

# One LLM Does Not Simulate All Students: Ability-Aware Student Simulation via Cognitive Diagnosis Guided LLM Assignment

Anonymous ACL submission

## Abstract

Large Language Models (LLMs) have become integral to personalized education systems, particularly in the realm of student behavior simulation. By predicting fine-grained learning behaviors, these simulations enable intelligent systems to provide tailored instructional support. However, most existing methods rely on a single high-capacity LLM to represent an entire population of diverse learners. In this work, we demonstrate that this “one-size-fits-all” approach induces a systematic *ability-dependent bias*, where high-capacity models tend to overestimate low-ability students while lower-capacity models underestimate high-ability ones. To mitigate this distortion, we propose an **ability-aware student simulation framework** that dynamically matches students with appropriate LLM backbones through cognitive alignment. We leverage Neural Cognitive Diagnosis (NeuralCD) to extract multidimensional cognitive profiles for both human students and LLM agents within a shared skill space, subsequently pairing each student with the most cognitively representative model. Extensive experiments demonstrate that our approach substantially reduces simulation bias and consistently outperforms single-model baselines across the entire proficiency spectrum. Our findings suggest that faithful behavior simulation necessitates the **alignment of model capacity with student ability**, establishing cognitive diagnosis as a principled mechanism for model assignment in educational AI.

## 1 Introduction

Personalized learning is widely recognized as an effective paradigm for improving learning outcomes by adapting instructional content and strategies to individual students (Bernacki et al., 2021; Shemshack and Spector, 2020; Jian, 2023). With the rapid advancement of large language models (LLMs), recent studies have increasingly incorporated LLMs into personalized learning systems

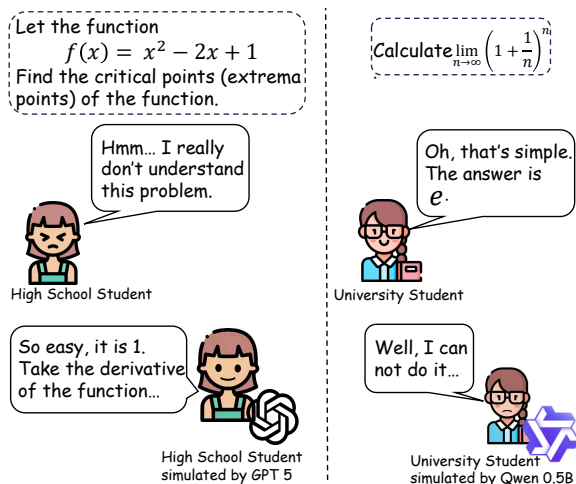


Figure 1: Different ways to simulate students. The left shows high-capacity LLMs simulating low cognitive ability students, while the right shows low-capacity LLMs simulating high cognitive ability students.

(Zhang et al., 2025; Wen et al., 2024; Neumann et al., 2024). Among these applications, *student behavior simulation* has attracted growing attention (Gao et al., 2025; Lu and Wang, 2024; Xu et al., 2024), where LLMs act as agents to predict students’ future learning behaviors from historical interactions.

Despite promising empirical results, existing LLM-based student simulation methods suffer from a fundamental limitation. Most prior work relies on a *single LLM*—typically a high-capacity model such as ChatGPT (Xu et al., 2025; Abasiantaeb et al., 2024)—to simulate all students. This approach assumes that a strong LLM can faithfully represent learners across the entire ability spectrum, which is inherently problematic. Due to their pre-trained knowledge, LLMs exhibit an uncontrollable *competence prior* that can leak into simulations and undermine cognitive fidelity. As illustrated in Figure 1, when simulating a struggling high school student solving a basic math problem, a powerful model such as GPT-5 tends to produce

066	correct answers, even in cases where the student’s	<b>Contributions</b>	Our main contributions are sum-	118
067	historical interaction data suggest that the relevant		marized as follows:	119
068	concept has not been mastered. Conversely, low-	• <b>Systematic Bias Analysis:</b>	We empirically	120
069	capacity LLMs frequently lack the domain knowl-		demonstrate that using a single LLM to simu-	121
070	edge required to simulate advanced learners. This		late diverse students induces systematic abil-	122
071	mismatch leads to systematic ability bias, overes-		ity bias, failing to capture the heterogeneity of	123
072	timating low-ability students while underestimat-		learners.	124
073	ing high-ability ones, thereby distorting simulated	• <b>Ability-Aware Assignment Framework:</b>	We	125
074	learning behaviors.		propose a novel framework that leverages	126
075	Given the diverse landscape of LLMs, ranging		NeuralCD to facilitate student–LLM match-	127
076	from lightweight to large-scale systems (Bai et al.,		ing, moving beyond task-driven selection to-	128
077	2023; Achiam et al., 2023; Liu et al., 2024; Touvron		ward <i>cognitive-ability alignment</i> .	129
078	et al., 2023; Team et al., 2023), a natural alterna-	• <b>Improved Simulation Fidelity:</b>	Extensive ex-	130
079	tive is to leverage a heterogeneous pool of LLMs		periments show that our approach consistently	131
080	to represent students with varying cognitive abil-		outperforms single-LLM baselines, achieving	132
081	ities. This raises a fundamental question: <b>How can</b>		higher simulation fidelity.	133
082	<b>a student’s latent cognitive state be systemati-</b>	<b>2 Related Work</b>		134
083	<b>cally aligned with the most representative LLM</b>	<b>2.1 LLM-based Agent Simulation in</b>		135
084	<b>backbone to ensure a faithful simulation?</b> While	<b>Education</b>		136
085	research on <i>LLM assignment</i> (Ding et al., 2024; Xia		The simulation of LLM-based agents is gradually	137
086	et al., 2024; Feng et al., 2024) is growing, it primar-		gaining momentum (Park et al., 2023; Man et al.,	138
087	ily focuses on task-driven optimization, aiming to		2025; Yang et al., 2025). In the educational domain,	139
088	find a cost-effective model combination that max-		these simulations have proliferated to encompass	140
089	imizes task success within a budget (Panda et al.,		diverse scenarios such as learning process simula-	141
090	2025; Song et al., 2025; Ashury Tahan et al., 2024).		tion (Mannekote et al., 2025; Gao et al., 2025; Xu	142
091	This objective is misaligned with educational simu-		et al., 2024) and pedagogical interactions (Zheng	143
092	lation, where the goal is not <i>utility-optimal perfor-</i>		et al., 2025; Lv et al., 2025). To achieve human-like	144
093	<i>mance</i> but <i>cognitive ability alignment</i> . A faithful		behavior, most existing agent-based simulations	145
094	simulation must reflect a student’s specific cog-		adopt a <i>modular architecture</i> paradigm (Chu et al.,	146
095	gnitive constraints, persistent misconceptions, and		2025; Bhowmik et al., 2024). Researchers typically	147
096	actual skill mastery, even when it leads to incorrect		employ a single, high-ability LLM as the central	148
097	answers. Therefore, the selection criterion must		brain and augment it with specialized components,	149
098	shift from maximizing task success to minimizing		such as memory modules, planning modules and	150
099	the cognitive ability gap between the LLM agent		reflection components for self-correction (Zheng	151
100	and the human learner.		et al., 2025; Arana et al., 2025). However, these	152
101	To address this challenge, we propose an <b>ability-</b>		structural augmentations do not mitigate the in-	153
102	<b>aware student simulation framework</b> grounded		herent competence prior of the underlying LLMs,	154
103	in <b>Cognitive Diagnosis (CD)</b> theory (Templin and		often leading to simulated behaviors that remain	155
104	Henson, 2006; Leighton and Gierl, 2007). Rooted		unaligned with the actual cognitive state of students.	156
105	in psychometrics, CD models infer learners’ la-		In contrast, our work shifts the focus from struc-	157
106	tent cognitive states by mapping observable behav-		tural modularity to <b>ability-aware alignment</b> . In-	158
107	iors to fine-grained skill mastery. Building on this		stead of relying on a single model with complex	159
108	foundation, we employ <i>Neural Cognitive Diagnosis</i>		external modules, we dynamically select an appro-	160
109	( <i>NeuralCD</i> ) (Wang et al., 2022), a representative		prate LLM for each student based on their diag-	161
110	deep learning-based CD model, to estimate stu-		gnosed cognitive profile. This approach ensures that	162
111	dents’ cognitive ability profiles from historical in-		the simulation is grounded in the intrinsic capacity	163
112	teraction data, while simultaneously deriving skill-		of the agent itself, thereby facilitating a more au-	164
113	level performance profiles for candidate LLMs. By		thentic reflection of student-specific behaviors and	165
114	matching students and LLMs based on their sim-		learning outcomes.	166
115	ilarity, we assign each student the LLM that best			
116	reflects their cognitive ability, rather than their task-			
117	solving potential.			

