# MetaCLASS: Metacognitive Coaching for Learning with Adaptive Self-regulation Support

**Naiming Liu**
Rice University
nl35@rice.edu

**Richard Baraniuk**
Rice University
richb@rice.edu

**Shashank Sonkar**
University of Central Florida
shashank.sonkar@ucf.edu

## Abstract

Large language models can generate fluent explanations, but effective tutoring requires supporting the learner's *thought process*, not just delivering content. Metacognitive tutoring targets this gap by prompting *planning*, *monitoring*, *debugging*, and *evaluation*, and crucially, deciding when to be active versus minimally present, based on learner signals and trajectory. We introduce **MetaCLASS**, a learning-science grounded framework that formulates metacognitive tutoring as *move selection* over 11 interpretable actions aligned to self-regulated learning processes. MetaCLASS uses a two-phase framework that first plans a pedagogical trajectory conditioned on learner profiles (calibration, help-seeking) and then generates natural dialogue consistent with that plan. This yields a dataset of 1,015 conversations (7,711 turns) annotated with turn-level metacognitive labels, and validated for pedagogical contingency and trajectory adherence. We benchmark nine LLMs on predicting the next coach move given the problem and dialogue context. The best model achieves only 43.2% accuracy, and models exhibit *compulsive intervention bias*: in turns where effective metacognitive tutoring requires silent (41.7% of cases), models predict "no intervention" only 4.2% of the time, while severely over-predicting high-intervention moves. These results show that traditional content-based tutoring ability does not translate to metacognitive tutoring competence, positioning MetaCLASS as a testbed for developing intelligent tutors that promote self-regulated learning.

## 1 Introduction

While modern AI tutoring systems can generate accurate explanations across domains, they often fail at a core tutoring skill: scaffolding the learner's *thought process*. This deficiency is pedagogical, not conceptual. Effective tutoring extends beyond merely delivering correct content; it is the art of steering the learning process by clarifying goals, monitoring understanding, recovering from misconceptions, and reflecting to improve. Crucially, effective coaching is not defined solely by intervention. Some of the most powerful pedagogical moments are moments of *restraint*, which means allowing a learner to struggle productively rather than disrupting with a hint (Kapur, 2008). Human tutors use this strategy naturally by prompting students to articulate what they know, when they are not sure, and how they would verify a step (Chi et al., 2001). For instance, whereas a content tutor might simply state " use $F = mg$", a metacognitive coach asks "What are you trying to find?" or "What would convince you this step is correct?"

These teachable process skills fall under **metacognition** and **self-regulated learning (SRL)** (Flavell, 1979; Zimmerman, 2002). Within educational psychology, metacognitive regulation comprises four key components: **Planning**, where learners set goals and select strategies before acting; **Monitoring**, which involves tracking comprehension and surfacing uncertainty in real-time; **Debugging**, the process of repairing detected problems through alternative strategies or resources; and **Evaluation**, which entails reflecting on the progress and supporting transfer to new problems.

Most LLM-based tutoring systems implicitly optimizes for *helpfulness-as-output*: generating the next plausible hint or response under instructions (e.g., scaffolding, socratic questioning, mistake correction) (Macina et al., 2025), or selecting from high-level supportive behaviors (Sonkar et al., 2023). This focus obscures the core tutoring decision: determining *what should a tutor do right now, given the learner's signals and trajectory*. Consequently, progress is bottlenecked by the lack of (1) a shared, interpretable *action space* for metacognitive tutoring; (2) scalable supervision with validity checks for those actions; and (3) a benchmark that explicitly tests metacognitive tutoring as *decision-*

*making under pedagogy*.

To bridge this gap, MetaCLASS reframes "AI tutoring" from content delivery to *coaching the learning process*. We operationalize these metacognitive processes as an interpretable *action space* of coach moves (Table 1), treating *"No_intervention"* as a first-class action to formalize the role of restraints. Additionally, by formulating metacognitive tutoring into a problem of *pedagogical action selection*, we model the explicit decision of which metacognitive move to take next, conditioned on the learner's evolving state.

In this paper, we make **three key contributions**:

1. **MetaCLASS Framework.** We introduce the first learning-science grounded framework for LLM-based metacognitive tutoring designed to foster SRL. We operationalize Metacognitive Awareness Inventory (MAI) into **11 interpretable coach moves**, and systematically map them to LLMs' internal reasoning phases. The framework treats *No_intervention* as a first-class pedagogical action, formalizing the pedagogical value of restraint.

2. **MetaCLASS Dataset and Analysis.** We generate **1,015** metacognitive coaching conversations (**7,711** turns) across GSM8K, MATH, and AIME datasets, annotated with detailed turn-level coach moves and learner profiles (calibration, help-seeking). We also propose a **rigorous validation analysis**, demonstrating that the generated dialogues exhibit high pedagogical quality, such as contingency and trajectory adherence.

3. **Benchmark and Findings.** We introduce **Coach Move Prediction** task, where models must predict the optimal metacognitive move given the problem and dialogue context. Our evaluation of nine LLMs reveals a **compulsive intervention bias**, where models systematically over-produce high-intervention moves and severely under-predict *No_intervention*, This result highlights a fundamental gap between the ability to generate explanations and the ability to perform metacognitive coaching.

## 2 Related Work

### 2.1 Metacognition / Self-Regulated Learning

Metacognition encompasses learners' awareness and control over their own cognition processes (Flavell, 1979). Within the framework of self-regulated learning (SRL), learners set goals, se-

lect strategies, monitor understanding, repair breakdowns, and evaluate outcomes (Zimmerman, 2002). Consequently, effective metacognitive tutoring must support learners' *regulation* of problem-solving, not just deliver domain content. However, a key challenge is that high-level SRL constructs (e.g., "monitoring") are difficult to translate into the immediate, turn-level decisions required by a dialogue agent. MetaCLASS addresses this by grounding metacognitive coaching in three established frameworks: regulation processes as actionable interventions, help-seeking as a coached skill, and calibration as a signal for contingent support.

**Help-seeking as Coached Self-regulation.** Help-seeking is a critical metacognitive skill rather than just a mechanism for getting answers. Learners must learn to regulate their own resource use by deciding *when* help is needed, *what* specifically to ask for, and *how* to apply the assistance to their current knowledge state (Aleven et al., 2006a; Roll et al., 2011a).

