

# Dynamic Student Cognitive Modeling and Performance Prediction in Multi-Agent Teaching Simulations

Anonymous ACL submission

## Abstract

In multi-agent teaching simulation scenarios, large language models (LLMs) exhibit an inherent assistant-oriented bias, leading them to generate overly advanced responses when acting as student agents, which limits their ability to accurately reflect real students' cognitive states and learning behaviors. To address this limitation, we model students' cognitive graphs and propose **GraphLR-MPP**, a **Graph**-structured **Learning Report** Enhanced **Math Performance Predictor**, which leverages GraphRAG-generated learning reports to train the model for predicting students' math problem-solving performance. Experimental results demonstrate that our method outperforms existing in-context learning (ICL) approaches and other supervised fine-tuning (SFT) methods. Furthermore, we introduce a Multi-Agent Teaching Intervention Trial that simulates the dynamic updating of students' cognition under instructional interventions, providing a scalable foundation for future agent-based teaching simulation experiments.

## 1 Introduction

To date, LLM-driven agents have been proven effective in simulating human behavior across various domains, including medicine, law, education, and gaming. In educational contexts, agent-based simulations typically rely on interactions among multiple roles (e.g., teachers, students, teaching assistants, and peers) to help students solve specific problems. Due to their inherently assistant-oriented nature, LLMs perform particularly well when simulating teachers or assistants. However, this same characteristic leads to suboptimal performance when LLMs are tasked with modeling students, especially lower-performing students. In such cases, they often fail to realistically reproduce the mistakes that real students are likely to make during problem solving, instead tending to generate consistently correct responses (He-Yueya et al.,

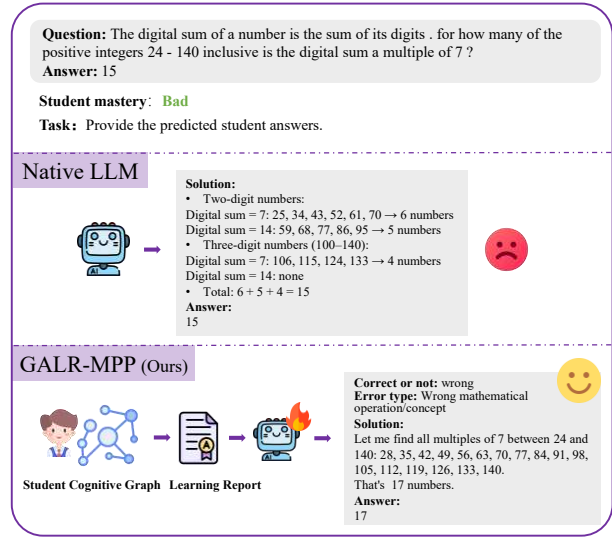


Figure 1: Compare between native LLM and our method when predicting the answer of a bad performance student.

2024; Markel et al., 2023; Aher et al., 2023), as show in Figure 1.

Building on these observations, we propose a GraphRAG-based dynamic cognitive modeling approach for student simulation. Specifically, We first construct individualized student cognitive graphs from their historical problem-solving records, where each node represents a knowledge concept along with the student's mastery level of that concept, while edges encode relationships between concepts. Related knowledge points are then clustered into communities, and a summary of each community is generated as a learning report, reflecting the student's overall mastery of that subset of knowledge. The resulting cognitive graph serves as the foundation for predicting students' subsequent problem-solving performance.

Current LLM-based approaches for student problem-solving prediction primarily rely on ICL, where predictions are made using either students' mastery of relevant knowledge components or

064	records retrieved based on textual similarity. How-	model for student math problem-solving pre-	115
065	ever, approaches based on isolated knowledge	diction, achieving superior performance over	116
066	points fail to capture the intrinsic relationships	existing baselines.	117
067	and interactions among concepts, leading to high		
068	variance in prediction. In contrast, similarity-	• We design a multi-agent teaching intervention	118
069	based retrieval relies on surface-level textual fea-	simulation, where the cognitive graph is up-	119
070	tures without a deeper understanding of underlying	dated to reflect students' dynamic changes un-	120
071	knowledge structures, making it prone to selecting	der interventions, providing insights for future	121
072	records that are textually similar but conceptually	research on adaptive educational strategies.	122
073	irrelevant, thereby introducing noise into the pre-		
074	diction process.		
075	Both approaches rely on limited and largely	<b>2 Related Work</b>	123
076	static features, which restrict the contextual infor-	<b>2.1 Agent-based Educational Simulations</b>	124
077	mation available to the model and often lead to	Large language model (LLM)-driven agents have	125
078	oversimplified or biased predictions. To address	demonstrated human-like behaviors across diverse	126
079	this limitation, we propose <b>GraphLR-MPP</b> , a	domains, including social interaction, medicine,	127
080	math performance prediction framework enhanced	law, software development, and gaming (Park et al.,	128
081	by graph-structured learning reports generated	2023; Li et al., 2024; Fan et al., 2024; Qian et al.,	129
082	from students' cognitive graphs. These reports	2024; Wang et al., 2023). In education, recent work	130
083	enable the model to capture rich student charac-	leverages multi-agent and cognitively grounded	131
084	teristics, including knowledge summaries, practice	LLM simulations to model realistic classroom	132
085	performance, and evolving mastery states. Through	dynamics, including long-term learner modeling,	133
086	supervised fine-tuning, GraphLR-MPP further al-	Theory-of-Mind-based peer interaction, adaptive	134
087	lows the LLM to learn implicit relationships be-	teaching strategies, personality-aware student sim-	135
088	tween learning reports and student responses, re-	ulation, and interactive teacher training environ-	136
089	sulting in more personalized, accurate, and stable	ments (Yuan et al., 2025; Gao et al., 2025; Sanyal	137
090	predictions. Experimental results demonstrate that	et al., 2025; Liu et al., 2024b; Markel et al., 2023;	138
091	our method consistently outperforms existing ICL-	Jinxin et al., 2023; Zhang et al., 2024). Building	139
092	based approaches as well as other SFT baselines in	on these general educational settings, several stud-	140
093	prediction accuracy, with particularly strong gains	ies further focus on mathematics education, where	141
094	in error-type prediction, making the results more	agent-based tutoring systems explore diverse in-	142
095	informative and practically valuable.	structional strategies for problem solving and dia-	143
096	Our fine-grained approach to student modeling	logue generation (Liu et al., 2024a; Yue et al., 2024;	144
097	and performance prediction provides a foundation	Liu et al., 2025; Macina et al., 2023; Ding et al.,	145
098	for multi-agent interactions in educational settings.	2024). In this paper, we focus on student cogni-	146
099	To this end, we developed a teaching simulation en-	tive modeling and teacher-student interactions for	147
100	vironment to predict the effects of two instructional	mathematical word problems.	148
101	interventions (Explicit, Systematic Instruction and		
102	Metacognitive Strategies) on student performance.	<b>2.2 Student Cognitive Modeling</b>	149
103	Students' cognitive states are dynamically updated	Student cognition modeling aims to infer learners'	150
104	throughout the intervention process, and the results	knowledge states and cognitive processes to sup-	151
105	indicate that both strategies effectively enhance	port personalized learning and performance predic-	152
106	student outcomes, consistent with findings from	tion. Early probabilistic models such as Bayesian	153
107	real-world educational research.	Knowledge Tracing (BKT) (Corbett and Anderson,	154
108	Our main contributions are as follows:	1994) and later deep learning-based approaches, in-	155
109	• We propose a GraphRAG-based method for	cluding DKT (Piech et al., 2015), DKVMN (Zhang	156
110	constructing dynamic student cognitive mod-	et al., 2017), SAKT (Pandey and Karypis, 2019),	157
111	els, which structures the relationships between	and graph-based knowledge tracing (Nakagawa	158
112	knowledge concepts and captures students'	et al., 2019), estimate knowledge mastery from	159
113	overall cognitive profiles.	structured interaction records. However, these	160
114	• We propose GraphLR-MPP, a graph-enhanced	methods rely on predefined concepts and highly	161
		structured data, limiting their ability to capture rich	162

reasoning processes and error patterns. With the advent of large language models, role-playing agents enable more expressive student modeling by capturing both cognitive and non-cognitive dimensions through semantic representations, supporting solution simulation in addition to outcome prediction. Recent work has introduced structured representations of student cognition, such as tree-based personality modeling (Jinxin et al., 2023) and graph-based cognition modeling (Wu et al., 2025). Nevertheless, most existing approaches remain largely static and fail to model the dynamic evolution of student cognition. To address this limitation, we represent student cognition as a continuously updated knowledge graph.