## 2.2 Cognitive Diagnosis Models

Cognitive Diagnosis (CD) models aim to infer the latent cognitive states of learners, particularly within the context of educational assessment. Classical paradigms including DINA (De La Torre, 2009), IRT (Lord, 2012), and MIRT (Reckase, 2009) estimate student abilities through probabilistic frameworks. While theoretically well-founded, these methods often rely on predefined parametric functions such as 1PL and 2PL models (DeMars, 2010). This reliance constrains their flexibility and scalability when processing complex or large-scale educational datasets. To address these limitations, deep learning-based CD models, such as NeuralCD (Wang et al., 2022) and RCD (Gao et al., 2021), utilize neural networks to capture complex student-exercise interactions from large-scale response data. Building on this paradigm, we employ NeuralCD to map both human students and LLM agents into a unified latent space by treating LLMs as “artificial learners”. This shared cognitive dimensionality allows us to directly quantify the mastery gap between students and models, facilitating precise ability-aware alignment.

## 2.3 Large Language Model Assignment

Research on LLM model assignment has gained momentum as a strategy to balance inference costs with task performance (Mei et al., 2025; Song et al., 2025; Ding et al., 2024). These methods typically function as routers that assign queries to appropriate models based on task difficulty. Such approaches are primarily task-driven (Dai et al., 2024; Zhao et al., 2024; Ashury Tahan et al., 2024), prioritizing the minimization of computational overhead or latency while maintaining a specific threshold of solution quality (Hu et al., 2024; Jitkrittum et al., 2025).

Our work differs fundamentally from this paradigm by shifting the focus from task-oriented routing to profile-oriented alignment for simulation. Rather than optimizing for efficiency, we prioritize behavioral fidelity by matching the intrinsic capabilities of an LLM with the diagnosed cognitive profiles of the target learners.

## 3 Method

### 3.1 Task Formulation

Let  $\mathcal{S}$  denote the set of students,  $\mathcal{E}$  the set of exercises, and  $\mathcal{K}$  the set of knowledge concepts. For each student  $s \in \mathcal{S}$ , the histor-

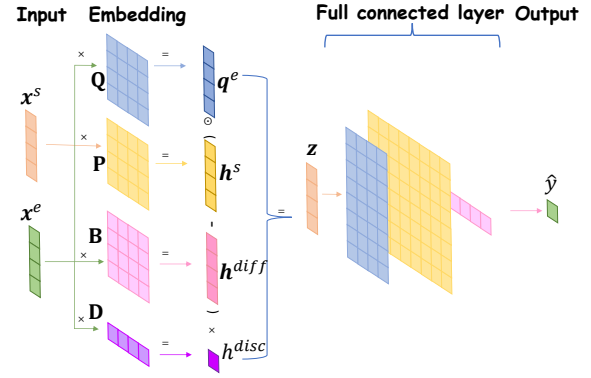


Figure 2: The architecture of NeuralCD.

ical interaction sequence is denoted as  $I_s = \{(e_1, y_{s,e_1}), \dots, (e_n, y_{s,e_n})\}$ , where each exercise  $e_i$  is a triplet consisting of textual content  $e_{i,\text{text}}$ , associated concepts  $e_{i,\text{concept}} \subseteq \mathcal{K}$ , and the ground-truth answer  $e_{i,\text{ans}}$ . The response  $y_{s,e_i} \in \{0, 1\}$  indicates whether the student’s answer was correct.

In this work, we redefine the student simulation task as a cognitive-aligned model assignment problem. Beyond the student data, we introduce a heterogeneous pool of Large Language Models  $\mathcal{M} = \{m_1, m_2, \dots, m_k\}$  with varying capacities. Our objective is to develop a framework that consists of two key components:

**Cognitive Profiling:** A diagnostic function  $f_D : I_s \rightarrow \alpha_s$  that maps a student’s history into a latent cognitive profile  $\alpha_s \in [0, 1]^{|\mathcal{K}|}$ , representing their mastery levels across all concepts.

**Ability-Aware Assignment:** An alignment function  $f_A : (\alpha_s, \mathcal{M}) \rightarrow m^*$  that selects the optimal LLM  $m^* \in \mathcal{M}$  whose intrinsic capability profile most closely mirrors the student’s cognitive state  $\alpha_s$ .