**Calibration and Confidence Signals.** Calibration, defined as the accuracy of a learner's self-assessment (Kruger and Dunning, 1999), provides observable signals for coaching. Over-confident learners may miss errors, while under-confident learners may abandon correct reasoning. Well-calibrated learners localize their uncertainty (e.g., "I understand X but I'm stuck on Y"), enabling the tutor to provide targeted monitoring and debugging interventions.

**Restraint as a Pedagogical Decision.** Metacognitive support does not require constant intervention. Research on human tutoring highlights the importance of selective prompting and timing (Chi et al., 2001), and productive failure research shows that initial struggle benefits deep learning (Kapur, 2008). Accordingly, MetaCLASS formalizes *"No_intervention"* as a meaningful pedagogical action and an explicit benchmark label.

### 2.2 Intelligent Tutors

Classic Intelligent Tutoring Systems (ITS) (Graesser et al., 2004; Anderson et al., 1995) demonstrated that adaptive feedback improves learning outcomes. However, adaptation in these systems is typically driven by correctness tracing rather than explicit metacognitive regulation. Recent LLM-based tutors excel at explanation generation (Macina et al., 2025; Sonkar et al.,
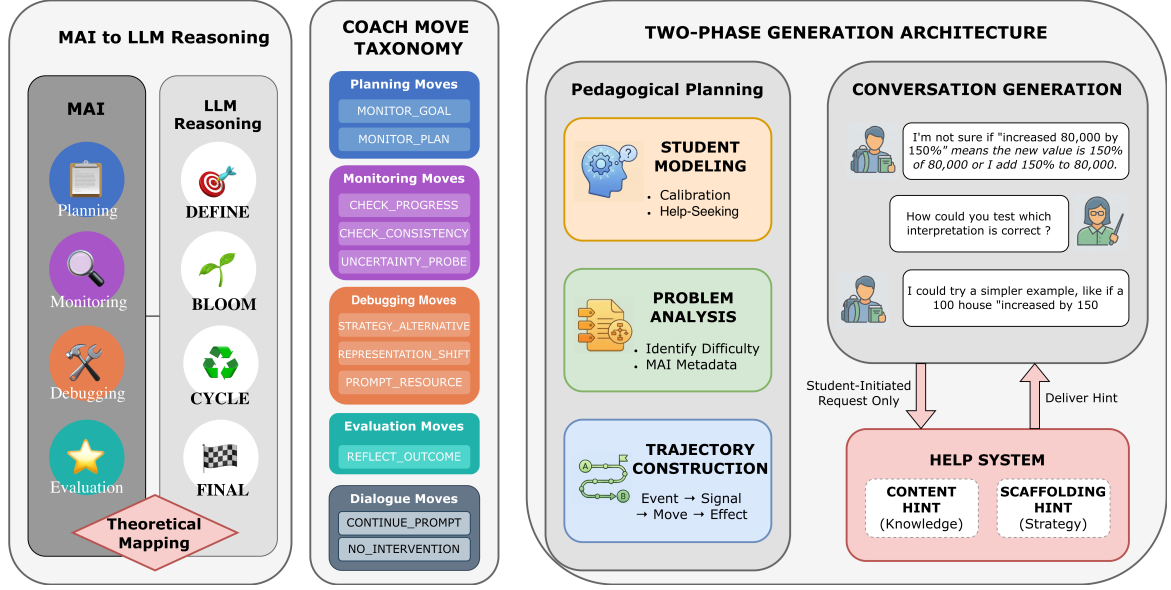
Figure 1: An overview of the MetaCLASS Framework. **(Left)** We ground the framework in learning science by mapping MAI factors to LLM reasoning phases. **(Center)** We define a taxonomy of 11 interpretable coach moves, categorized by regulatory process. **(Right)** A two-phase generation architecture first plans a pedagogical trajectory together with learner profiles and problem analysis, which then guides the generation of metacognitive coaching dialogues.

2023), which often makes them prioritize *utterance quality* (fluency, helpfulness) over *pedagogical decision-making*. MetaCLASS fills this gap by treating metacognition as a decision process. By selecting the optimal coach move (including strategic silence) conditioned on the learner's state, the framework explicitly scaffolds self-regulation rather than merely correcting errors.

## 3 MetaCLASS Framework

In this section, we first ground the framework in learning science by aligning MAI theory with established LLM reasoning patterns, deriving 11 interpretable coach moves. We then detail the two-phase architecture that utilizes these moves to construct contingent, metacognitive coaching dialogues (Figure 1).

### 3.1 Connecting MAI to LLMs Reasoning

#### 3.1.1 Metacognitive Awareness Inventory

MetaCLASS framework is theoretically grounded in the Metacognitive Awareness Inventory (MAI) (Schraw and Dennison, 1994), which characterizes how learners control their cognitive process during learning through four empirically validated factors.

- **Planning** includes pre-task activities, such as goal setting, strategy selection, and resource allocation (e.g., *I set specific goals before I begin a task*).

- **Monitoring** captures the real-time tracking of comprehension during learning, enabling the detection of confusion, impasses, or inconsistencies (e.g., *I find myself pausing regularly to check my comprehension*).

- **Debugging** encompasses error correction strategies triggered when Monitoring detects a breakdown (e.g., *I re-evaluate my assumptions when I get confused*).

- **Evaluation** involves post-task reflection on performance and strategy effectiveness, supporting transfer to future tasks (e.g., *I ask myself how well I accomplished my goals once I'm finished*).

These four factors form a validated taxonomy of metacognitive regulation, providing MetaCLASS with a principled basis for designing coach moves.