### 3 Methods

Our research focuses on modeling students’ dynamic cognition and predicting their mathematical problem-solving performance. We adopt a GraphRAG-based approach to construct student cognition and generate learning reports by clustering related knowledge concepts, which are then used to train GraphLR-MPP. Finally, we demonstrate a multi-agent teaching simulation to validate the effectiveness of instructional interventions through dynamically updated student cognition. The overall framework is shown in Figure 2.

#### 3.1 Student Cognitive Graph Construction

Given a student’s historical problem-solving records, we construct a personalized knowledge graph for the student. This process is divided into three steps: entity extraction, relation extraction, and community clustering.

**Entity Extraction.** Firstly, an LLM is applied to each historical record to extract all knowledge points involved in the corresponding problem. Each knowledge point is assigned one of the predefined mathematical knowledge types: *Concept*, *Theorem*, *Method*, *Skill*, or *Model*. Based on the student’s performance on the problem, the student’s mastery of each knowledge point is further inferred as either *Good* or *Bad*, and a textual description of the mastery status is generated accordingly. Each knowledge point, together with its mathematical knowledge type and the student’s mastery information, forms a single node in the knowledge graph.

This process is formally defined as:

$$\mathbf{V}_i = \Pi_v(\mathbf{R}_i) = \{v_{i1}, v_{i2}, \dots, v_{ik}\} \quad (1)$$

where  $R_i$  denotes the  $i$ -th historical problem-solving record of a student, including the problem statement, the student’s solution process and answer, and the corresponding error analysis. Each node  $v_{ij}$  is represented as a tuple  $v_{ij} = (n_{ij}, t_{ij}, m_{ij})$ , where  $n_{ij}$  is the name of the knowledge point,  $t_{ij}$  denotes its mathematical knowledge type, and  $m_{ij}$  represents the student’s mastery information, including the mastery level and its textual description.

**Relation Extraction.** After extracting the knowledge points corresponding to each record, we further employ an LLM to infer the relationships between these knowledge points. According to definitions in mathematics education, we categorize the relationships into the following five types:

- **Prerequisite\_of:** One knowledge point is a prerequisite for learning another.
- **Equivalent\_of:** Two concepts describe the same mathematical object or property.
- **Contains\_in:** One knowledge point is a special case or subset of another.
- **Parallel\_to:** Two knowledge points belong to different branches of the same category.
- **Application\_of:** A knowledge point is an application of another.

Formally, the set of extracted relationships can be expressed as:

$$\mathbf{E}_i = \Pi_e(\mathbf{V}_i) = \{e_1, e_2, \dots, e_m\} \quad (2)$$

Each relationship  $e_j$  is represented as a tuple  $e_j = (s_j, t_j, r_j, d_j, w_j)$ , where  $s_j$  and  $t_j$  denote the source and target entities,  $r_j$  indicates the relationship type,  $d_j$  provides a textual description of the relationship, and  $w_j \in [0, 1]$  represents the strength of the relationship between the source and target entities.

**Community Clustering.** After constructing the entity–relation graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , we perform community clustering to identify coherent groups of knowledge points. Specifically, following the community partition method in Edge et al. (2024), we apply a hierarchical Leiden algorithm to the largest connected component of the graph to obtain multi-level community partitions. This process can be expressed as:

$$\mathcal{C} = \Phi_{\text{Leiden}}(\mathcal{G}) \quad (3)$$

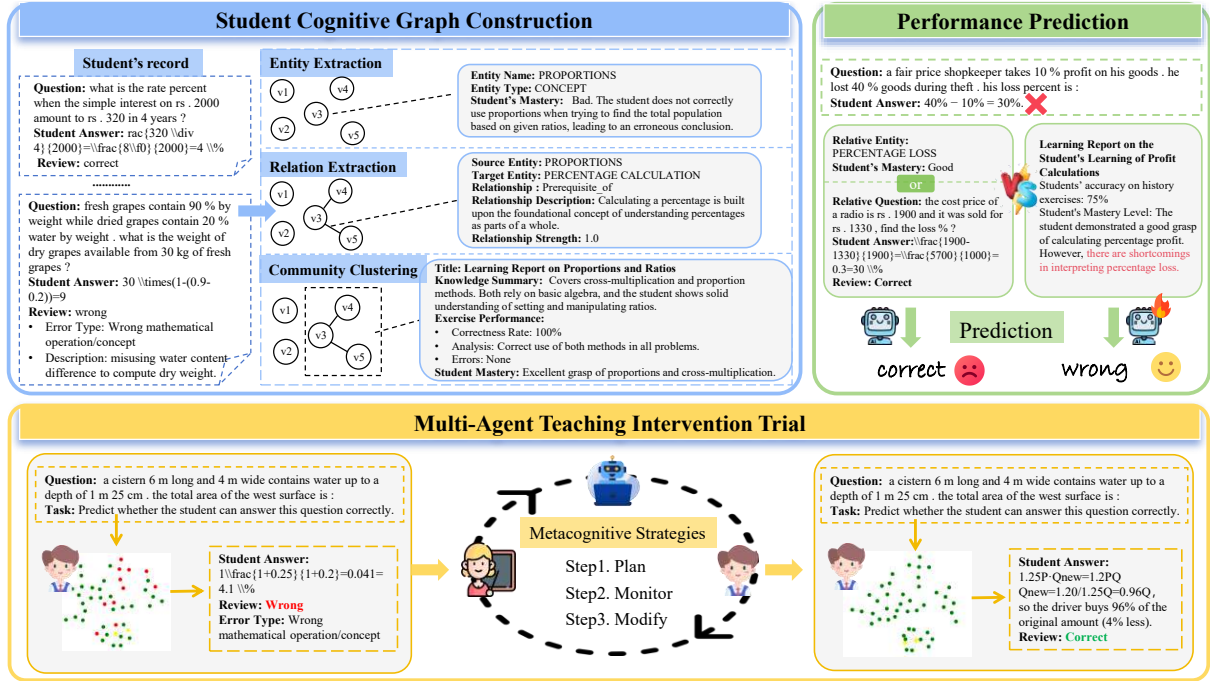


Figure 2: The complete framework of this study, consisting of three components: Student Cognitive Graph Construction, Performance Prediction, and Multi-Agent Teaching Intervention Trial.

where  $\mathcal{C} = \{c_{l,k}\}$  denotes the set of communities across different hierarchical levels  $l$ , where  $k$  denotes the community index, and each  $c_{l,k}$  corresponds to a subset of knowledge point nodes. A single node may belong to different communities at different levels, enabling multi-granularity representations of the student's cognitive structure.

Based on the resulting community structure, we further generate a learning report for each community. For a given community  $c_{l,k}$ , its internal entities and relations are first aggregated and serialized into a textual description, which is then fed into an LLM to produce a structured learning report:

$$\mathcal{R}_{l,k} = \Pi_r(c_{l,k}) \quad (4)$$

The reports are generated in a top-down manner across hierarchy levels, allowing higher-level community reports to incorporate information from previously generated lower-level reports. Each learning report summarizes the student's past performance on exercises related to the corresponding knowledge points, including the overall accuracy, detailed results for each knowledge point, and the main types of errors. Based on this, a concise summary of the student's mastery for the knowledge section is generated, providing teachers with a systematic evaluation of the student's learning status.

### 3.2 Student Performance Prediction

Existing LLM-based student performance prediction methods rely on ICL with limited and discrete retrieved information, which often introduces bias. As shown in Figure 2, such retrieval may misleadingly suggest strong mastery and lead to incorrect predictions. In contrast, learning reports provide a structured representation of students' knowledge states, but directly applying ICL with them tends to produce overly pessimistic predictions. Therefore, we adopt a SFT strategy for more reliable performance prediction.

Specifically, we fine-tune a general LLM on a student performance prediction dataset  $D_s$ . Each training example consists of a structured input sequence  $X = \{X_p, X_r\}$ , where  $X_p$  denotes the problem statement together with the corresponding learning report, and  $X_r = (x_1, x_2, \dots, x_L)$  represents the student performance annotations (e.g., answer correctness and error types), serialized as target tokens. The model is trained to maximize the likelihood of the student response tokens conditioned on the preceding context. The training objective is defined as:

$$\mathcal{L} = - \sum_{i=1}^L \log p_{\theta}(x_i | X_p, X_{r,<i}), \quad x_i \in X_r, \quad (5)$$

where  $L$  is the length of the target sequence,  $x_i$  denotes the  $i$ -th token in the student performance sequence, and  $X_{r, < i}$  represents the previously generated student performance tokens before position  $i$ .