The final goal is to utilize the assigned model  $m^*$  to simulate the student’s future behaviors. A successful simulation ensures that the response  $\hat{y}_{s,e_{new}}$  generated by  $m^*$  is not only accurate in terms of prediction but also reflects the underlying cognitive constraints and misconceptions of student  $s$ .

### 3.2 Neural Cognitive Diagnosis

To project both human students and LLM agents into a unified mastery space, we employ Neural Cognitive Diagnosis (NeuralCD) (Wang et al., 2022) as our underlying diagnostic engine. This model allows us to transform observable response data into latent, multidimensional ability profiles. The NeuralCD model architecture is show in Figure 2.

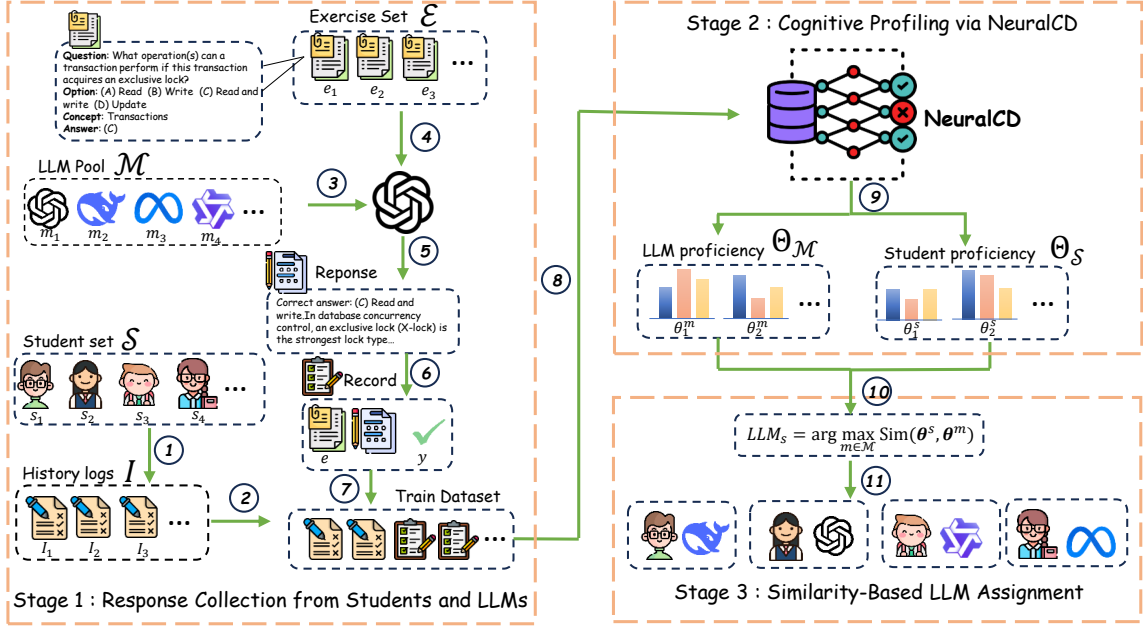


Figure 3: The overview of adaptive LLM assignment for student simulation.

Each student  $s$  is encoded as a binary *one-hot* vector  $\mathbf{x}^s \in \{0, 1\}^{1 \times |S|}$ , where exactly one entry is set to 1 to indicate the student's identity. Based on this representation, NeuralCD learns a student-specific knowledge concept proficiency vector  $\mathbf{h}^s \in (0, 1)^{1 \times |\mathcal{K}|}$  via:

$$\mathbf{h}^s = \sigma(\mathbf{x}^s \mathbf{P}) \quad (1)$$

where  $\sigma(\cdot)$  denotes the sigmoid function and  $\mathbf{P} \in \mathbb{R}^{|S| \times |\mathcal{K}|}$  is a trainable parameter matrix.

Similarly, each exercise  $e \in \mathcal{E}$  is represented by a binary *one-hot* vector  $\mathbf{x}^e \in \{0, 1\}^{1 \times |\mathcal{E}|}$ . We further adopt a predefined Q-matrix  $\mathbf{Q} \in \{0, 1\}^{|\mathcal{E}| \times |\mathcal{K}|}$  to encode the associations between exercises and knowledge concepts, where  $\mathbf{Q}_{ij} = 1$  indicates that exercise  $i$  involves the  $j$ -th knowledge concept. Using the Q-matrix, the concept-level representation of exercise  $e$  is computed as:

$$\mathbf{q}^e = \mathbf{x}^e \mathbf{Q} \quad (2)$$

To characterize exercise properties, NeuralCD estimates both *concept-level difficulty* and *exercise discrimination*. Specifically, the difficulty vector  $\mathbf{h}^{diff} \in (0, 1)^{1 \times |\mathcal{K}|}$  and the discrimination scalar  $h^{disc} \in (0, 1)$  are computed as:

$$\mathbf{h}^{diff} = \sigma(\mathbf{x}^e \mathbf{B}), \quad h^{disc} = \sigma(\mathbf{x}^e \mathbf{D}) \quad (3)$$

where  $\mathbf{B} \in \mathbb{R}^{|\mathcal{E}| \times |\mathcal{K}|}$  and  $\mathbf{D} \in \mathbb{R}^{|\mathcal{E}| \times 1}$  are trainable parameter matrices.

Following MIRT (Chalmers, 2012), NeuralCD integrates student proficiency and exercise characteristics to construct the input to the prediction network:

$$\mathbf{z} = \mathbf{q}^e \odot (\mathbf{h}^s - \mathbf{h}^{diff}) \times h^{disc} \quad (4)$$

where  $\odot$  denotes element-wise product.

The resulting vector  $\mathbf{z}$  is then fed into a multi-layer fully connected neural network:

$$\mathbf{f}_1 = \phi(\mathbf{W}_1 \mathbf{z}^T + \mathbf{b}_1) \quad (5)$$

$$\mathbf{f}_2 = \phi(\mathbf{W}_2 \mathbf{f}_1 + \mathbf{b}_2) \quad (6)$$

$$\hat{y} = \phi(\mathbf{W}_3 \mathbf{f}_2 + \mathbf{b}_3) \quad (7)$$

where  $\phi(\cdot)$  denotes the activation function and  $\hat{y}$  represents the predicted probability that student  $s$  correctly answers exercise  $e$ .

The model is trained by minimizing the binary cross-entropy loss:

$$\mathcal{L} = - \sum_i [r_i \log \hat{y}_i + (1 - r_i) \log(1 - \hat{y}_i)] \quad (8)$$

where  $r_i \in \{0, 1\}$  denotes the ground-truth response correctness.

After training, the learned parameter matrix  $\mathbf{P}$  enables us to infer each student's knowledge concept proficiency vector  $\mathbf{h}^s$  via Equation 1, which serves as the cognitive profile for subsequent LLM assignment and student behavior simulation.