#### 3.1.2 LLMs Reasoning Structure

Recent work (Marjanović et al., 2025) analyze reasoning chains from DeepSeek-R1 (Guo et al., 2025) and identified four consistent reasoning structure: **DEFINE** (reformulating the problem and identifying goals), **BLOOM** (decomposing and executing initial steps), **CYCLE** (reconsidering assumptions and detecting errors through self-monitoring), and **FINAL** (committing to an answer with confidence assessment).

| MAI | Move | Function | Example Response |
|---|---|---|---|
| **Planning** | MONITOR_GOAL (MG)<br>MONITOR_PLAN (MP) | Elicit goal awareness<br>Elicit strategy awareness | "What are you trying to find?"<br>"What's your approach?" |
| **Monitoring** | CHECK_PROGRESS (CP)<br>CHECK_CONSISTENCY (CC)<br>UNCERTAINTY_PROBE (UP) | Prompt progress tracking<br>Surface contradictions<br>Localize confusion | "How's it going?"<br>"Does that fit with what you said?"<br>"What's making you hesitate?" |
| **Debugging** | STRATEGY_ALTERNATIVE (SA)<br>REPRESENTATION_SHIFT (RS)<br>PROMPT_RESOURCE (PR) | Invite different approach<br>Suggest re-framing<br>Prompt help-seeking awareness | "What else could you try?"<br>"Would a diagram help?"<br>"What would help you move forward?" |
| **Evaluation** | REFLECT_OUTCOME (RO) | Prompt retrospective reflection | "What worked here?" |
| **Dialogue** | CONTINUE_PROMPT (CT)<br>NO_INTERVENTION (NI) | Neutral nudge<br>Preserve productive flow | "Keep going."<br>(silence) |

Table 1: MetaCLASS Coach Move Taxonomy: a mapping of 11 interpretable coach moves to MAI regulatory processes. These moves are designed to prompt metacognitive reflection rather than deliver domain content.

### 3.1.3 Connecting LLM Reasoning to MAI

MetaCLASS's central insight is recognizing that the empirically observed reasoning structures of LLMs maps systematically onto the regulatory factors of the MAI:

- **DEFINE → Planning**: Both involve establishing specific goals prior to execution.
- **BLOOM → Planning**: Both entail problem decomposition and strategy selection.
- **CYCLE → Monitoring + Debugging**: Both involve self-monitoring and the adjustment of strategies when stuck.
- **FINAL → Evaluation**: Both require performance assessment and reflection on the final outcome.

This mapping has three key implications. First, it grounds LLM reasoning patterns in validated learning science. Second, it enables principled coach design, where each MAI factor prescribes interventions for the corresponding reasoning phase. Third, it positions MetaCLASS coach moves as *external scaffolding* for LLMs' internal metacognitive processes. If models naturally execute these regulatory processes when solving problems, effective tutoring means prompting learners to follow the same process.

### 3.2 Metacognitive Coach Moves

We operationalize MAI's four regulatory factors into eleven coach moves (Table 1). Each move translates an internal metacognitive process into an external pedagogical prompt. Critically, these moves are designed to elicit reflection on the problem-solving process itself, which should *never* provide domain knowledge, evaluate correctness, or reveal solutions. A special coach move is PROMPT_RESOURCE (PR), derived from help-seeking research (Roll et al., 2011a), which handles hints delivery, but only when students initiate requests. Additionally, we include two dialogue moves support natural flow: CONTINUE_PROMPT (CT)[1] provides neutral nudges (e.g., *"Keep going"*), and NO_INTERVENTION (NI) facilitates productive struggle through silence.

### 3.3 Two-Phase Generation Architecture

MetaCLASS uses a two-phase architecture that conducts pedagogical planning prior to the generation of the coaching conversation.

#### 3.3.1 Phase 1: Pedagogical Planning

Before generating the coaching dialogue, we constructs a complete pedagogical plan, ensuring that every coach move is theoretically grounded in the student's demonstrated state and trajectory.

**Student Modeling** Effective coaching requires adapting to individual differences. MetaCLASS models students along two orthogonal dimensions: Calibration, Help-Seeking.

**Calibration** measures the accuracy of self-assessment, specifically whether students' perceived understanding matches their actual understanding (Kruger and Dunning, 1999).

- **Over-Confident** students overestimate their understanding, using confident language (e.g., *"obviously," "just"*) and dismissing challenges without detecting errors.

---

[1] CONTINUE_PROMPT is abbreviated as CT to avoid confusion with CHECK_PROGRESS (CP).

- **Under-Confident** students doubt correct reasoning, hedging constantly (e.g., *"maybe," "I'm not sure"*) and hesitating to commit to answers.
- **Well-Calibrated** students accurately assess understanding, articulating specific confusions (e.g., *"I understand X but I'm stuck on Y"*).

**Help-Seeking**   captures patterns of resource utilization during learning (Aleven et al., 2006a; Roll et al., 2011a).

- **Avoidant** students resist asking for help even when stuck, viewing help requests as failure and persist unproductively. They are primary targets for resource prompts.
- **Executive** students use resources strategically, attempting problems first and integrating help into understanding. They exhibit the target behavior for self-regulated learning.
- **Dependent** students request help prematurely and broadly, seeking answers rather than understanding. They should not receive immediate resource prompts.

Crossing these dimensions yields nine theoretical profiles, of which **eight are valid** (as the combination of Over-confident and Executive is contradictory). Each valid profile demands different coaching strategies.

**Problem Analysis**   Each problem is analyzed to identify potential learning obstacles. **Knowledge Gaps** (missing domain knowledge) are addressed through content hints; **Strategy Gaps** (knowledge present but cannot apply) are addressed through scaffolding hints; **Monitoring Gaps** (unnoticed errors) and **Execution Gaps** (careless mistakes) are addressed through prompts that highlight inconsistency or encourage verification without revealing the solution. Both content and scaffolding hints are pre-generated as problem-specific resources during planning, but delivered only when students explicitly request them.

Additionally, we generate metacognitive support structures aligned with the MAI framework for each problem. This includes *planning support* (clarifying goals and strategies), *monitoring support* (highlighting potential ambiguities and verification checkpoints), *debugging support* (introducing simpler instances and common errors), and *evaluation support* (providing metacognitive insights). This structured decomposition serves as the pedagogical foundation for the trajectory construction.

**Trajectory Construction**   The core of planning is the trajectory, a sequence of pedagogical events that serves as the causal backbone of the conversation. We structure each event as a four-part chain:

- **Event**: What happens in the problem-solving process (e.g., encountering an ambiguity, getting stuck on a concept)
- **Signal**: What observable behavior the student produces (e.g., hedging language, long pause)
- **Move**: What coach intervention is appropriate given the student's profile (e.g., uncertainty probe, consistency check)
- **Effect**: What change should result from the intervention (e.g., articulating confusion, requesting for help)

### 3.3.2   Phase 2: Dialogue Generation

This phase transforms the established LLM reasoning and trajectories into a coherent student-coach dialogues. The generation process strictly enforces the pedagogical plan: students exhibit behaviors consistent with their assigned profiles, encounter anticipated learning obstacles, and progress through the trajectory's event sequence. While the interaction logic is structurally fixed, the specific phrasing and conversational flow are dynamically generated by LLMs, allowing the interaction to sound natural rather than scripted.