By minimizing this negative log-likelihood over all samples in  $D_s$ , the model learns an explicit mapping from historical exercise accuracy and learning reports to actual student responses, resulting in the final predictive model, **GraphLR-MPP**, which achieves more accurate and calibrated predictions than ICL-based approaches.

### 3.3 Multi-Agent Teaching Intervention Trial

Most existing student cognitive modeling approaches in teaching scenario simulation with LLM-Based agents rely on fixed and static representations. However, in real-world settings, students' cognition always evolves over time. Motivated by this observation, we design a multi-agent teaching intervention trial to capture such dynamic changes.

#### 3.3.1 Task Design and Hypotheses

Specifically, we first construct the student cognitive model using the method described in Section 3.1. Subsequently, the teacher agent applies two instructional intervention strategies (Explicit, Systematic Instruction and Metacognitive Strategies) to re-teach the problems that the student previously answered incorrectly. The former delivers concepts or procedures in a structured and sequential manner, consisting of lesson orientation, initial instruction, guided practice, and independent practice. While the latter guide students to regulate their problem-solving through planning, monitoring, and modification.

After the whole intervention process is completed, the student's cognitive graph is updated accordingly. To evaluate the effectiveness of the intervention strategies, we assess the student's performance on a test set of previously unseen problems before and after the intervention. Since both the instructional intervention strategies are evidence-based practices (EBPs), we expect that they can effectively improve student performance. Accordingly, we formulate the following two hypotheses:

- H1: Explicit, Systematic Instruction enhances students' cognitive states and improves their accuracy in mathematical problem solving.

---

#### Algorithm 1 Multi-agent collaboration algorithm

**Require:** Student  $S$ , Teacher  $T$ , Controller  $C$ , Maxturn  $m$ , Teacher's response  $r^T$ , Student's response  $r^S$ , Controller's response  $r^C$ , Dialogue history  $H$ , Current intervention step  $k$

**Input:** Intervention Strategy  $I$ , Question  $q$ , Correct Answer  $CA$ , Student Answer  $SA$ , Review  $R$

**Output:** Dialogue History  $H$

```

1:  $T \leftarrow (I, q, CA, SA, R)$ 
2:  $S \leftarrow (q, SA)$ 
3:  $k = 0$ 
4:  $T \leftarrow I^{(0)}$ 
5:  $T$  generate first turn guidance  $r_0^T$ 
6:  $H \leftarrow H \cup \{r_0^T\}$ 
7: for  $i = 1$  to  $m$  do
8:    $S$  generate  $r_i^S$ 
9:    $H \leftarrow H \cup \{r_i^S\}$ 
10:   $C \leftarrow H^{(k)}$ 
11:   $C$  decides whether to stay, move to next step, or end.
12:  if "end" in  $r_i^C$  then
13:    break
14:  else if "next step" in  $r_i^C$  then
15:     $k = k + 1$ 
16:     $T \leftarrow I^{(k)}$ 
17:  end if
18:   $T$  generate  $i$ -th turn guidance  $r_i^T$ 
19:   $H \leftarrow H \cup \{r_i^T\}$ 
20: end for
21: return  $H$ 

```

---

- H2: Metacognitive Strategies enhance students' cognitive states and improve their accuracy in mathematical problem solving.

By testing whether these hypotheses hold, we aim to demonstrate the effectiveness of our student modeling and simulation framework, thereby providing empirical support for the reliability of using this framework for future simulated teaching trial.

#### 3.3.2 Multi-agent Collaboration Framework

We simulate the instructional intervention process through multi-round interactions between teacher and student agents. In addition, we introduce a controller to manage the session flow, which determines whether to proceed to the next phase based on the current interaction. The overall interaction procedure is presented in Algorithm 1.

Table 1: Statistics of Student Performance in MathEDU-Plus.

Student	Total	Incorrect	Accuracy
Student1	668	186	72.16%
Student2	685	85	87.59%
Student3	676	194	71.30%
Student4	660	163	75.30%
Student5	652	195	70.09%
Student6	630	99	84.29%

## 4 Experimental setup

### 4.1 Datasets

Our experiments are conducted using data derived from the MathEDU dataset (Hsu et al., 2025), which collects the authentic problem-solving processes of six real students as they solve mathematics problems from the MathQA dataset (Amini et al., 2019), covering topics such as Geometry and Probability. Each student solution is annotated by three mathematics experts with respect to correctness, error type, and error equation, along with expert advice for revising the solution. The error type include *Wrong Mathematical Operation/Concept*, *Calculation Error*, and *Incomplete Answer*, etc.

It is worth noting that the MathEDU dataset contains only problem IDs and does not include the textual descriptions of the problems. To facilitate subsequent processing, we retrieve the corresponding problem statements and correct answers from the original MathQA dataset and merge them with the student solution processes in MathEDU. During preprocessing, we remove duplicate problems, exclude entries for which the corresponding problems cannot be found, and eliminate unnecessary Chinese fields. The resulting processed dataset, denoted as MathEDU-Plus, contains cleaned records for all six students. The number of records and answer accuracy for each student are reported in Table 1.

For each student, we split the data into training and test sets using a 7:3 ratio. The training set is used for knowledge graph construction and fine-tuning the error prediction model, while the test set is used to evaluate the model’s accuracy in predicting the student’s problem-solving behavior.

### 4.2 Metrics

Our evaluation is conducted in two main parts. In the Student Performance Prediction setting, we report the overall accuracy of different methods in

predicting students’ answer outcomes. We further evaluate performance on correct answers and incorrect answers using *Precision*, *Recall*, and *F1-score*, respectively. In addition, we place particular emphasis on the accuracy of error type prediction. To this end, we employ an LLM-based evaluation to assess whether the predicted error types are consistent with the students’ ground-truth errors. In the Multi-Agent Teaching Intervention Trial, we focus on examining how students’ knowledge graphs are updated after instructional interventions and whether their answer accuracy improves accordingly. These results serve to validate the effectiveness of our dynamic student cognitive graph construction method as well as the proposed instructional intervention strategies.

### 4.3 Baselines

We compare our method with five ICL baselines and two supervised SFT methods. Among the ICL approaches, one follows Wu et al. (2025), which retrieves five relevant entities and selects a learning record covering the most entities as context. We also include Random and Similarity baselines from the same study, as well as an Entities-only variant to isolate the effects of entities and records. In addition, we evaluate an ICL baseline that directly uses raw learning reports as context to assess the necessity of fine-tuning. All ICL methods are tested with both gpt-4o-mini and gpt-4o. Beyond ICL, we further compare our approach with two SFT baselines based on five relevant entities or one relevant record, demonstrating the advantage of fine-tuning on learning reports.

### 4.4 Implementation Details

All experiments are conducted on the MathEDU-Plus dataset. For each selected student, their learning records are divided into training and test sets at a 7:3 ratio. Based on the training data, we construct student cognitive knowledge graphs using the GraphRAG framework (Edge et al., 2024), which enables structured representation of students’ knowledge states. Specifically, GPT-4o-mini is employed to extract entities and relations, and the Leiden algorithm is used to cluster related knowledge concepts into communities. For each community, GPT-4o-mini is further used to generate a corresponding learning report. During the testing phase, for each given problem, we retrieve the five most relevant entities using the text-embedding-3-small model. The

most relevant learning report is then identified based on the cluster to which these entities belong.

For performance prediction, We fine-tune the model using the training data. The model is initialized with Qwen2.5-7B-Instruct and fine-tuned for five epochs using the LoRA approach (Hu et al., 2022), with a rank of 8 and a scaling factor of  $\alpha = 32$ . Training is conducted with a batch size of 2, a learning rate of  $2 \times 10^{-4}$  using a cosine scheduler, and a warm-up ratio of 0.1.

For multi-agent simulation, both the student agent and the teacher agent are powered by GPT-4o, while the controller agent is driven by GPT-4o-mini, given the relatively simple nature of the control task. We set the generation temperature to 0.7 for both the student and teacher agents. In contrast, the controller’s generation temperature is set to 0.3 to maintain a higher level of determinism.

## 5 Results

### 5.1 Student Performance Prediction

#### 5.1.1 Overall Prediction

The average performance of our method and seven baselines in predicting the behaviors of the six students is reported in Table 2. Overall, our method achieves higher average accuracy than the competing approaches.