### 3.3 Adaptive LLM Assignment for Student Simulation

The core component of our framework is an **adaptive LLM assignment module**, which assigns each student an appropriate LLM to serve as the simulation backbone. Instead of relying on a single LLM to simulate all students, our approach explicitly accounts for heterogeneity in student proficiency by matching students with LLMs of comparable cognitive profiles.

The assignment procedure consists of three sequential stages: (1) collecting response records from both students and LLMs, (2) estimating their knowledge concept proficiency vectors via NeuralCD, and (3) matching each student with the most suitable LLM based on proficiency similarity.

**Stage 1: Response Collection from Students and Candidate LLMs.** For students, we directly collect their historical response records. In addition, we define a candidate set of LLMs  $\mathcal{M}$  to serve as potential simulation backbones. Since LLMs do not possess prior interaction histories, we prompt each LLM  $m \in \mathcal{M}$  to answer all exercises in the dataset. We employ a chain-of-thought (CoT) prompting strategy (Wei et al., 2022) to elicit structured reasoning and record the final answer correctness. As a result, we obtain a train dataset that includes both real student responses and LLM-generated responses over the same set of exercises.

**Stage 2: Cognitive Profiling via NeuralCD.** Given the combined response records from students and candidate LLMs, we train a NeuralCD model to infer latent cognitive profiles. By treating LLMs as pseudo-students during training, NeuralCD provides a unified and interpretable representation of both human learners and LLMs under the same diagnostic framework. To maintain consistency with the parameterization of cognitive diagnosis models (Baker, 2001; Liu et al., 2023), we denote the student ability representation  $h^s$  produced by NeuralCD as  $\theta$ . Specifically, NeuralCD produces a student proficiency set  $\Theta_S = \{\theta^{s_1}, \theta^{s_2}, \dots, \theta^{s_{|S|}}\}$  and an LLM proficiency set  $\Theta_{\mathcal{M}} = \{\theta^{m_1}, \theta^{m_2}, \dots, \theta^{m_{|\mathcal{M}|}}\}$ .

**Stage 3: Similarity-Based LLM Assignment.** For each student  $s$ , we compare the student’s proficiency vector with those of all candidate LLMs in  $\mathcal{M}$  and select the most similar one LLM $_s$  as the simulation backbone. We adopt a similarity-based

criterion to select LLM $_s$ :

$$\text{LLM}_s = \arg \max_{m \in \mathcal{M}} \text{Sim}(\theta^s, \theta^m) \quad (9)$$

where  $\theta^s$  and  $\theta^m$  denote the knowledge concept proficiency vectors of student  $s$  and LLM  $m$ , respectively, and  $\text{Sim}(\cdot)$  represents a similarity function. In this work, we instantiate  $\text{Sim}(\cdot)$  as cosine similarity.

This assignment strategy ensures that each student is simulated by an LLM whose cognitive proficiency profile is most compatible with the student’s own abilities.

### 3.4 Student Simulation Agent

As the primary focus of this work is *adaptive LLM selection* rather than student agent design, we adopt an LLM-based student simulation framework from prior studies (Gao et al., 2025; Mannekote et al., 2025; Xu et al., 2024). Each student agent consists of a memory module, an action module, and a reflection module. Details of the agent framework are provided in Appendix A.

**Memory Module.** The memory module includes both short-term and long-term memory. Short-term memory stores the student’s recent response records. Long-term memory contains the student’s ability profile inferred via NeuralCD as well as reinforced short-term memories. If short-term records that occur more frequently than a predefined threshold are promoted into long-term memory.

**Action Module.** The action module models the student’s exercise-solving behavior. Before answering an exercise, the agent retrieves relevant records from short-term memory and further extracts related long-term memories based on similarity. This process includes identifying the knowledge concepts involved in the exercise to assist problem solving. The agent then generates an answer and performs a self-assessment to judge whether the response is correct.

**Reflection Module.** The reflection module enables the student agent to reflect on its behavior by comparing the LLM-generated response with the ground-truth student response. Reflection is triggered only when the two are inconsistent. In such cases, the reflection outcome, together with the response record, is written into short-term memory.

By adopting a standard student agent architecture, our framework isolates the effect of adap-

tive LLM selection, ensuring that performance differences arise from the suitability of the selected LLMs rather than agent design choices.

## 4 Experiments

### 4.1 Dataset Description

We conduct experiments on the public DBE-KT22 dataset (Abdelrahman et al., 2022), which contains rich educational interaction records including exercise texts, associated knowledge concepts, and student responses. The dataset consists of 1,361 students, 212 exercises, and 98 knowledge concepts, with a total of 167,222 student–exercise interaction records. Such characteristics make DBE-KT22 suitable for evaluating student behavior simulation under heterogeneous cognitive profiles.

### 4.2 Experiment Set Up

**NeuralCD Training Configuration.** For each student, their interaction records are chronologically split, with the first 80% used to train the NeuralCD model. During training, all LLM-generated responses corresponding to these exercises are included in the training set, allowing the model to learn from both student behaviors and LLM outputs. The remaining 20% of student interactions are reserved as ground-truth for evaluating LLM simulation. This setup ensures that the evaluation is performed on unseen student behaviors, avoiding information leakage. Detailed training hyperparameters are reported in Appendix B.

**LLM Inference Settings.** For all LLMs, we set the temperature to 0 to ensure reproducibility, except for GPT-5-Mini<sup>1</sup>, which does not support a configurable temperature and uses a fixed default temperature of 1.

**Student Selection for LLM Simulation.** Due to the computational and cost constraints associated with large-scale LLM-based simulation, we conduct experiments on a carefully selected subset of students. For every student  $s$ , we compute the mean of  $\theta^s$  as student’s overall proficiency score. Then we partition students into three groups according to its empirical quantiles: students whose scores fall within the 0–33rd percentile are categorized as low-proficiency, those within the 33–66th percentile as medium-proficiency, and those within the 66–100th percentile as high-proficiency. From

<sup>1</sup>GPT-5-Mini (2025-08-07). The specific model types and configurations are detailed in Appendix D.

each group, we randomly sample 100 students, resulting in a total of 300 students for evaluation. The IDs of the selected students are reported in Appendix C for reproducibility.

### 4.3 LLM Pool and Baselines

We construct an LLM pool consisting of 35 language models with diverse capabilities. Detailed model configurations are provided in Appendix D.