### 3.3.3   Help System Design

The help system models student-initiated resource utilization. In contrast to traditional tutoring systems where hints are offered by the tutor, Meta-CLASS generates hints during planning phase but delivers them only when students explicitly request assistance.

**Hint Types.** We distinguish between two categories: *Content hints* provide domain knowledge (e.g., formulas, facts) to address Knowledge Gaps, while *scaffolding hints* offer procedural guidance (e.g., strategies, approaches) to resolve Strategy Gaps.

**Help Quality.** We evaluate the pedagogical appropriateness of help requests: *appropriate* (asked when genuinely stuck, chose correct type), *premature* (asked prior to a valid attempt), *delayed* (asked after a period of excessive unproductive struggle), or *mismatched* (asked for wrong hint type).

| Dataset | # Conv | # Turn | Contingent | Trajectory |
|---|---|---|---|---|
| GSM8K | 496 | 3,180 | 98.6% | 99.4% |
| MATH | 490 | 4,142 | 96.2% | 100% |
| AIME | 29 | 389 | 97.3% | 100% |
| **Total / Avg** | 1,015 | 7,711 | 97.3% | 99.8% |

Table 2: MetaCLASS dataset Overview with evaluation results on Contingent Scaffolding and Trajectory Following.

# 4 MetaCLASS Dataset

## 4.1 Dataset Statistics

Following MetaCLASS framework, we generate student-coach conversations with GPT-5.1 (Achiam et al., 2023) model, across three math datasets of varying difficulty: GSM8K (grade school math word problems) (Cobbe et al., 2021), MATH (high school competition) (Hendrycks et al., 2021), and AIME (advanced olympiad questions). As shown in Table 2, the dataset comprises 1,015 conversations totaling 7,711 turns. An example of generated conversation can be found in Appendix C.

## 4.2 Dataset Validation

**Contingent Scaffolding** We first examine contingent scaffolding (Wood et al., 1976). This analysis identifies student difficulty signals (e.g., "I'm not sure," "I'm stuck,", etc.) and evaluated whether coach responded moves based on conversation phase. Of the 828 conversations that show such signals, coaches select pedagogically appropriate moves in 97.3% of cases (Table 2), demonstrating that interventions adapt to learner state rather than following predetermined sequences.

**Trajectory Alignment.** We analyze trajectory-conversation alignment to evaluate whether the LLM executes its planned trajectory. Table 2 demonstrate 99.8% plan-execution alignment across all datasets, providing evidence that MetaCLASS-generated conversations exhibit coherent, intentional pedagogy according to the planned trajectory.

**MAI Factor Coverage** We measure MAI factor coverage of whether each MAI factor appears at least once per conversation. Coverage rates are planning (80.9%), monitoring (97.5%), debugging (37.4%), and evaluation (100%). The moderate debugging rate is appropriate since debugging requires errors, which not all conversations contain.

| Model | GSM8K | | MATH | | AIME | |
|---|---|---|---|---|---|---|
| | Full | Min | Full | Min | Full | Min |
| *Large Models* | | | | | | |
| GPT-oss-120b | 36.2 | **32.9** | 30.8 | 27.1 | 25.2 | 19.3 |
| Qwen3-80b-Think | **43.2** | 32.3 | **37.5** | **27.6** | **29.6** | 20.3 |
| Qwen3-80b-Inst | 36.2 | 29.7 | 30.4 | 26.4 | 22.1 | **24.7** |
| LLaMA-70b | 36.0 | 29.3 | 30.1 | 23.0 | 22.4 | 18.3 |
| *Small Models* | | | | | | |
| GPT-oss-20b | **34.5** | **29.8** | **29.3** | **23.8** | 21.9 | **17.5** |
| Phi4-14b | 33.7 | 25.1 | 29.1 | 21.3 | 21.1 | 18.8 |
| LLaMA-8b | 28.5 | 20.8 | 25.9 | 17.9 | 15.7 | 14.4 |
| Qwen3-4b | 34.4 | 5.1 | 28.4 | 5.9 | **22.6** | 6.4 |
| Phi4-4b | 27.0 | 19.8 | 22.3 | 16.8 | 13.1 | 9.8 |

Table 3: Model performance on Coach Move Prediction across three mathematical reasoning datasets, comparing Full versus Minimal prompting strategies. The best accuracy in each column is in highlighted **bold**.

Critically, debugging exhibits profile differentiation: AVOIDANT students engage in debugging 81.2% of the time compared to only 9.5% for EXECUTIVE and 12.4% for DEPENDENT students. This pattern is theoretically grounded because AVOIDANT students struggle longer before seeking help and thus encounter more errors.

# 5 Coach Move Prediction Benchmark

## 5.1 Task Definition

We define **Coach Move Prediction** as inference with a prompted, fixed LLM: given a math problem and the tutoring dialogue up to the current student turn, the model must predict the next metacognitive coaching action. Let $x = (P, H_t)$ denote the input, where $P$ is the math problem and $H_t = \{(u_1, m_1, r_1), \ldots, (u_{t-1}, m_{t-1}, r_{t-1}), u_t\}$ is the dialogue history up to the current student utterance $u_t$. Here, $u_i$ is the student utterance, $m_i$ is the coach move label, and $r_i$ is the coach response text at turn $i$.

The model generates a *structured output* containing: (i) a discrete coach move $m_t$ and (ii) an optional natural-language coach response $r_t$. We evaluate only the move selection: the predicted move $\hat{m}_t$ is extracted by parsing the model output and must belong to a finite action space $\mathcal{M}$ of 11 metacognitive coach moves (Table 1).

Conceptually, the prompting and decoding strategy induce an implicit probability distribution over moves, and we decode a single prediction:

$$\hat{m}_t = \arg\max_{m \in \mathcal{M}} P_\theta(m \mid x), \qquad (1)$$

where $P_\theta$ represents the distribution of the fixed pretrained model under our prompt template and decoding procedure. We report accuracy of $\hat{m}_t$ against the ground-truth move $m_t$ specified by the MetaCLASS pedagogical trajectory. If the model outputs an invalid or unparsable label (not in $\mathcal{M}$), we count the prediction as incorrect.

