| GPT-4o-mini            | 0.455  | 0.682     | 1.000          | 0.848   | 1.000     | 1.000           | 0.944     |
| Claude-Sonnet-4.6-Low  | 0.255  | 0.477     | 1.000          | 0.761   | 0.317     | 0.667           | 0.694     |
| Claude-Sonnet-4.6-High | 0.164  | 0.500     | 0.000          | 0.674   | 0.220     | 0.667           | 0.611     |
| Claude-3.5-Sonnet-Img  | 0.527  | 0.773     | 1.000          | 0.717   | 0.439     | 1.000           | 0.639     |
| Claude-3.5-Sonnet      | 0.455  | 0.614     | 1.000          | 0.674   | 0.341     | 0.667           | 0.667     |
| Gemini-3.1-Low         | 0.255  | 0.273     | 0.500          | 0.739   | 0.561     | 1.000           | 0.611     |
| Gemini-3.1-High        | 0.255  | 0.295     | 0.500          | 0.674   | 0.537     | 1.000           | 0.528     |
| Gemini-1.5-Flash-Img   | 0.236  | 0.295     | 0.500          | 0.652   | 0.561     | 0.667           | 0.694     |
| Gemini-1.5-pro         | 0.309  | 0.409     | 1.000          | 0.804   | 0.512     | 1.000           | 0.556     |
| Gemini-1.5-Flash       | 0.236  | 0.409     | 1.000          | 0.696   | 0.585     | 1.000           | 0.611     |
| Gemini-1.5-Flash-8b    | 0.218  | 0.386     | 0.000          | 0.804   | 0.659     | 0.667           | 0.750     |
| Deepseek-R1            | 0.564  | 0.727     | 1.000          | 0.739   | 0.927     | 1.000           | 0.944     |
| Deepseek-V3            | 0.327  | 0.545     | 0.500          | 0.804   | 0.659     | 1.000           | 0.806     |
| Llama-Maverick-Low     | 0.291  | 0.364     | 1.000          | 0.717   | 0.610     | 0.667           | 0.861     |
| Llama-Maverick-High    | 0.309  | 0.364     | 1.000          | 0.739   | 0.683     | 1.000           | 0.722     |
| Llama-3.1-405B         | 0.200  | 0.568     | 1.000          | 0.587   | 0.366     | 1.000           | 0.639     |
| Llama-3.3-70B          | 0.218  | 0.295     | 1.000          | 0.630   | 0.488     | 0.333           | 0.611     |
| Llama-3.1-8B-128K      | 0.273  | 0.795     | 1.000          | 0.891   | 0.878     | 1.000           | 0.778     |
| Qwen-2.5-72B           | 0.236  | 0.432     | 1.000          | 0.848   | 0.659     | 1.000           | 0.722     |
| Qwen-2.5-7B            | 0.527  | 0.614     | 1.000          | 0.913   | 0.878     | 0.667           | 0.944     |

Table 11: Over-Attribution by model (rows) and skill (columns)