Specifically, we evaluate three representative settings: (1) a *single strong LLM* with high model capacity, (2) a *single weak LLM* with limited capacity, and (3) our proposed *NeuralCD-guided multi-LLM selection framework*. By comparing the bias patterns exhibited by these settings, we aim to verify both the limitations of single-LLM simulation and the effectiveness of adaptive LLM assignment.

### 4.4 Evaluation Metrics

We evaluate simulation quality by comparing LLM-generated responses with students’ ground-truth answers. A naive accuracy metric, however, is insufficient due to the imbalance between correct and incorrect student responses.

Let  $\mathcal{D}_{\text{correct}}^s$  and  $\mathcal{D}_{\text{incorrect}}^s$  denote the sets of exercises that a student answered correctly and incorrectly, respectively. Similarly, for the student Agent (LLM simulation), let  $\mathcal{D}_{\text{correct}}^m$  and  $\mathcal{D}_{\text{incorrect}}^m$  denote the sets of exercises answered correctly and incorrectly by the Agent.

Based on these sets, the conditional accuracies are computed as

$$Acc^+ = \frac{|\mathcal{D}_{\text{correct}}^s \cap \mathcal{D}_{\text{correct}}^m|}{|\mathcal{D}_{\text{correct}}^s|} \quad (10)$$

$$Acc^- = \frac{|\mathcal{D}_{\text{incorrect}}^s \cap \mathcal{D}_{\text{incorrect}}^m|}{|\mathcal{D}_{\text{incorrect}}^s|} \quad (11)$$

Here,  $Acc^+$  measures the proportion of exercises that the Agent correctly answers among the exercises that the student answered correctly, and  $Acc^-$  measures the proportion of exercises that the Agent incorrectly answers among the exercises that the student answered incorrectly.

To evaluate student simulation while accounting for bias between correct and incorrect responses, we adopt the *Bias-Aware Accuracy (BAA)*, inspired by Balanced Accuracy (Brodersen et al., 2010):

$$BAA = \frac{Acc^+ + Acc^-}{2} \cdot (1 - |Acc^+ - Acc^-|) \quad (12)$$

An effective student simulation framework should not only achieve high overall accuracy but also avoid systematic bias towards predicting either correct or incorrect responses. **The first term** of the metric captures the overall simulation accuracy across all response types, while **the second term** serves as an explicit penalty for imbalance between correct and incorrect predictions. When the accuracies of the two response types diverge, the penalty term decreases accordingly, reflecting degraded simulation quality.

## 5 Result

### 5.1 Bias Analysis of Single LLM Simulation

To better understand the intrinsic simulation bias of LLMs, we analyze how model capacity affects the accuracy of simulating correct and incorrect student responses. Specifically, we compare  $Acc^+$  and  $Acc^-$  across LLMs with different parameter scales<sup>2</sup>. We also conduct case study in E.

As shown in Figure 4, we observe a clear capacity-dependent trend. As model size increases—corresponding to stronger model capability— $Acc^+$  consistently improves. In contrast,  $Acc^-$  steadily decreases. This indicates that LLM-based student simulation is strongly influenced by the intrinsic capability of the underlying model.

Intuitively, larger and more capable LLMs tend to overestimate students’ mastery and are biased toward producing correct answers, making them less effective at reproducing incorrect student behaviors. Conversely, smaller LLMs are more likely to generate incorrect responses, resulting in higher  $Acc^-$  but lower  $Acc^+$ . Therefore, selecting either a weak or a strong LLM as a single simulation backbone inevitably introduces bias: weaker models better capture incorrect behaviors, while stronger models better capture correct behaviors.

We further compare LLM-based simulation results with the LLMs’ ground-truth (we calculate it through Section 3.3). The ground truth is defined as the results of the LLM’s direct answers to all non-duplicated exercises involved in the student simulation. Specifically, the ground truth corresponds to the answering accuracy in  $Acc^+$  table and the answering error rate in  $Acc^-$  table. We observe that  $Acc^+$  does not exceed the model’s

<sup>2</sup>Qwen2.5-3B-Instruct, Qwen2.5-7B-Instruct, Qwen2.5-14B-Instruct, LLaMA-3.2-1B-Instruct, LLaMA-3.2-3B-Instruct, LLaMA-3.1-8B-Instruct. The specific model types and configurations are detailed in Appendix D.

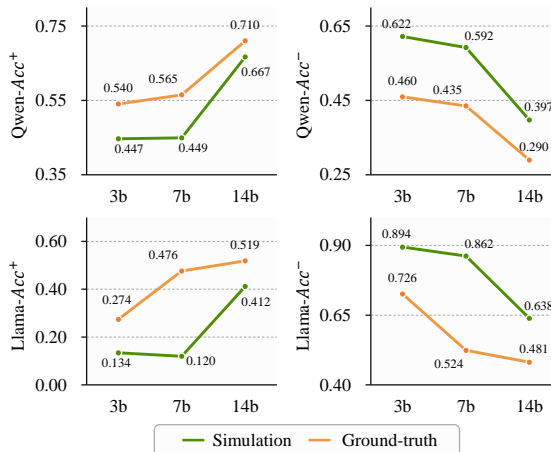


Figure 4:  $Acc^+$  and  $Acc^-$  w.r.t. ground truth across model scales for Qwen and LLaMA.

Table 1: Simulation performance across different student ability levels.

Model	$BAA_{low}$	$BAA_{mid}$	$BAA_{high}$	Overall
<i>Weak LLMs</i>				
Qwen-3B	0.427	0.417	0.463	0.441
Qwen-7B	0.476	0.440	0.439	0.446
Qwen-14B	0.360	0.398	0.394	0.389
LLaMA-1B	0.121	0.120	0.126	0.123
LLaMA-3B	0.131	0.124	0.126	0.127
LLaMA-8B	0.441	0.428	0.377	0.406
<i>Strong LLMs</i>				
DeepSeek	0.242	0.297	0.321	0.296
GPT-5-Mini	0.142	0.174	0.172	0.167
<i>CD-based Method</i>				
<b>Ours</b>	<b>0.523</b>	<b>0.516</b>	<b>0.475</b>	<b>0.505</b>

own ground truth, while  $Acc^-$  is higher than the corresponding ground truth. This suggests that an LLM’s intrinsic capability upper-bounds its simulation performance.

Overall, these findings imply that student simulation bias originates from the inherent capability constraints of LLMs. A single LLM cannot simultaneously and faithfully simulate students across diverse proficiency levels, as its own strengths and weaknesses are inevitably reflected in the simulated behaviors.

### 5.2 Effectiveness of Adaptive LLM Assignment

We further demonstrates the effectiveness of our proposed adaptive LLM assignment strategy. We report  $BAA$  scores for students with low, medium,

Table 2: Detailed  $Acc^+$  and  $Acc^-$  of representative LLMs across different student ability groups, where  $\Delta = |Acc^+ - Acc^-|$ .