## 5.2 Models and Evaluations

**Models.** We evaluate nine open-source LLMs spanning four families. Large models include GPT-oss-120B (Agarwal et al., 2025), Qwen3-80B-Instruct (Yang et al., 2025), and LLaMA-3.3-70B (Touvron et al., 2023). Small models include GPT-oss-20B, Qwen3-4B, Phi-4-14B (Abdin et al., 2024), Phi-4-4B, and Llama-3.1-8B. We use hyperparameters following model recommendations (see Appendix B).

**Prompting Strategy.** We evaluate model performance under two prompt conditions. The *Full prompt* provides (1) definitions and examples for all 11 coach moves; (2) move selection guidelines specifying when each is appropriate; and (3) coaching rules emphasizing metacognitive support over direct content delivery. In contrast, the *Minimal prompt* provides only move names without definitions, examples, or guidelines.

**Evaluation Metrics.** Our primary evaluation is *Accuracy*, reported both at the aggregate level and on a per-move basis. We also report the *NI Detection rate* to assess the model's ability to support students with appropriate silent restraint (e.g., knowing when not to intervene).

## 6 Results and Discussion

Our analysis reveals three fundamental failures distinguishing metacognitive coaching from standard content tutoring: (1) models generally struggle with metacognitive move prediction; (2) models exhibit a *compulsive intervention bias* that explicit prompt instruction cannot override; and (3) models tend to favor consistency checks over strategic redirection.

### 6.1 Performance Across Models and Prompts

Table 3 shows overall accuracy for coach move prediction. Model performance plateaus at 43.2%, revealing that open-source LLMs continue to struggle with identifying the appropriate timing and

method for metacognitive intervention. This contrasts with their documented proficiency on content-based tutoring (Scarlatos et al., 2025; Sonkar et al., 2023), showing that knowing how to explain does not translate to knowing when to intervene. Furthermore, while larger models generally achieve higher accuracy, smaller models are inconsistent. For instance, Qwen3-4b occasionally matches larger models yet exhibits catastrophic drops under certain configurations (e.g., falling to 5.1% on GSM8K with Minimal prompt).

**Impact of Context (Full vs Minimal).** Full prompts with explicit MAI definitions consistently outperform Minimal prompts. Qwen3-4b shows this dependency most starkly, dropping over 20% in the absence of explicit instruction. This result reveals that LLMs lack inherent metacognitive tutoring ability and are unable to autonomously deploy effecting coaching strategies without detailed instructions.

### 6.2 Move-level Analysis

**Accuracy Variation.** Accuracy varies dramatically across moves (Table 4): near-perfect on RO ($\sim$99%) versus near-zero on NI ($\leq$29.0%) and CT ($\leq$15.4%). High-accuracy moves follow clear surface patterns (RO at completion, UP after hesitation, MP after problem statements). Mid-accuracy moves (CC $\leq$59.5%, MP $\leq$68%) show high model-to-model variance, suggesting some LLMs capture these relevant patterns.

**The Compulsive Intervention Bias** The most critical finding is models' systematic inability to withhold intervention. As shown in Table 5, NI, the skill of recognizing when to remain silent, proves to be the most difficult, with accuracy frequently below 5% for instruction-tuned models despite being the most common ground truth label (41.7%). This occurs even though prompts explicitly specify *"NI should be $\sim$35-50% of coach turns."* Models exhibit average NI accuracy of 4.2% with massive negative prediction bias of $-37.5\%$, dramatically under-predicting silence. This pattern suggests instruction-tuning of LLMs creates deeply ingrained intervention patterns that prompting cannot override. Even thinking-enabled models, which achieve 5$\times$ better NI detection than instruction-tuned counterparts (29.0% vs. 5.7%), still intervene in over 70% of cases where restraint is appropriate. Conversely, models over-predict high-intervention moves like CP with significant

| Model | MG | MP | CP | CC | UP | SA | RS | PR | RO | CT | NI | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *Large Models* | | | | | | | | | | | | |
| gpt-oss-120b | 23.0 | 60.5 | 33.1 | 31.7 | 69.6 | 23.8 | 33.0 | 9.0 | 98.6 | 15.4 | 4.5 | 32.7 |
| Qwen3-80b-Think | 10.1 | 54.4 | 5.4 | 33.6 | 67.2 | 12.6 | 7.2 | 8.7 | 98.1 | 0.0 | 29.0 | 39.5 |
| Qwen3-80b-Inst | 51.8 | 38.1 | 19.2 | 33.0 | 71.3 | 12.1 | 8.2 | 37.5 | 99.6 | 0.0 | 5.7 | 32.4 |
| LLaMA-70b | 11.3 | 63.4 | 19.7 | 16.7 | 71.5 | 34.4 | 7.2 | 23.3 | 99.7 | 7.7 | 6.9 | 32.1 |
| *Small Models* | | | | | | | | | | | | |
| gpt-oss-20b | 14.7 | 59.6 | 17.6 | 23.6 | 80.8 | 19.1 | 6.2 | 11.8 | 98.0 | 0.0 | 5.0 | 31.0 |
| Phi4-14b | 32.8 | 55.3 | 15.3 | 46.5 | 36.4 | 38.8 | 1.0 | 24.3 | 99.2 | 7.7 | 1.4 | 30.6 |
| LLaMA-8b | 18.1 | 46.8 | 23.2 | 28.6 | 64.9 | 22.1 | 3.1 | 15.3 | 79.9 | 7.7 | 2.0 | 26.4 |
| Qwen3-4b | 66.0 | 25.2 | 4.0 | 45.1 | 72.5 | 25.9 | 15.5 | 35.8 | 98.4 | 0.0 | 0.3 | 30.6 |
| Phi4-4b | 67.5 | 22.9 | 18.8 | 30.6 | 30.9 | 1.2 | 0.0 | 6.2 | 93.6 | 0.0 | 0.1 | 23.8 |