With the exception of the method that uses learning reports as context, the remaining four ICL baselines exhibit similar prediction patterns, characterized by strong performance in predicting correct student answers (F1 Score > 0.75) and comparatively weaker performance in predicting incorrect answers (F1 Score fluctuating between 0.3 and 0.4). Interestingly, the ICL method that leverages learning reports as context demonstrates an opposite trend: it consistently outperforms the other ICL methods in predicting incorrect answers, particularly in terms of recall, despite its inferior performance in predicting correct answers. By contrast, our fine-tuned model combines the strengths of both approaches, effectively leveraging the richer information contained in learning reports to achieve a more balanced and desirable performance in predicting both correct and incorrect answers. Similarly, our approach consistently surpasses the other two SFT baselines on all metrics.

Notably, for the ICL methods, to avoid potential bias from any single model, we further perform a Pearson’s correlation analysis on the predictions of GPT-4o and GPT-4o-mini, which yields a very

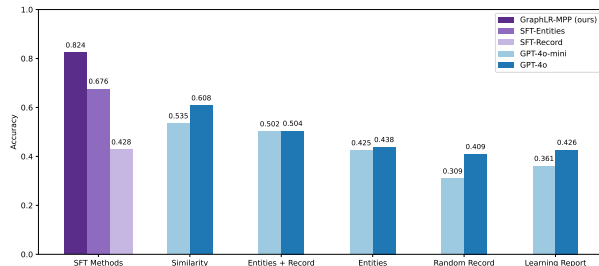


Figure 3: the proportion of consistent error types.

strong positive correlation ( $\rho \approx 0.99$ ), demonstrating the robustness and cross-model consistency of all ICL methods.

#### 5.1.2 Error-type Prediction

To further evaluate the effectiveness of different methods in predicting students’ incorrect answers, we examine the proportion of correctly retrieved answers for which the predicted error type matches the ground-truth error type. The results are shown in Figure 3. Our method achieves an accuracy of 82.4%, substantially outperforming all baseline approaches. In contrast, the method that uses a randomly selected record as context attains an accuracy of only around 0.35, indicating limited reliability in error-type prediction. These results demonstrate that, compared with random prediction, our approach provides a more faithful and reliable approximation of students’ error behaviors.

### 5.2 Simulated Teaching Intervention

We simulate the teaching intervention process within an interactive agent-based environment, taking Student1 and Student5 as illustrative examples. In the original answer records, student1 contained 125 incorrect responses, while student5 had 127 incorrect responses. We allow intelligent agents to sequentially assume the roles of these two students. For each incorrect item, a teacher agent provides re-instruction using two pedagogical intervention strategies: Explicit, Systematic Instruction (ESI) and Metacognitive Strategies (MS).

We then evaluate the simulated students’ post-intervention performance. The results show that, for Student1, 93 items were corrected after ESI, while 114 items were corrected following MS. Similarly, Student5 corrected 102 and 112 items under the two respective intervention strategies. Based on the updated student response records, we reconstruct the corresponding knowledge graphs. As illustrated in Figure 4, red nodes denote misconcep-

Methods	SFT						ICL							
	GraphLR-MPP	Entities	Record	Similarity		Entities+Record		Entities		Random Record		Learning Report		
Base Model	Qwen2.5-7B-Instruct			4o-mini	4o	4o-mini	4o	4o-mini	4o	4o-mini	4o	4o-mini	4o	
Overall Acc.	<b>0.748</b>	0.677	0.707	0.696	0.668	0.696	<u>0.713</u>	0.701	0.710	0.670	0.667	0.57	0.468	
Correct	Prec.	0.787	0.759	0.756	<u>0.799</u>	0.786	0.764	0.770	0.767	0.770	0.757	0.766	0.792	<b>0.801</b>
	Rec.	<b>0.888</b>	0.820	<u>0.887</u>	0.786	0.757	0.828	0.853	0.839	0.852	0.807	0.778	0.554	0.379
	F1	<b>0.833</b>	0.786	<u>0.813</u>	0.787	0.766	0.793	0.808	0.799	0.807	0.778	0.769	0.649	0.496
Wrong	Prec.	<b>0.480</b>	0.327	0.341	<u>0.438</u>	0.391	0.395	0.429	0.407	0.420	0.352	0.364	0.320	0.310
	Rec.	0.316	0.250	0.165	0.466	0.428	0.304	0.309	0.303	0.306	0.303	0.359	<u>0.580</u>	<b>0.735</b>
	F1	0.380	0.276	0.221	<b>0.431</b>	0.390	0.342	0.357	0.344	0.348	0.309	0.353	0.405	<u>0.415</u>

Table 2: Performance comparison between SFT and various ICL strategies. Best results are in bold, second best are underlined.

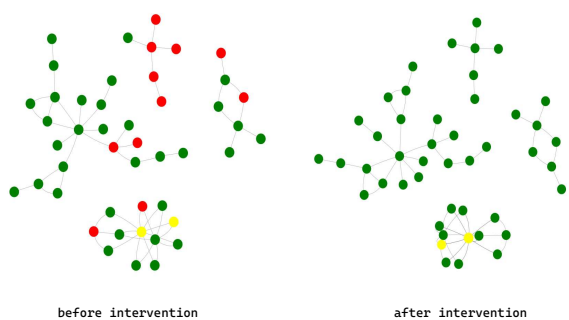


Figure 4: Student’s cognitive graph before and after intervention. Red indicates not mastered, green indicates mastered, and yellow indicates partially mastered.

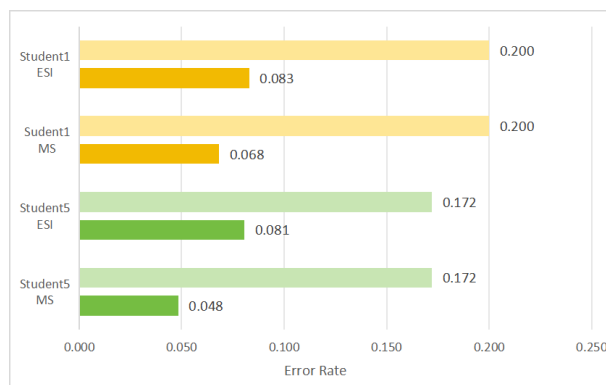


Figure 5: Error rate changes for Student1 and Student5 before and after ESI and MS interventions (light: before; dark: after).

tions or incorrect knowledge states, whereas green nodes indicate mastered concepts. The transition of nodes from red to green reflects the updating and refinement of students’ cognitive states as a result of the teaching interventions.

Given the test set, we retrieved the problem-relevant learning reports based on both the initial knowledge graph and the post-intervention knowledge graph, and used the fine-tuned model to predict students’ responses. The results, shown in Figure 5, indicate that for Student1, the error rates decreased by 58.54% and 65.85% after ESI and MS, respectively. For Student5, the corresponding reductions were 53.12% and 71.88%. These findings demonstrate that both instructional interventions positively impact students’ mathematics problem-solving performance, with Metacognitive Strategies outperforming Explicit, Systematic Instruction. This conclusion provides empirical support and reference for instructional design in real classroom settings.

## 6 Conclusion

This paper constructs student cognitive graphs based on the GraphRAG framework, leveraging LLM to extract knowledge concepts and the relationships among them from students’ historical learning records, which are represented as nodes and edges in the graph. Closely related concepts are further clustered to generate learning reports that reflect students’ mastery of specific knowledge areas. Based on these learning reports and students’ answer records, we fine-tune a base model to obtain a student performance prediction model GraphLR-MPP, achieving superior performance over existing baselines, particularly in error-type prediction. Finally, through a Multi-Agent Teaching Intervention Trial, we simulate the dynamic evolution of students’ cognition under instructional interventions by updating the knowledge graph, thereby providing a unified framework and empirical reference for future simulated teaching experiments.

## 592 Limitations

593 Most existing studies on multi-agent educational  
594 simulation rely on synthetic datasets, while real-  
595 world student response data remain scarce, raising  
596 concerns about the fidelity of simulated behaviors.  
597 To better approximate realistic educational settings,  
598 we perform student cognitive modeling using the  
599 MathEDU dataset, which is collected from authentic  
600 student problem-solving records. But these stu-  
601 dents achieve relatively high accuracy rates (above  
602 70%), leading to a shortage of error cases and lim-  
603 iting the diversity of performance prediction and  
604 error-type analysis. Despite these constraints, ex-  
605 tensive experiments and comprehensive compar-  
606 isons show that our method consistently outper-  
607 forms all baselines. We anticipate that evaluations  
608 on larger and higher-quality datasets will further  
609 strengthen the robustness and generalizability of  
610 our approach.