| Model                  | Recall | Recognize | Classify/Order | Compute | Determine | Represent/Model | Implement |
|------------------------|--------|-----------|----------------|---------|-----------|-----------------|-----------|
| GPT-5.2-Low            | 0.190  | 0.281     | 0.429          | 0.682   | 0.371     | 0.800           | 0.664     |
| GPT-5.2-High           | 0.277  | 0.286     | 0.333          | 0.633   | 0.293     | 1.000           | 0.621     |
| GPT-4o-Img             | 0.509  | 0.560     | 0.875          | 0.818   | 0.750     | 0.778           | 0.742     |
| GPT-4o                 | 0.470  | 0.567     | 0.500          | 0.822   | 0.677     | 0.636           | 0.769     |
| GPT-4o-mini            | 0.640  | 0.820     | 0.904          | 0.787   | 0.854     | 0.895           | 0.821     |
| Claude-Sonnet-4.6-Low  | 0.368  | 0.404     | 0.500          | 0.766   | 0.529     | 0.500           | 0.656     |
| Claude-Sonnet-4.6-High | 0.438  | 0.438     | 0.889          | 0.844   | 0.346     | 0.579           | 0.607     |
| Claude-3.5-Sonnet-Img  | 0.425  | 0.526     | 1.000          | 0.741   | 0.427     | 0.737           | 0.594     |
| Claude-3.5-Sonnet      | 0.424  | 0.414     | 0.364          | 0.743   | 0.429     | 0.833           | 0.717     |
| Gemini-3.1-Low         | 0.349  | 0.531     | 0.308          | 0.752   | 0.633     | 1.000           | 0.800     |
| Gemini-3.1-High        | 0.358  | 0.542     | 0.222          | 0.752   | 0.583     | 1.000           | 0.845     |
| Gemini-1.5-Flash-Img   | 0.191  | 0.274     | 0.500          | 0.726   | 0.541     | 0.522           | 0.688     |
| Gemini-1.5-pro         | 0.301  | 0.500     | 0.556          | 0.811   | 0.547     | 0.765           | 0.817     |
| Gemini-1.5-Flash       | 0.201  | 0.370     | 0.375          | 0.675   | 0.486     | 0.667           | 0.628     |
| Gemini-1.5-Flash-8b    | 0.218  | 0.378     | 0.604          | 0.813   | 0.562     | 0.769           | 0.703     |
| Deepseek-R1            | 0.737  | 0.897     | 0.727          | 0.846   | 0.865     | 1.000           | 0.948     |
| Deepseek-V3            | 0.368  | 0.447     | 0.556          | 0.805   | 0.476     | 0.875           | 0.728     |
| Llama-Maverick-Low     | 0.472  | 0.483     | 0.556          | 0.806   | 0.588     | 0.682           | 0.656     |
| Llama-Maverick-High    | 0.427  | 0.467     | 0.556          | 0.774   | 0.640     | 0.842           | 0.688     |
| Llama-3.1-405B         | 0.091  | 0.237     | 0.800          | 0.451   | 0.306     | 0.333           | 0.531     |
| Llama-3.3-70B          | 0.182  | 0.394     | 0.333          | 0.551   | 0.471     | 0.542           | 0.519     |
| Llama-3.1-8B-128K      | 0.345  | 0.563     | 0.727          | 0.750   | 0.625     | 0.967           | 0.690     |
| Qwen-2.5-72B           | 0.175  | 0.493     | 0.571          | 0.619   | 0.557     | 0.722           | 0.680     |
| Qwen-2.5-7B            | 0.674  | 0.702     | 0.889          | 0.867   | 0.824     | 0.926           | 0.883     |

Table 12: False-Attribution by model (rows) and skill (columns)