Table 4: Per-move prediction accuracy (%) for the 11 coach moves (Full prompt and aggregate across all datasets).
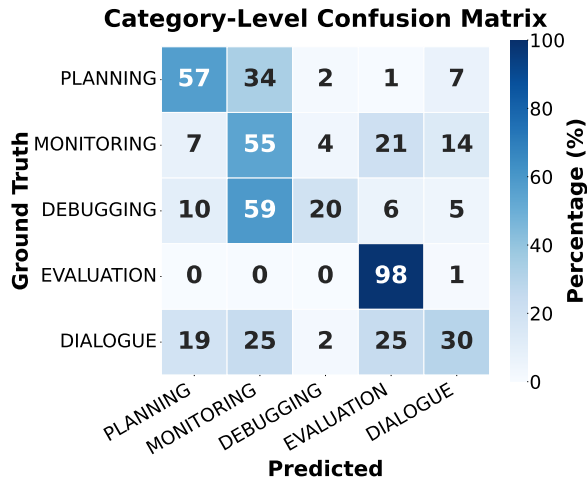


Figure 2: Category-level confusion matrix for Qwen3-80b-Think. DEBUGGING moves are predominantly misclassified as MONITORING (59%), suggesting the model struggles to distinguish exploratory probing from corrective intervention.
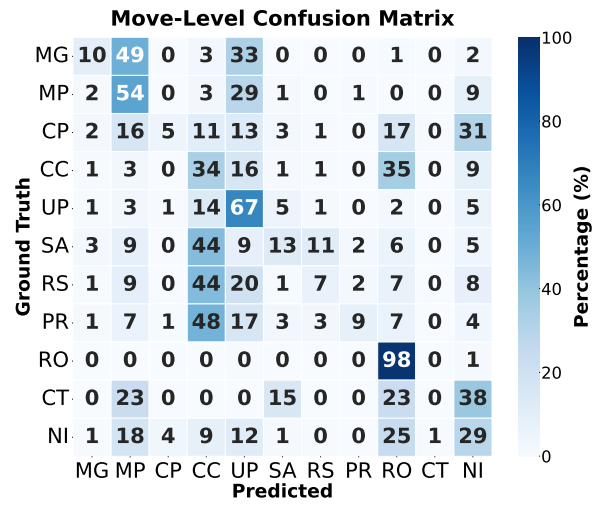


Figure 3: Move-level confusion matrix for Qwen3-80b-Think. NI reaches only 29% accuracy, despite comprising 42% of ground truth, revealing compulsive intervention bias.

positive bias (+19.3%) and near-perfect accuracy (99%).

### 6.3 Failure Mode Analysis

To understand systematic misprediction patterns, we analyzed confusion matrices for Qwen3-80b-Think in Figures 2 and 3.

**Planning-Monitoring conflation.** MG is frequently confused with MP (49%), with both exhibiting high cross-predictions with UP (33% and 29% respectively). This suggests models cannot distinguish goal-setting from strategy-setting or from localizing confusions.

**Strategy redirection misdiagnosed as verification.** SA and RS show low accuracy with 44% confusion with CC. When students need alterna-

tive approaches, models favor consistency checks over strategic redirection. Figure 2 confirms this: 59% of Debugging moves are mispredicted as Monitoring, showing models recognize stuck students but systematically choose verification over strategy change.

**Resource prompting nearly invisible.** PR achieves only 9% accuracy, most often substituted with CC (48%). Models cannot distinguish students needing verification support from those requiring help-seeking scaffolding, leaving help-avoidant students without appropriate intervention.

### 6.4 Implications for AI Tutoring Systems.

These findings have critical implications for LLM-based tutoring systems. Current models produce 8-10 times more interventions than appropriate,

| Rank | Move | Acc (%) | N | GT % | Pred % | Bias |
|------|------|---------|------|------|--------|------|
| 1 | NI | 7.4 | 22,533 | 41.7 | 4.2 | -37.5 |
| 2 | CT | 7.7 | 91 | 0.2 | 3.6 | +3.4 |
| 3 | RS | 13.8 | 679 | 1.3 | 2.2 | +0.9 |
| 4 | PR | 15.8 | 2,016 | 3.7 | 2.3 | -1.4 |
| 5 | SA | 20.5 | 2,380 | 4.4 | 2.9 | -1.5 |
| 6 | MG | 20.6 | 2,282 | 4.2 | 3.7 | -0.5 |
| 7 | CC | 27.4 | 5,586 | 10.3 | 10.4 | +0.1 |
| 8 | CP | 29.3 | 2,982 | 5.5 | 11.2 | +5.7 |
| 9 | MP | 47.9 | 4,606 | 8.5 | 12.9 | +4.4 |
| 10 | UP | 58.8 | 3,710 | 6.9 | 12.8 | +5.9 |
| 11 | RO | 99.0 | 7,112 | 13.2 | 32.5 | +19.3 |

Table 5: Move-level accuracy and prediction bias (large models with full prompt for all datasets). Ranked by accuracy with N showing ground truth instances. Bias (Pred % - GT %) reveals under-prediction of restraint (NI: -37.5) and over-prediction of high-intervention moves (RO: +19.3).

training students to wait for hints rather than develop transferable metacognitive skills like planning, monitoring, and debugging. Additionally, this compulsive intervention bias is invisible in standard evaluation. Metrics assessing helpfulness or satisfaction cannot detect over-intervention and students may even prefer excessive scaffolding despite its harm to long-term development. Without explicit ground truth for restraint, as MetaCLASS provides, this pedagogical failure remains hidden behind metrics that reward responsiveness over developmental appropriateness.

## 7 Conclusion

We presented **MetaCLASS**, a learning-science grounded framework and benchmark for *metacognitive tutoring*. MetaCLASS operationalizes metacognitive regulation into 11 interpretable coach moves spanning planning, monitoring, debugging, and evaluation. To provide scalable, auditable supervisio, we utilize a two-phase generation procedure first plans pedagogical trajectories conditioned on learner profiles, then generates dialogue consistent with that plan, yielding a dataset with 7,711 conversation turns and explicit turn-level metacognitive labels and validation analyses testing contingency and trajectory following.

Our Coach Move Prediction benchmark reveals that current LLMs struggle with pedagogical decision-making: the best model achieves only 43.2% accuracy, and models consistently exhibit compulsive intervention bias by over-predicting high-intervention moves while severely under-

predicting restraint (no intervention). These results suggest standard instruction-following tuning does not produce metacognitive tutoring competence. MetaCLASS provides a concrete testbed for training and evaluating decision-level metacognitive support.