## 611 Ethics Statements

612 The data used in this work are publicly available  
613 open-source data, and thus do not raise any ethical  
614 or moral concerns.

## 615 Use of AI Assistants

616 We primarily used AI tools to assist with writing  
617 polishing and language refinement.

## 618 References

619 Gati V Aher, Rosa I Arriaga, and Adam Tauman Kalai.  
620 2023. Using large language models to simulate mul-  
621 tiple humans and replicate human subject studies. In  
622 *International conference on machine learning*, pages  
623 337–371. PMLR.

624 Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik  
625 Koncel-Kedziorski, Yejin Choi, and Hannaneh Ha-  
626 jishirzi. 2019. Mathqa: Towards interpretable math  
627 word problem solving with operation-based formalisms. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: Human language technologies, volume 1 (long and short papers)*, pages 2357–2367.

633 Albert T Corbett and John R Anderson. 1994. Knowl-  
634 edge tracing: Modeling the acquisition of procedural  
635 knowledge. *User modeling and user-adapted inter-*  
636 *action*, 4(4):253–278.

637 Yuyang Ding, Hanglei Hu, Jie Zhou, Qin Chen,  
638 Bo Jiang, and Liang He. 2024. Boosting large lan-  
639 guage models with socratic method for conversa-  
640 tional mathematics teaching. In *Proceedings of the*

*33rd ACM International Conference on Information  
and Knowledge Management*, pages 3730–3735. 641 642

Darren Edge, Ha Trinh, Newman Cheng, Joshua  
Bradley, Alex Chao, Apurva Mody, Steven Truitt,  
Dasha Metropolitansky, Robert Osazuwa Ness, and  
Jonathan Larson. 2024. From local to global: A  
graph rag approach to query-focused summarization.  
*arXiv preprint arXiv:2404.16130*. 643 644 645 646 647 648

Zhihao Fan, Jialong Tang, Wei Chen, Siyuan Wang,  
Zhongyu Wei, Jun Xi, Fei Huang, and Jingren Zhou.  
2024. Ai hospital: Interactive evaluation and collabo-  
ration of llms as intern doctors for clinical diagnosis.  
*arXiv preprint arXiv:2402.09742*. 649 650 651 652 653

Xian Gao, Zongyun Zhang, Ting Liu, and Yuzhuo Fu.  
2025. Onlinemate: An llm-based multi-agent com-  
panion system for cognitive support in online learn-  
ing. *arXiv preprint arXiv:2509.14803*. 654 655 656 657

Joy He-Yueya, Noah D Goodman, and Emma Brunskill.  
2024. Evaluating and optimizing educational content  
with large language model judgments. *arXiv preprint  
arXiv:2403.02795*. 658 659 660 661

Wei-Ling Hsu, Yu-Chien Tang, and An-Zi Yen. 2025.  
Mathedu: Towards adaptive feedback for student  
mathematical problem-solving. *arXiv preprint  
arXiv:2505.18056*. 662 663 664 665

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan  
Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,  
Weizhu Chen, and 1 others. 2022. Lora: Low-rank  
adaptation of large language models. *ICLR*, 1(2):3. 666 667 668 669

Shi Jinxin, Zhao Jiabao, Wang Yilei, Wu Xingjiao, Li Ji-  
awen, and He Liang. 2023. Cgmi: Configurable  
general multi-agent interaction framework. *arXiv  
preprint arXiv:2308.12503*. 670 671 672 673

Junkai Li, Siyu Wang, Meng Zhang, Weitao Li, Yungh-  
wei Lai, Xinhui Kang, Weizhi Ma, and Yang Liu.  
2024. Agent hospital: A simulacrum of hospi-  
tal with evolvable medical agents. *arXiv preprint  
arXiv:2405.02957*. 674 675 676 677 678

Ben Liu, Jihai Zhang, Fangquan Lin, Xu Jia, and Min  
Peng. 2025. One size doesn't fit all: A personal-  
ized conversational tutoring agent for mathematics  
instruction. In *Companion Proceedings of the ACM  
on Web Conference 2025*, pages 2401–2410. 679 680 681 682 683

Jiayu Liu, Zhenya Huang, Tong Xiao, Jing Sha, Jinze  
Wu, Qi Liu, Shijin Wang, and Enhong Chen. 2024a.  
Socraticlm: Exploring socratic personalized teaching  
with large language models. *Advances in Neural  
Information Processing Systems*, 37:85693–85721. 684 685 686 687 688

Zhengyuan Liu, Stella Xin Yin, Geyu Lin, and Nancy  
Chen. 2024b. Personality-aware student simulation  
for conversational intelligent tutoring systems. In  
*Proceedings of the 2024 Conference on Empirical  
Methods in Natural Language Processing*, pages 626–  
642. 689 690 691 692 693 694

695	Jakub Macina, Nico Daheim, Sankalan Chowdhury, Tanmay Sinha, Manu Kapur, Iryna Gurevych, and Mrinmaya Sachan. 2023. Mathdial: A dialogue tutoring dataset with rich pedagogical properties grounded in math reasoning problems. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 5602–5621.	human-like learning dynamics with llm-empowered agents. <i>arXiv preprint arXiv:2508.05622</i> .	751 752
696			
697			
698			
699			
700			
701			
702	Julia M Markel, Steven G Opferman, James A Landay, and Chris Piech. 2023. Gpteach: Interactive training with gpt-based students. In <i>Proceedings of the tenth acm conference on learning@ scale</i> , pages 226–236.	Jiani Zhang, Xingjian Shi, Irwin King, and Dit-Yan Yeung. 2017. Dynamic key-value memory networks for knowledge tracing. In <i>Proceedings of the 26th international conference on World Wide Web</i> , pages 765–774.	753 754 755 756
703			
704			
705			
706			
707	Hiromi Nakagawa, Yusuke Iwasawa, and Yutaka Matsuo. 2019. Graph-based knowledge tracing: modeling student proficiency using graph neural network. In <i>IEEE/WIC/ACM international conference on web intelligence</i> , pages 156–163.	Zheyuan Zhang, Daniel Zhang-Li, Jifan Yu, Linlu Gong, Jinchang Zhou, Zhiyuan Liu, Lei Hou, and Juanzi Li. 2024. Simulating classroom education with llm-empowered agents. <i>arXiv preprint arXiv:2406.19226</i> .	757 758 759 760 761
708			
709			
710			
711			
712	Shalini Pandey and George Karypis. 2019. A self-attentive model for knowledge tracing. <i>arXiv preprint arXiv:1907.06837</i> .		
713			
714			
715	Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In <i>Proceedings of the 36th annual acm symposium on user interface software and technology</i> , pages 1–22.		
716			
717			
718			
719			
720			
721	Chris Piech, Jonathan Bassen, Jonathan Huang, Surya Ganguli, Mehran Sahami, Leonidas J Guibas, and Jascha Sohl-Dickstein. 2015. Deep knowledge tracing. <i>Advances in neural information processing systems</i> , 28.		
722			
723			
724			
725			
726	Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, and 1 others. 2024. Chatdev: Communicative agents for software development. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15174–15186.		
727			
728			
729			
730			
731			
732			
733	Debdeep Sanyal, Agniva Maiti, Umakanta Maharana, Dhruv Kumar, Ankur Mali, C Lee Giles, and Murari Mandal. 2025. Investigating pedagogical teacher and student llm agents: Genetic adaptation meets retrieval augmented generation across learning style. <i>arXiv preprint arXiv:2505.19173</i> .		
734			
735			
736			
737			
738			
739	Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. 2023. Voyager: An open-ended embodied agent with large language models. <i>arXiv preprint arXiv:2305.16291</i> .		
740			
741			
742			
743			
744	Tao Wu, Jingyuan Chen, Wang Lin, Mengze Li, Yumeng Zhu, Ang Li, Kun Kuang, and Fei Wu. 2025. Embracing imperfection: Simulating students with diverse cognitive levels using llm-based agents. <i>arXiv preprint arXiv:2505.19997</i> .		
745			
746			
747			
748			
749	Yu Yuan, Lili Zhao, Wei Chen, Guangting Zheng, Kai Zhang, Mengdi Zhang, and Qi Liu. 2025. Simulating		
750			

## A Experimental Details

### A.1 Dataset Example

The data used to construct the student cognitive graph are obtained from MathEDU-Plus, which is adapted from the MathEDU dataset. It consists of a large number of problem-solving records from six students. Each record contains the problem and its ground-truth answer, the student’s solution process, teacher reviews, etc. An illustrative example is provided in Figure 6.