Model	Low Ability			Medium Ability			High Ability			Overall		
	$Acc^+$	$Acc^-$	$\Delta_{low}$	$Acc^+$	$Acc^-$	$\Delta_{mid}$	$Acc^+$	$Acc^-$	$\Delta_{high}$	$Acc^+$	$Acc^-$	$\Delta$
<i>Weak LLMs</i>												
Qwen-3B	0.435	0.641	0.206	0.423	0.631	0.208	0.468	0.602	0.134	0.447	0.622	0.175
Qwen-7B	0.479	0.571	0.092	0.444	0.602	0.158	0.442	0.598	0.156	0.449	0.592	0.143
Qwen-14B	0.727	0.378	0.349	0.648	0.405	0.243	0.661	0.402	0.259	0.667	0.397	0.270
LLaMA-1B	0.131	0.894	0.763	0.130	0.896	0.766	0.138	0.892	0.754	0.134	0.894	0.760
LLaMA-3B	0.119	0.846	0.727	0.121	0.870	0.749	0.119	0.863	0.744	0.120	0.862	0.742
LLaMA-8B	0.446	0.610	0.164	0.434	0.627	0.193	0.384	0.664	0.280	0.412	0.638	0.226
<i>Strong LLMs</i>												
GPT-5-Mini	0.871	0.149	0.722	0.821	0.172	0.649	0.834	0.175	0.659	0.835	0.168	0.667
DeepSeek	0.822	0.267	0.555	0.752	0.311	0.441	0.736	0.336	0.400	0.755	0.311	0.444
<i>CD-based Method</i>												
Ours	0.567	0.525	<b>0.042</b>	0.517	0.516	<b>0.001</b>	0.476	0.551	<b>0.075</b>	0.505	0.532	<b>0.027</b>

and high proficiency levels, as well as the overall performance. We compare three types of simulation settings: six weak-capacity LLMs (Qwen-3B, Qwen-7B, Qwen-14B, LLaMA-1B, LLaMA-3B, and LLaMA-8B), two strong-capacity LLMs (DeepSeek and GPT-5-Mini)<sup>3</sup>, and our NeuralCD-guided adaptive LLM assignment method.

As illustrated in Table 1, our method achieves the superior *BAA* scores across all proficiency tiers as well as in the overall evaluation. This consistent performance gain demonstrates that adaptively matching students with cognitively aligned LLMs substantially enhances simulation fidelity, outperforming any single-model baseline across the entire learner spectrum.

We further observe a clear interaction between model capacity and student proficiency level. Most weak-capacity models perform relatively better when simulating low-proficiency students, but their *BAA* scores gradually decrease as the simulated student proficiency increases. This trend can be observed for models such as Qwen-7B and LLaMA-8B. In contrast, strong-capacity models demonstrate stronger performance when simulating high-proficiency students, with *BAA* scores increasing as student proficiency rises, as seen in DeepSeek and GPT-5-Mini.

Notably, **stronger model capacity does not necessarily translate into better simulation performance**. This observation is closely related to the simulation bias discussed in Section 5.1. As shown

<sup>3</sup>GPT-5-Mini (2025-08-07), DeepSeek-V3.2. The specific model types and configurations are detailed in Appendix D.

in Table 2, extremely strong models (e.g., GPT-5-Mini) tend to exhibit high  $Acc^+$  but low  $Acc^-$ , while very weak models (e.g., LLaMA-3B) show the opposite pattern. The large discrepancy between  $Acc^+$  and  $Acc^-$  for such models leads to suboptimal *BAA*, despite their strengths in one aspect of simulation.

These results underscore that no single LLM can serve as a universally optimal simulator across heterogeneous student populations. By leveraging cognitive diagnosis to adaptively assign LLMs based on student proficiency, our framework effectively mitigates the competence bias inherent in fixed-model paradigms, achieving superior simulation fidelity across the entire ability spectrum.

## 6 Conclusion

In this work, we identified the inherent biases of the single-LLM paradigm in student simulation and proposed an ability-aware framework grounded in NeuralCD. By dynamically matching human learners with appropriately capable LLM backbones, our method effectively mitigates competence-related biases and enhances simulation fidelity across diverse proficiency levels. This shift from performance-driven to alignment-oriented model selection establishes a principled foundation for more authentic educational simulations.

## 7 Limitations

Despite its effectiveness, our method has several limitations. First, due to cost constraints, the LLM

pool used in this study does not exhaustively cover all existing open-source and closed-source models.

Second, our approach relies on students' historical response records to estimate their cognitive states via cognitive diagnosis models. As a result, it is not directly applicable to cold-start scenarios where no prior student interaction data is available. Addressing this limitation may require incorporating additional sources of information, such as demographic features.

Third, our evaluation primarily focuses on the accuracy and bias of simulated student responses to problem-solving tasks. We do not explicitly assess the behavioral plausibility or rationality of the simulated actions, nor do we consider other important dimensions of student behavior, such as learning strategy development, metacognitive processes, or peer collaboration. Extending the evaluation beyond answer correctness to cover richer and more realistic learning behaviors is an important direction for future work.

Finally, because student simulation itself is not the primary focus of this work, we adopt a commonly used student agent framework from prior studies to isolate the impact of LLM selection. However, these existing simulation architectures are not guaranteed to be optimal. As shown in Table 2, our method does not achieve the best overall accuracy. Improvements in student agent design may further enhance the overall performance of our framework. We leave the exploration of more advanced student simulation mechanisms to future work.

Notably, model fine-tuning (Touvron et al., 2023) is not considered as a baseline in our experiments. Fine-tuning essentially trains a new model to fit student abilities. In contrast, our approach focuses on *model assignment* rather than *model retraining*, allowing us to explicitly analyze and mitigate the inherent capacity-dependent biases of existing LLMs while preserving their original capabilities. This design choice also improves reproducibility and reduces training costs.

## References

Zahra Abbasiantaeb, Yifei Yuan, Evangelos Kanoulas, and Mohammad Aliannejadi. 2024. Let the llms talk: Simulating human-to-human conversational qa via zero-shot llm-to-llm interactions. In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*, pages 8–17.

Ghodai Abdelrahman, Sherif Abdelfattah, Qing Wang, and Yu Lin. 2022. Dbe-kt22: A knowledge tracing dataset based on online student evaluation. *arXiv preprint arXiv:2208.12651*.