## Limitations

MetaCLASS is designed as a *diagnostic benchmark* for metacognitive tutoring decisions. We operationalize metacognitive support as move selection in student-coach dialogues within math problem-solving settings (GSM8K, MATH, AIME). This focus prioritizes interpretability and controlled evaluation over broad coverage of domains and modalities. In addition, our benchmark uses a single target coach move prediction derived from an explicit pedagogical trajectory. While this design enables consistent supervision, certain conversational contexts may admit multiple pedagogically reasonable moves that are not credited under strict matching. Accordingly, our results should be interpreted specifically as evidence of current LLMs' ability to follow theory-grounded metacognitive intervention policies under a controlled setting.

## Acknowledgments

## References

Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J Hewett, Mojan Javaheripi, Piero Kauffmann, and 1 others. 2024. Phi-4 technical report. *arXiv preprint arXiv:2412.08905*.

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

Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Arbus, Rahul K Arora, Yu Bai, Bowen Baker, Haiming Bao, and 1 others. 2025. gpt-oss-120b & gpt-oss-20b model card. *arXiv preprint arXiv:2508.10925*.

Vincent Aleven, Bruce Mclaren, Ido Roll, and Kenneth Koedinger. 2006a. Toward meta-cognitive tutoring: A model of help seeking with a cognitive tutor. *International Journal of Artificial Intelligence in Education*, 16(2):101–128.

Vincent Aleven, Bruce Mclaren, Ido Roll, and Kenneth Koedinger. 2006b. Toward meta-cognitive tutoring: A model of help seeking with a cognitive tutor. *International Journal of Artificial Intelligence in Education*, 16(2):101–128.

Vincent Aleven, Ido Roll, Bruce M McLaren, and Kenneth R Koedinger. 2016. Help helps, but only so much: Research on help seeking with intelligent tutoring systems. *International Journal of Artificial Intelligence in Education*, 26(1):205–223.

John R Anderson, Albert T Corbett, Kenneth R Koedinger, and Ray Pelletier. 1995. Cognitive tutors: Lessons learned. *The journal of the learning sciences*, 4(2):167–207.

Michelene TH Chi, Stephanie A Siler, Heisawn Jeong, Takashi Yamauchi, and Robert G Hausmann. 2001. Learning from human tutoring. *Cognitive science*, 25(4):471–533.

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, and 1 others. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.

John H Flavell. 1979. Metacognition and cognitive monitoring: A new area of cognitive–developmental inquiry. *American psychologist*, 34(10):906.

Arthur C Graesser, Shulan Lu, George Tanner Jackson, Heather Hite Mitchell, Mathew Ventura, Andrew Olney, and Max M Louwerse. 2004. Autotutor: A tutor with dialogue in natural language. *Behavior Research Methods, Instruments, & Computers*, 36(2):180–192.

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

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*.

Manu Kapur. 2008. Productive failure. *Cognition and instruction*, 26(3):379–424.

Kenneth R Koedinger and Vincent Aleven. 2007. Exploring the assistance dilemma in experiments with cognitive tutors. *Educational psychology review*, 19(3):239–264.

Justin Kruger and David Dunning. 1999. Unskilled and unaware of it: how difficulties in recognizing one's own incompetence lead to inflated self-assessments. *Journal of personality and social psychology*, 77(6):1121.

Jakub Macina, Nico Daheim, Ido Hakimi, Manu Kapur, Iryna Gurevych, and Mrinmaya Sachan. 2025. Mathtutorbench: A benchmark for measuring open-ended pedagogical capabilities of llm tutors. *arXiv preprint arXiv:2502.18940*.

Sara Vera Marjanović, Arkil Patel, Vaibhav Adlakha, Milad Aghajohari, Parishad BehnamGhader, Mehar Bhatia, Aditi Khandelwal, Austin Kraft, Benno Krojer, Xing Han Lù, and 1 others. 2025. Deepseek-r1 thoughtology: Let's think about llm reasoning. *arXiv preprint arXiv:2504.07128*.

Alexander Renkl, Robert K Atkinson, and Cornelia S Große. 2004. How fading worked solution steps works–a cognitive load perspective. *Instructional science*, 32(1):59–82.

Ido Roll, Vincent Aleven, Bruce M McLaren, and Kenneth R Koedinger. 2011a. Improving students' help-seeking skills using metacognitive feedback in an intelligent tutoring system. *Learning and instruction*, 21(2):267–280.

Ido Roll, Vincent Aleven, Bruce M McLaren, and Kenneth R Koedinger. 2011b. Improving students' help-seeking skills using metacognitive feedback in an intelligent tutoring system. *Learning and instruction*, 21(2):267–280.

Alexander Scarlatos, Naiming Liu, Jaewook Lee, Richard Baraniuk, and Andrew Lan. 2025. Training llm-based tutors to improve student learning outcomes in dialogues. In *International Conference on Artificial Intelligence in Education*, pages 251–266. Springer.

Gregory Schraw and Rayne Sperling Dennison. 1994. Assessing metacognitive awareness. *Contemporary educational psychology*, 19(4):460–475.

Daniel L Schwartz and John D Bransford. 1998. A time for telling. *Cognition and Instruction*, 16(4):475–5223.

Shashank Sonkar, Naiming Liu, Debshila Mallick, and Richard Baraniuk. 2023. Class: A design framework for building intelligent tutoring systems based on learning science principles. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1941–1961.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, and 1 others. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Janneke Van de Pol, Monique Volman, and Jos Beishuizen. 2010. Scaffolding in teacher–student interaction: A decade of research. *Educational psychology review*, 22(3):271–296.

David Wood, Jerome S Bruner, and Gail Ross. 1976. The role of tutoring in problem solving. *Journal of child psychology and psychiatry*, 17(2):89–100.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*.

Barry J Zimmerman. 2002. Becoming a self-regulated learner: An overview. *Theory into practice*, 41(2):64–70.

## A  MetaCLASS Design Principles Grounded in Learning Science

MetaCLASS translates findings from tutoring and SRL into *generation constraints* and evaluation targets. These principles govern *when* support should be offered, *how much* to offer, and when the best move is to withhold help.

**Contingent scaffolding (adaptive support).**  Tutoring effectiveness depends on providing support that is contingent on the learner's observed state rather than a fixed script. Classic scaffolding work characterizes effective tutoring as dynamic assistance responsive to learner needs (Wood et al., 1976). A synthesis emphasizes three defining features of scaffolding: *contingency*, *fading*, and *transfer of responsibility* (Van de Pol et al., 2010). MetaCLASS operationalizes contingency by tying each coach move to an explicit *Event→Signal→Move→Effect* trajectory, ensuring interventions are triggered by learner signals rather than tutor preference.