Problem id: 34421 Problem: of the total amount that jill spent on a shopping trip , excluding taxes , she spent 25 percent on clothing , 25 percent on food , and 50 percent on other items . if jill paid a 10 percent tax on the clothing , no tax on the food , and an 2 percent tax on all other items , then the total tax that she paid was what percent of the total amount that she spent , excluding taxes ?
Student id: 1 Student answer: 0.06 Student process: $0.5 \times 0.1 + 0.5 \times 0.02 = 0.06$
Review: Correct or not: wrong Error counts: 1 Error type: Comprehension error Error equation: $0.5 \times 0.1$ Teacher advice: The student did not notice that Jill only paid tax on clothing, so the student incorrectly calculated the tax on food, leading to an error. The correct answer should be $0.25 \times 1 + 0.5 \times 0.02 = 0.035 = 3.5\%$

Figure 6: Sample from the MathEDU-Plus Dataset.

### A.2 Metrics Details

In the Error Prediction setting, we evaluate model performance from multiple perspectives, including overall answer outcome prediction and fine-grained error type prediction. The following describes each evaluation metric in detail.

#### A.2.1 Overall Accuracy

Overall accuracy measures the proportion of student answers whose correctness (correct or incorrect) is predicted correctly. It is defined as:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (6)$$

where  $TP$  and  $TN$  denote the numbers of true positives and true negatives,  $FP$  and  $FN$  represent false positives and false negatives, respectively.

#### A.2.2 Precision, Recall, and F1-score

We report Precision, Recall, and F1-score to evaluate model performance on both correct answer prediction and incorrect answer prediction. These metrics are computed in a binary classification setting by alternately treating correct answers and incorrect answers as the positive class. Precision is defined as:

$$\text{Precision} = \frac{TP}{TP + FP} \quad (7)$$

where  $TP$  denotes the number of correctly predicted positive instances and  $FP$  denotes the number of incorrectly predicted positive instances.

Recall is defined as:

$$\text{Recall} = \frac{TP}{TP + FN} \quad (8)$$

where  $FN$  denotes the number of positive instances that are incorrectly predicted as negative.

The F1-score is the harmonic mean of Precision and Recall:

$$F1 = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (9)$$

When correct answers are treated as the positive class, these metrics measure the model’s ability to identify correct responses. Conversely, when incorrect answers are treated as the positive class, they evaluate the model’s effectiveness in detecting student errors.

#### A.2.3 Error Type Precision Accuracy

Beyond predicting whether an answer is correct, we place particular emphasis on the accuracy of error type prediction. Let  $E = \{e_1, e_2, \dots, e_N\}$  denote the set of student responses with annotated ground-truth error types. The error type prediction accuracy is defined as:

$$\text{ErrorTypeAcc} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(\hat{e}_i = e_i) \quad (10)$$

where  $\hat{e}_i$  denotes the predicted error type for the  $i$ -th response,  $e_i$  is the corresponding ground-truth error type, and  $\mathbb{I}(\cdot)$  is the indicator function. Since error types often involve nuanced reasoning and implicit misconceptions, we adopt an LLM-based evaluation protocol to assess whether the predicted error types are semantically consistent with the ground-truth annotations.



- Step1. Plan. Guide students to analyze the problem, identify what the problem is asking, and plan a strategy to solve the problem by themselves."
- Step2. Monitor. Prompt students to check their steps as they work and verify that their approach and reasoning make sense."
- Step3. Modify. Support students in revising their plan or choosing an alternative strategy if their current approach is ineffective, or if their answer seems unreasonable. Encourage students to explain why a modification is needed.

## B More Experimental Results and Analysis

### B.1 Complete Student Performance Prediction Results

#### B.1.1 Complete Overall Prediction Results

The complete prediction results of our method (GraphLR-MPP) compared with five ICL baselines and two additional SFT baselines on six students are reported in Tables 3-8. As shown, GraphLR-MPP consistently outperforms all baseline methods in terms of overall accuracy across all six students.

#### B.1.2 Complete Error-type Prediction Results

For an important aspect of our analysis, Table 9 reports the accuracy of error-type prediction for six students achieved by GraphLR-MPP and the baseline methods. To provide a more intuitive comparison, we further visualize the results using a heatmap in Figure 9. As shown, our method significantly outperforms the other approaches.

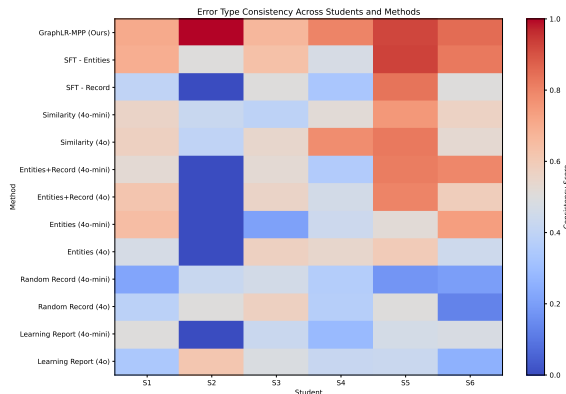


Figure 9: Heatmap of Error Type Predictions for Six Students.

Methods	S1	S2	S3	S4	S5	S6
SFT (Qwen2.5-7B-Instruct)						
GraphLR-MPP	<b>0.706</b>	<b>1.000</b>	<b>0.667</b>	<b>0.800</b>	<b>0.921</b>	<b>0.850</b>
Entities	<u>0.692</u>	0.500	<u>0.636</u>	0.474	<b>0.929</b>	<u>0.824</u>
Record	0.400	0.000	0.500	0.333	0.833	0.500
ICL (GPT-4o-mini)						
Similarity	0.556	0.429	0.393	0.519	0.750	0.563
Entities+Record	0.526	0.000	0.526	0.357	0.813	0.790
Entities	0.647	0.000	0.211	0.438	0.516	0.737
Random Record	0.222	0.500	0.462	0.364	0.179	0.200
Learning Report	0.500	0.000	0.432	0.286	0.467	0.478
ICL (GPT-4o)						
Similarity	0.571	0.400	0.546	<u>0.778</u>	0.822	0.531
Entities+Record	0.619	0.000	0.556	0.462	0.800	0.588
Entities	0.471	0.429	0.571	0.546	0.594	0.444
Random Record	0.385	0.000	0.571	0.368	0.500	0.129
Learning Report	0.340	<u>0.615</u>	0.487	0.426	0.432	0.257

Table 9: Comparison of Error Type Consistency Across Different Methods for Six Students, S1–S6 denote Student1-6.

### B.2 Intervention Results

#### B.2.1 Intervention case

Here, we provide examples of teachers employing Explicit, Systematic Instruction and Metacognitive Strategies to guide students in re-solving questions they previously answered incorrectly, as shown in Figures 10 and 11, respectively.

#### B.2.2 Cognitive Graph Update

We visualize the complete cognitive graph of Student1 before and after the metacognitive strategy intervention in Figure 12. The proportion of green nodes is noticeably higher than that of red nodes, indicating an improvement in the student's cognitive level after the intervention.

Methods		SFT				ICL								
		GraphLR-MPP	Entities	Record	Similarity	Entities+Record		Entities		Random Record		Learning Report		
Base Model		Qwen2.5-7B-Instruct		4o-mini	4o	4o-mini	4o	4o-mini	4o	4o-mini	4o	4o-mini	4o	
Overall	Acc.	<b>0.71</b>	0.645	0.655	0.649	0.595	0.615	<u>0.683</u>	0.6	0.688	0.645	0.634	0.551	0.439
Correct	Prec.	<b>0.765</b>	0.713	0.705	<u>0.746</u>	0.706	0.709	0.739	0.697	0.729	0.727	0.731	0.74	0.729
	Rec.	0.845	0.824	<b>0.873</b>	0.746	0.711	0.754	0.838	0.746	<b>0.873</b>	0.789	0.746	0.542	0.303
	F1	<b>0.803</b>	0.765	0.78	0.746	0.709	0.73	0.785	0.721	<u>0.795</u>	0.757	0.739	0.626	0.428
Wrong	Prec.	<b>0.519</b>	0.359	0.333	0.429	0.339	0.352	0.477	0.321	<u>0.486</u>	0.388	0.4	0.356	0.322
	Rec.	0.392	0.23	0.148	0.429	0.333	0.302	0.333	0.27	0.27	0.311	0.381	<u>0.571</u>	<b>0.746</b>
	F1	<u>0.446</u>	0.28	0.205	0.429	0.336	0.325	0.393	0.293	0.347	0.345	0.39	0.439	<b>0.45</b>

Table 3: Comparison of methods in predicting Student1’s performance. Best results are in bold, second best are underlined.