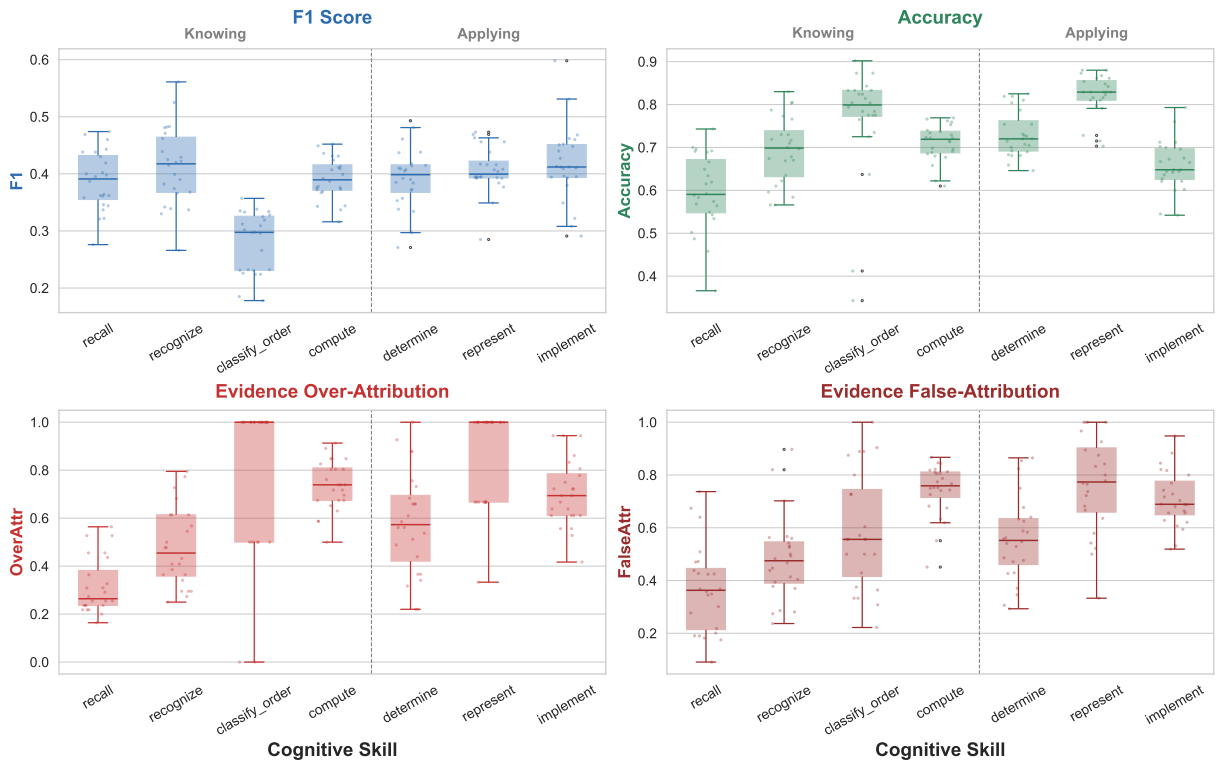


Figure 9: Performance and evidential attribution errors across cognitive skills, measured by F1 score, accuracy, evidence over-attribution, and false-attribution.

**[Compute]** After substituting the value of  $x$  into the second equation, were the four arithmetic operations performed correctly in calculating the value of  $a$  or  $m$ ?

**a**

$$3x + 5 = -x + 1$$

$$3x + x = 1 - 5$$

$$4x = -4$$

$$x = -1$$

$$\frac{-1 + 3}{2} = -2 + a$$

$$2 + 2 = a$$

**Teachers:**  
Evident No

**Claude-3-5-Sonnet:**  
Evident No **X**

**Claude-3-5-Sonnet-img:**  
Evident No **O**

**[Compute]** When calculating the prime factorization of a constant, was it performed accurately without arithmetic errors?

**b**

$$4\sqrt{3} \times a$$

} }  
21

**Teachers:**  
Evident No

**DeepSeek-R1:**  
Vague No **X**

**DeepSeek-V3:**  
Vague No **X**

**[Recognize]** Does the student notice that, to simplify the equation, both sides need to be divided by the coefficient of the highest order term, if necessary?

**c**

$$2x^2 + 4x - 5 = 0$$

$$\frac{2x^2 + 4x - 5}{2} = \frac{0}{2}$$

$$x^2 + 2x - \frac{5}{2} = 0$$

$$x^2 + 4x = \frac{5}{2} + 4$$

$$\left(\frac{x}{2}\right)^2 (x+2)^2 = \frac{13}{2}$$

$$x+2 = \pm \sqrt{\frac{13}{2}}$$

$$x = -2 \pm \sqrt{\frac{13}{2}}$$

**Teachers:**  
Evident Yes

**Gemini-1.5-Pro:**  
Evident No **X**

**GPT-4o:**  
Evident No **X**

Figure 10: Illustrative examples of diagnosis check items and student responses that LLMs failed to diagnose correctly. Evidence for human judgment is marked with a red box.