**Productive struggle and timing of instruction.**  Learning can benefit when students attempt sensemaking before receiving explanations. Productive Failure shows that attempting solutions (and even failing) can prepare learners to learn more deeply from subsequent instruction (Kapur, 2008). Related work argues there is a "time for telling": direct explanation is more effective after learners have constructed differentiated knowledge structures through prior activity (Schwartz and Bransford, 1998). In MetaCLASS, this motivates NO_INTERVENTION and low-directiveness moves early in a struggle episode, preserving space for learner-generated progress before steering.

**Restraint and the assistance dilemma.**  Tutors face the *assistance dilemma*: too little help leads to floundering, while too much help can reduce effort, shallow processing, and self-regulation (Koedinger and Aleven, 2007). Analyses of human tutoring similarly show that effective tutors elicit student reasoning and selectively intervene (Chi et al., 2001). MetaCLASS encodes restraint with an explicit target distribution (35–50% NO_INTERVENTION) and by treating minimal dialogue acts as first-class pedagogical actions rather than failures to respond.

**Prioritizing internal regulation before external resources.**  Support should diminish and responsibility should shift toward the learner as problem solving progresses. Work on faded guidance (e.g., fading worked solution steps) supports the idea that gradually removing scaffolds can promote independent problem solving and deeper engagement (Renkl et al., 2004). MetaCLASS reflects this principle by prioritizing internal regulation moves (e.g., STRATEGY_ALTERNATIVE, REPRESENTATION_SHIFT) before external resource prompting, discouraging premature outsourcing when the learner can still make productive progress.

**Strategic help-seeking as a target competency.**  Help-seeking is a central SRL behavior: learners must decide *when* to seek help, *what* help to request, and *how* to use it (Aleven et al., 2006b). Empirical ITS work demonstrates that metacognitive feedback can improve learners' help-seeking actions and their use of assistance (Roll et al., 2011b), and broader reviews summarize both the promise and limits of help features in ITSs (Aleven et al., 2016). MetaCLASS generates hints during planning but delivers them only upon student request; additionally, PROMPT_RESOURCE is constrained to occur *after* internal strategy moves, promoting disciplined, student-initiated resource use.

Together, these principles ground MetaCLASS in learning science: coaching is contingent, preserves productive struggle, avoids over-assistance, and treats help-seeking as a coached skill rather than a shortcut.

## B  Model Parameters

Following model card recommendations, we use the following sampling parameters. For Qwen3 models in thinking mode, we use temperature 0.6, top-$p$ 0.95, and top-$k$ 20. For Qwen3 models in non-thinking mode, we utilize temperature 0.7 and top-$p$ 0.8. For gpt-oss models, we use temperature 1.0 and top-$p$ 1.0. Additionally, for Phi-4 models, we use temperature 0.8 and top-$p$ 0.95.

> **Example MetaCLASS Conversation**
>
> **Turn 1 [STUDENT]:** Let me restate this: Josh buys a house for $80,000, spends $50,000 on repairs, and those repairs increase the value of the house by 150%. I need his profit, so essentially selling price minus total cost.
> **Turn 2 [COACH]:** *[NO_INTERVENTION]*
> **Turn 3 [STUDENT]:** Total cost first: 80,000 plus 50,000 is 130,000. So he has 130k invested in the house.
> **Turn 4 [COACH]:** *[NO_INTERVENTION]*
> **Turn 5 [STUDENT]:** Now, "increased the value of the house by 150%." I'm pausing here, because I'm not totally sure what the 150% is relative to. Is it 150% of the original 80k, or 150% of the 50k repairs, or something else?
> **Turn 6 [COACH]:** What are the different reasonable ways you could interpret that phrase? *[UNCERTAINTY_PROBE]*
> **Turn 7 [STUDENT]:** I see at least two. One: the house's value increased by 150% of its original value, so the new value would be some multiple of 80k. Two: the repairs somehow generated a 150% return on the 50k, so the increase is based on 50k. I can reason through it, but to be efficient, I'd like a scaffolding hint specifically about how to interpret "increased the value of the house by 150%."
> **Turn 8 [SYSTEM]:** *Scaffolding hint:* Focus on the phrase "increased the value of the house by 150%." Try rewriting it as an equation in words: new value = original value + 150% of original value. If that still feels abstract, imagine a $100 house and ask yourself what "increased by 150%" would make its new value, then map that pattern back to this problem.
> **Turn 9 [COACH]:** *[NO_INTERVENTION]*
> **Turn 10 [STUDENT]:** Okay, so "new value = original + 150% of original." If I test with a $100 house, 150% of 100 is 150, so the new value would be 100 + 150 = 250 dollars, i.e., 2.5 times the original. Applying that here, the house's new value should be 2.5 times 80,000, so that's 200,000 as the selling price. Then profit would be selling price minus total cost: 200,000 - 130,000 = 70,000.
> **Turn 11 [COACH]:** How does the interpretation you used line up with the exact wording of the problem compared to the alternative you had in mind? *[CHECK_CONSISTENCY]*
> **Turn 12 [STUDENT]:** The problem says "increased the value of the house by 150%," which sounds like a percentage of the house's value, not a percentage return on the repair cost. My interpretation—base is the original 80k—matches that: the house's value itself is what's being compared before and after. The alternative, using 50k as the base, would really be about the repairs' return, but that's not what the sentence names as increasing.
> **Turn 13 [COACH]:** What will you carry forward from this problem for the next time you see something like "increased by X%"? *[REFLECT_OUTCOME]*
> **Turn 14 [STUDENT]:** Two things: mathematically, "increased by X%" means final = original + X% of original, so the new amount is $(1 + X/100)$ times the original. Metacognitively, when the wording feels ambiguous, it helps to explicitly list possible bases for the percent, try a simpler number like 100 to see the pattern, and then, if I'm still unsure, ask for a focused hint on the interpretation, not the whole solution.

Figure 4: Example of generated MetaCLASS conversation.

# C   Example of MetaCLASS dataset

An example of generated MetaCLASS conversation is shown in Figure 4.