Methods		SFT				ICL								
		GraphLR-MPP	Entities	Record	Similarity	Entities+Record		Entities		Random Record		Learning Report		
Base Model		Qwen2.5-7B-Instruct		4o-mini	4o	4o-mini	4o	4o-mini	4o	4o-mini	4o	4o-mini	4o	
Overall	Acc.	<b>0.845</b>	0.801	0.825	0.757	0.728	0.825	<u>0.83</u>	0.825	<u>0.83</u>	0.796	0.796	0.607	0.447
Correct	Prec.	0.901	0.896	0.89	0.914	0.901	0.903	0.908	0.903	0.908	<b>0.918</b>	0.909	0.888	0.908
	Rec.	<b>0.93</b>	0.881	<u>0.919</u>	0.805	0.784	0.903	0.903	0.903	0.903	0.849	0.859	0.643	0.427
	F1	<b>0.915</b>	0.888	0.904	0.856	0.838	0.903	<u>0.905</u>	0.903	<u>0.905</u>	0.882	0.883	0.746	0.581
Wrong	Prec.	0.133	0.083	0	0.163	0.111	0.143	<u>0.182</u>	0.143	<u>0.182</u>	<b>0.2</b>	0.161	0.083	0.109
	Rec.	0.095	0.095	0	<u>0.333</u>	0.238	0.143	0.19	0.143	0.19	<u>0.333</u>	0.238	0.286	<b>0.619</b>
	F1	0.111	0.089	0	<u>0.219</u>	0.152	0.143	0.186	0.143	0.186	<b>0.25</b>	0.192	0.129	0.186

Table 4: Comparison of methods in predicting Student2’s performance. Best results are in bold, second best are underlined.

Methods		SFT				ICL								
		GraphLR-MPP	Entities	Record	Similarity	Entities+Record		Entities		Random Record		Learning Report		
Base Model		Qwen2.5-7B-Instruct		4o-mini	4o	4o-mini	4o	4o-mini	4o	4o-mini	4o	4o-mini	4o	
Overall	Acc.	<b>0.725</b>	0.667	0.676	<u>0.706</u>	0.676	0.691	0.696	0.686	0.681	0.627	0.672	0.461	0.48
Correct	Prec.	0.76	0.723	0.718	<u>0.789</u>	0.748	0.753	0.755	0.752	0.756	0.714	0.753	0.727	<b>0.791</b>
	Rec.	<b>0.897</b>	0.862	<b>0.897</b>	0.8	0.821	0.841	0.848	0.834	0.814	0.793	0.8	0.386	0.366
	F1	<b>0.823</b>	0.786	0.789	0.795	0.783	0.795	<u>0.799</u>	0.791	0.784	0.752	0.776	0.505	0.5
Wrong	Prec.	<b>0.545</b>	0.355	0.348	<u>0.491</u>	0.422	0.452	0.463	0.442	0.438	0.302	0.42	0.299	0.328
	Rec.	0.305	0.186	0.136	0.475	0.322	0.322	0.322	0.322	0.356	0.22	0.356	<u>0.644</u>	<b>0.763</b>
	F1	0.391	0.244	0.195	<b>0.483</b>	0.365	0.376	0.38	0.373	0.393	0.225	0.385	0.409	<u>0.459</u>

Table 5: Comparison of methods in predicting Student3’s performance. Best results are in bold, second best are underlined.

Methods		SFT				ICL								
		GraphLR-MPP	Entities	Record	Similarity	Entities+Record		Entities		Random Record		Learning Report		
Base Model		Qwen2.5-7B-Instruct		4o-mini	4o	4o-mini	4o	4o-mini	4o	4o-mini	4o	4o-mini	4o	
Overall	Acc.	<b>0.698</b>	0.593	0.633	<u>0.673</u>	0.618	0.648	<u>0.673</u>	0.668	0.658	0.653	0.658	0.603	0.387
Correct	Prec.	0.754	0.75	0.726	0.811	0.793	0.752	0.763	0.765	0.745	0.744	0.776	<u>0.854</u>	<b>0.857</b>
	Rec.	<b>0.878</b>	0.673	0.81	0.728	0.653	0.782	0.81	0.796	<u>0.816</u>	0.81	0.755	0.558	0.204
	F1	<b>0.811</b>	0.71	0.765	0.767	0.716	0.767	<u>0.785</u>	0.78	0.779	0.775	0.766	0.675	0.33
Wrong	Prec.	0.357	0.284	0.2	<b>0.403</b>	0.346	0.304	0.349	0.348	0.289	0.282	0.357	<u>0.369</u>	0.287
	Rec.	0.192	0.365	0.135	0.519	0.519	0.269	0.288	0.308	0.212	0.212	0.385	<u>0.731</u>	<b>0.904</b>
	F1	0.25	0.319	0.161	<u>0.454</u>	0.415	0.286	0.316	0.327	0.244	0.242	0.37	<b>0.49</b>	0.435

Table 6: Comparison of methods in predicting Student4’s performance. Best results are in bold, second best are underlined.

Methods		SFT				ICL								
		GraphLR-MPP	Entities	Record	Similarity	Entities+Record		Entities		Random Record		Learning Report		
Base Model		Qwen2.5-7B-Instruct		4o-mini	4o	4o-mini	4o	4o-mini	4o	4o-mini	4o	4o-mini	4o	
Overall	Acc.	<b>0.707</b>	0.608	0.656	0.649	<u>0.683</u>	0.605	0.62	0.649	0.649	0.62	0.566	0.588	0.629
Correct	Prec.	<b>0.699</b>	0.649	0.671	0.638	0.673	0.622	0.626	0.642	0.643	0.623	0.591	0.676	<u>0.685</u>
	Rec.	0.855	0.821	<b>0.889</b>	<b>0.889</b>	0.863	0.786	0.829	0.872	0.863	0.846	0.778	0.61	0.65
	F1	<b>0.769</b>	0.725	<u>0.765</u>	0.743	0.757	0.694	0.713	0.739	0.737	0.717	0.672	0.641	0.667
Wrong	Prec.	<b>0.726</b>	0.447	0.581	0.69	<u>0.709</u>	0.561	0.6	0.674	0.667	0.609	0.490	0.484	0.564
	Rec.	0.511	0.246	0.261	0.33	0.443	0.364	0.341	0.352	0.364	0.318	0.284	<u>0.556</u>	<b>0.602</b>
	F1	<b>0.6</b>	0.318	0.36	0.446	0.545	0.441	0.435	0.463	0.471	0.418	0.36	0.517	<u>0.582</u>

Table 7: Comparison of methods in predicting Student5’s performance. Best results are in bold, second best are underlined.

Methods		SFT				ICL								
		GraphLR-MPP	Entities	Record	Similarity	Entities+Record		Entities		Random Record		Learning Report		
Base Model		Qwen2.5-7B-Instruct		4o-mini	4o	4o-mini	4o	4o-mini	4o	4o-mini	4o	4o-mini	4o	
Overall	Acc.	<b>0.804</b>	0.749	<u>0.794</u>	0.739	0.709	<u>0.794</u>	0.774	0.779	0.754	0.678	0.678	0.608	0.427
Correct	Prec.	0.84	0.825	0.823	<b>0.898</b>	<u>0.893</u>	0.842	0.83	0.84	0.839	0.817	0.836	0.865	0.833
	Rec.	<u>0.922</u>	0.857	<b>0.935</b>	0.747	0.708	0.903	0.89	0.883	0.844	0.753	0.727	0.584	0.325
	F1	<b>0.879</b>	0.841	<u>0.875</u>	0.816	0.79	0.871	0.859	0.861	0.841	0.784	0.778	0.698	0.467
Wrong	Prec.	<b>0.6</b>	0.436	<u>0.583</u>	0.451	0.416	0.559	0.5	0.514	0.455	0.333	0.354	0.326	0.252
	Rec.	0.4	0.378	0.311	<u>0.711</u>	<u>0.711</u>	0.422	0.378	0.422	0.444	0.422	0.511	0.689	<b>0.778</b>
	F1	0.48	0.405	0.406	<b>0.552</b>	<u>0.525</u>	0.481	0.43	0.463	0.449	0.373	0.418	0.443	0.38

Table 8: Comparison of methods in predicting Student6’s performance. Best results are in bold, second best are underlined.