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

Jasper Meynard Arana, Kristine Ann M Carandang, Ethan Robert Casin, Christian Alis, Daniel Stanley Tan, Erika Fille Legara, and Christopher Monterola. 2025. Foundations of peers: Assessing llm role performance in educational simulations. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 4: Student Research Workshop)*, pages 908–918.

Shir Ashury Tahan, Ariel Gera, Benjamin Sznajder, Leshem Choshen, Liat Ein-Dor, and Eyal Shnarch. 2024. Label-efficient model selection for text generation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8384–8402, Bangkok, Thailand. Association for Computational Linguistics.

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, and 1 others. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.

Frank B Baker. 2001. *The basics of item response theory*. ERIC.

Matthew L Bernacki, Meghan J Greene, and Nikki G Lobczowski. 2021. A systematic review of research on personalized learning: Personalized by whom, to what, how, and for what purpose (s)? *Educational Psychology Review*, 33(4):1675–1715.

Saptarshi Bhowmik, Luke West, Alex Barrett, Nuodi Zhang, Chih-Pu Dai, Zlatko Sokolij, Sherry Southerland, Xin Yuan, and Fengfeng Ke. 2024. Evaluation of an llm-powered student agent for teacher training. In *European conference on technology enhanced learning*, pages 68–74. Springer.

Kay Henning Brodersen, Cheng Soon Ong, Klaas Enno Stephan, and Joachim M Buhmann. 2010. The balanced accuracy and its posterior distribution. In *2010 20th international conference on pattern recognition*, pages 3121–3124. IEEE.

R Philip Chalmers. 2012. mirt: A multidimensional item response theory package for the r environment. *Journal of statistical Software*, 48:1–29.

Zhendong Chu, Shen Wang, Jian Xie, Tinghui Zhu, Yibo Yan, Jinheng Ye, Aoxiao Zhong, Xuming Hu, Jing Liang, Philip S Yu, and 1 others. 2025. Llm agents for education: Advances and applications. *arXiv preprint arXiv:2503.11733*.



822	Wei Song, Zhenya Huang, Cheng Cheng, Weibo Gao,	Ern, Phey Ling Kit, Jenny Giam Xiuhui, John Pinto,	878
823	Bihan Xu, GuanHao Zhao, Fei Wang, and Runze Wu.	and Ee-Peng Lim. 2025. <a href="#">Consistent client simulation</a>	879
824	2025. Irt-router: Effective and interpretable multi-	<a href="#">for motivational interviewing-based counseling</a> . In	880
825	llm routing via item response theory. <i>arXiv preprint</i>	<i>Proceedings of the 63rd Annual Meeting of the As-</i>	881
826	<i>arXiv:2506.01048</i> .	<i>sociation for Computational Linguistics (Volume 1:</i>	882
		<i>Long Papers)</i> , pages 20959–20998, Vienna, Austria.	883
827	Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-	Association for Computational Linguistics.	884
828	Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan		
829	Schalkwyk, Andrew M Dai, Anja Hauth, Katie Mil-	Zheyuan Zhang, Daniel Zhang-Li, Jifan Yu, Linlu Gong,	885
830	llican, and 1 others. 2023. Gemini: a family of	Jinchang Zhou, Zhanxin Hao, Jianxiao Jiang, Jie	886
831	highly capable multimodal models. <i>arXiv preprint</i>	Cao, Huiqin Liu, Zhiyuan Liu, and 1 others. 2025.	887
832	<i>arXiv:2312.11805</i> .	Simulating classroom education with llm-empowered	888
		agents. In <i>Proceedings of the 2025 Conference of the</i>	889
833	Jonathan L Templin and Robert A Henson. 2006.	<i>Nations of the Americas Chapter of the Association</i>	890
834	Measurement of psychological disorders using cog-	<i>for Computational Linguistics: Human Language</i>	891
835	gnitive diagnosis models. <i>Psychological methods</i> ,	<i>Technologies (Volume 1: Long Papers)</i> , pages 10364–	892
836	11(3):287.	10379.	893
837	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	Zesen Zhao, Shuwei Jin, and Z Morley Mao. 2024.	894
838	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	Eagle: Efficient training-free router for multi-llm	895
839	Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti	inference. <i>arXiv preprint arXiv:2409.15518</i> .	896
840	Bhosale, and 1 others. 2023. Llama 2: Open founda-		
841	tion and fine-tuned chat models. <i>arXiv preprint</i>	Longwei Zheng, Fei Jiang, Xiaoqing Gu, Yuanyuan Li,	897
842	<i>arXiv:2307.09288</i> .	Gong Wang, and Haomin Zhang. 2025. Teaching via	898
		llm-enhanced simulations: Authenticity and barriers	899
843	Fei Wang, Qi Liu, Enhong Chen, Zhenya Huang,	to suspension of disbelief. <i>The Internet and Higher</i>	900
844	Yu Yin, Shijin Wang, and Yu Su. 2022. Neuralcd:	<i>Education</i> , 65:100990.	901
845	a general framework for cognitive diagnosis. <i>IEEE</i>		
846	<i>Transactions on Knowledge and Data Engineering</i> ,	<b>A Student Agent Framework</b>	902
847	35(8):8312–8327.		
		In this appendix, we present two representative	903
848	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	works that focus on the simulation of student be-	904
849	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	havior in prior studies. And then we provide a	905
850	and 1 others. 2022. Chain-of-thought prompting elic-	detailed description of the student simulation agent	906
851	its reasoning in large language models. <i>Advances</i>	adopted in our framework.	907
852	<i>in neural information processing systems</i> , 35:24824–	The first work is EduAgent (Xu et al., 2024).	908
853	24837.	EduAgent adopts a modular design with two core	909
		space. <b>Memory Space</b> : it hierarchically stores	910
854	Qingsong Wen, Jing Liang, Carles Sierra, Rose Luckin,	physiological data (gaze trajectories, mouse opera-	911
855	Richard Tong, Zitao Liu, Peng Cui, and Jiliang Tang.	tions), cognitive data (6 types of states like work-	912
856	2024. Ai for education (ai4edu): Advancing person-	load), and knowledge data (post-course test results),	913
857	alized education with llm and adaptive learning. In	integrating student personas and course-related in-	914
858	<i>Proceedings of the 30th ACM SIGKDD Conference</i>	formation. <b>Action Space</b> : Outputs gaze/motor be-	915
859	<i>on Knowledge Discovery and Data Mining</i> , pages	haviors mapped to Areas of Interest (AOIs), and	916
860	6743–6744.	personalized question-answering performance.	917
		The second work is Agent4Edu (Gao et al.,	918
861	Yu Xia, Fang Kong, Tong Yu, Liya Guo, Ryan A Rossi,	2025). The LLM-powered generative agent in	919
862	Sungchul Kim, and Shuai Li. 2024. Which llm to	Agent4Edu integrates three specialized modules	920
863	play? convergence-aware online model selection	for personalized learning simulation: <b>Learner</b>	921
864	with time-increasing bandits. In <i>Proceedings of the</i>	<b>Profile Module</b> , initialized with real-world re-	922
865	<i>ACM Web Conference 2024</i> , pages 4059–4070.	sponse data to capture explicit practice styles (e.g.,	923
		activity, diversity) and implicit cognitive factors	924
866	Songlin Xu, Hao-Ning Wen, Hongyi Pan, Dallas	(e.g., problem-solving ability); <b>Memory Mod-</b>	925
867	Dominguez, Dongyin Hu, and Xinyu Zhang. 2025.	<b>ule</b> , designed based on human learning mech-	926
868	Classroom simulacra: Building contextual student	anisms to include factual memory (reinforced	927
869	generative agents in online education for learning	response records), short-term memory (recent	928
870	behavioral simulation. In <i>Proceedings of the 2025</i>	practice details), and long-term memory (rein-	929
871	<i>CHI Conference on Human Factors in Computing</i>	forced facts, LLM-generated summaries, and a	930
872	<i>Systems</i> , pages 1–26.		
873	Songlin Xu, Xinyu Zhang, and Lianhui Qin. 2024. Edu-		
874	agent: Generative student agents in learning. <i>arXiv</i>		
875	<i>preprint arXiv:2404.07963</i> .		
876	Yizhe Yang, Palakorn Achananuparp, Heyan Huang,		
877	Jing Jiang, Nicholas Gabriel Lim, Cameron Tan Shi		