**Problem:** pipe a and pipe b fill water into a tank of capacity 5000 litres , at a rate of 200 l / min and 50 l / min . pipe c drains at a rate of 25 l / min . pipe a is open for 1 min and closed , then pipe b is open for 2 min and closed . further the pipe c is opened and drained for another 2 min . this process is repeated until the tank is filled . how long will it take to fill the tank ?  
Right answer: 100

**Student answer:** 20  
**Student process:** In every five minutes,  $200 + 50 - 25 = 225$  litres of water are added.  $5000 \div 225 = 22.22$  minutes.  $22.22 \times 5 = 111.11$  minutes.  
**Review:**  
**Correct or not: wrong**  
**Error counts:** 1  
**Error type:** Wrong mathematical operation/concept  
**Error equation:**  $5000 \div (200 + 50 - 25) = 22.22$   
**Teacher advice:** The teacher suggests that 250 is the amount of water changed every 5 minutes, so 20 should be multiplied by 5, which is 100, to get the correct total time spent.

**Step1. Plan**

Let's focus on what the problem is asking: how long it will take to fill the tank. You already calculated the amount of water changed every 5 minutes, which is 250 liters. Now, think about how many sets of 5-minute cycles are needed to fill the 5000-liter tank. **What should you do next to find the total time?**

To find the total time, I need to calculate how many 5-minute cycles it takes to fill the 5000-liter tank. Since each cycle adds 250 liters, I'll divide the total capacity by the amount added per cycle:  $5000 \div 250 = 20$ . So, it takes 20 cycles. **Since each cycle is 5 minutes long, the total time is  $20 \times 5 = 100$  minutes.**

**Step2. Monitor**

Great job! Now, take a moment to double-check if your answer makes sense. Does 100 minutes seem reasonable based on the process described in the problem? Consider the rates of filling and draining. Does the time align with the expected outcome?

Yes, 100 minutes seems reasonable. The pipes are filling the tank at a good rate, but pipe c is also draining water. Since the cycles involve both filling and draining, 100 minutes aligns with the expected outcome. The calculation also matches the tank's total capacity.

**Step3. Modify**

Your reasoning is spot on, and your calculation is correct! The answer of 100 minutes makes sense given the rates and the process described in the problem. **There's no need to modify your approach.** Well done!

Thank you! I'm glad the solution and reasoning are correct.

Figure 10: A Case of Teacher-Led *Metacognitive Strategy Intervention* in Students' Mathematical Problem Solving

**Problem:** two trains of equal length are running on parallel lines in the same direction at 46 km / hr and 36 km / hr . the faster train catches and completely passes the slower train in 72 seconds . what is the length of each train ( in meters ) ?  
Right answer: 100

**Student answer:** 20  
**Student process:**  $46 - 36 = 10$  km/hr.  $10 \times 72 = 720$  meters.  $720 \div 2 = 360$  meters.  
**Review:**  
**Correct or not: wrong**  
**Error counts:** 1  
**Error type:** Wrong mathematical operation/concept  
**Teacher advice:** Train A needs to pass the length of train B and its own length to completely pass train B. Since the trains are of equal length, the length of each train is  $200 \div 2 = 100$ .

**Step1. Orientation to the Lesson**

Let's start by reviewing how relative speed is used when two objects are moving in the same direction. Recall that the relative speed is the difference between their speeds. Also, remember that when one object passes another, it needs to cover the complete length of both objects to fully pass.

Got it! So, the relative speed is the difference between the speeds of the two trains:  $46 - 36 = 10$  km/hr. When one train passes the other, it must cover the length of both trains.

**Step2. Initial Instruction**

The challenging part of this problem is understanding that the faster train covers the length of both trains to completely pass the slower one. Let's ensure you know how to calculate the distance covered using the relative speed and time. Convert the relative speed from km/hr to m/s first, then use the time in seconds to find the total distance. Does this operation make sense to you?

Yes, that makes sense! I'll convert the relative speed from km/hr to m/s, then use the time in seconds to find the total distance.

**Step3. Teacher-Guided Practice**

Great! Let's begin. First, we need to convert the relative speed of 10 km/hr to meters per second. To do that, divide by 3.6. What do you get?

$10 \div 3.6 = 2.78$  m/s

Correct! Now, let's use this relative speed to find the total distance covered in 72 seconds. Multiply the relative speed of 2.78 m/s by the time of 72 seconds. What is the total distance the faster train traveled to pass the slower one?

$2.78 \times 72 = 200.16$  meters

Great! Now, remember that this distance is the total length the faster train needs to cover, which includes both trains. Since the trains are of equal length, we need to divide this distance by 2 to find the length of one train. What is the length of one train?

$200.16 \div 2 = 100.08$  meters. So, the length of one train is approximately 100 meters.

**Step4. Independent Practice**

Excellent work! Now, I want you to try solving the original problem independently, applying the steps we've discussed. Start by calculating the relative speed between the trains, converting it to meters per second, and using the given time to find the total distance. Once you have that, determine the length of each train. Let me know your answer when you're ready!

Sure! Here's my solution. The relative speed is  $46 - 36 = 10$  km/hr. Converting to meters per second:  $10 \div 3.6 = 2.78$  m/s. The total distance covered in 72 seconds is  $2.78 \times 72 = 200.16$  meters. The length of each train is  $200.16 \div 2 = 100.08$  meters. So, the length of each train is approximately 100 meters.

Figure 11: A Case of Teacher-Led *Explicit, Systematic Instruction* in Students' Mathematical Problem Solving

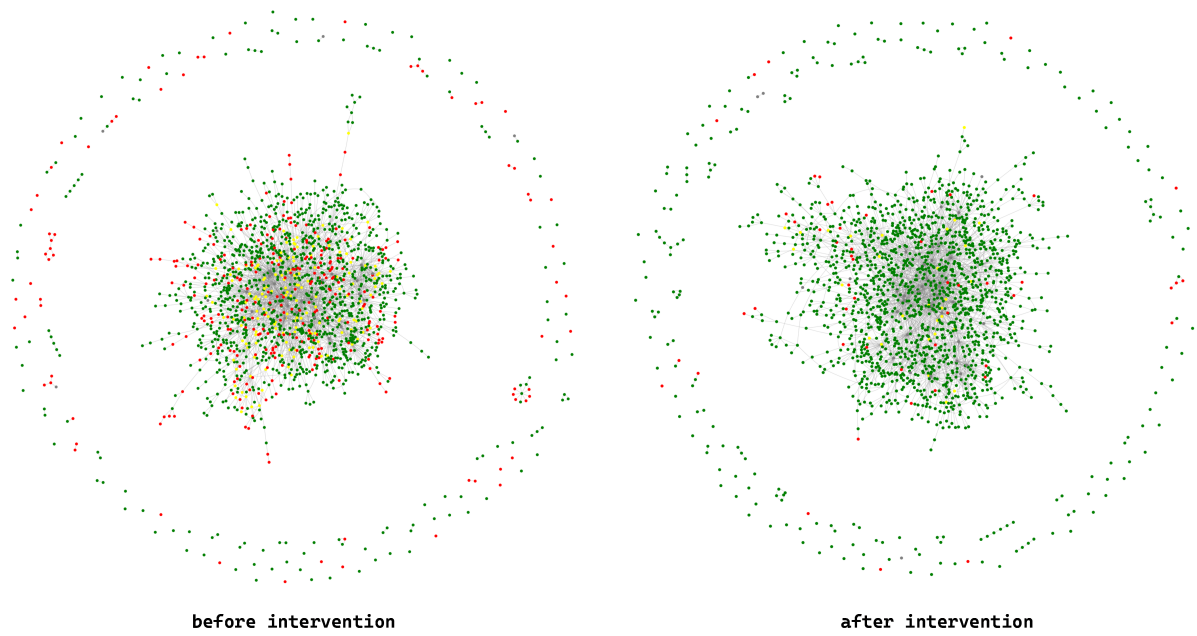


Figure 12: Complete Cognitive Graphs of Students Before and After Metacognitive Strategy Intervention. Red indicates not mastered, green indicates mastered, and yellow indicates partially mastered.