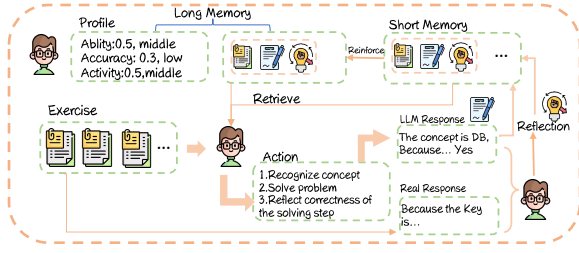


Figure 5: The Over View of Student Agent

forgetting curve) with retrieval, writing, and reflection capabilities; and **Action Module**, enabling human-like behaviors such as cognitive-driven exercise acceptance/rejection, exercise reading/understanding (with corrective reflection for mismatched knowledge concepts), and chain-of-thought analysis/solving (generating answers and correctness predictions, plus corrective reflection for inconsistencies).

Because our dataset only contains student-exercise interaction records and does not include fine-grained behavioral signals such as gaze trajectories or mouse operations, we adopt the second line of work as our student simulation framework.

As illustrated in Figure 5, the agent consists of three core modules: a memory module, a behavior module, and a reflection module. The agent simulates student learning by sequentially interacting with exercises in temporal order.

### A.1 Memory Module

The memory module maintains both short-term memory and long-term memory.

**Short-Term Memory.** Short-term memory stores the student agent’s most recent interaction records. Each memory entry is represented as  $r = (e, ans, reflection)$ , where  $e$  denotes the exercise,  $ans$  is the agent’s response, and  $reflection$  records the agent’s reflective analysis. Each exercise  $e$  includes both the exercise text and its associated knowledge concepts.

Short-term memory has a limited capacity. In our implementation, the maximum capacity is set to 5 recent interaction records. In addition, for each knowledge concept, we track its occurrence frequency within the short-term memory. If the occurrence count of a concept reaches a predefined threshold, the corresponding short-term memory entries are reinforced and transferred into long-term memory. We set the reinforcement threshold to 3 occurrences.

**Long-Term Memory.** Long-term memory stores more stable information about the student agent. It includes the student’s cognitive profile (i.e., the knowledge concept proficiency vector), historical response accuracy, and response activity level. Response activity is defined as the ratio between the number of exercises attempted by the student and the total number of exercises. These attributes allow the agent to condition its behavior on its overall ability and engagement level. Long-term memory also contains reinforced memories transferred from short-term memory.

### A.2 Action Module

The action module governs how the student agent answers an exercise. For a given exercise, the agent first identifies the knowledge concepts involved. Before answering, the agent retrieves all short-term memories as contextual input. It then searches long-term memory for records whose associated concepts overlap with those of the current exercise. If such records exist, they are collected as candidates, and the most recent three entries (in temporal order) are selected.

Based on the exercise content, extracted concepts, and retrieved memories, the agent generates a response to the exercise. After answering, the agent additionally judges whether it believes its own response is correct. The generated response and self-assessment together constitute the agent’s answer *ans*.

### A.3 Reflection Module

The reflection module is responsible for metacognitive correction. After the response is generated, the agent compares its self-assessment with the ground-truth outcome. If the agent correctly judges its response as correct, or correctly judges it as incorrect, no reflection is performed. Otherwise, the agent triggers a reflection process to analyze why the response or the self-assessment is incorrect.

The reflection result, together with the exercise and the generated response, is written into short-term memory as part of the interaction record. This mechanism allows the agent to adjust future behavior based on past inconsistencies.

### A.4 Overall Workflow

The student agent simulates learning by interacting with exercises sequentially in chronological order. For each exercise, the agent first retrieves relevant memories, then generates a response through the

behavior module, performs reflection when necessary, and finally stores the resulting record in short-term memory. Through the interaction of memory, behavior, and reflection modules, the agent produces coherent and temporally grounded learning behavior simulation.

## B Model Training Configuration for Cognitive Diagnosis

The NeuralCDM model was trained on an NVIDIA GeForce RTX 4090 GPU. We used a batch size of 32 and trained for 5 epochs with the Adam optimizer at a learning rate of 0.002. The prediction network consisted of two fully connected hidden layers with output dimensions of 512 and 256, respectively. Dropout regularization with a rate of 0.5 was applied to both fully connected layers.

## C Student IDs Used in Experiments

Due to computational and cost constraints, we do not conduct LLM-based simulation on the full student population. Instead, we select a representative subset of students for evaluation.

Specifically, based on the cognitive diagnosis results obtained from NeuralCD, we stratify students into three ability groups: *low-ability*, *medium-ability*, and *high-ability*. From each group, we randomly sample 100 students, resulting in a total of 300 students used in the experiments.

To facilitate reproducibility and fair comparison in future work, we explicitly list the student identifiers corresponding to each ability group in this appendix. These identifiers uniquely determine the subset of students used for evaluation and allow exact reconstruction of the experimental setting.

**Low-Ability Students.** The student IDs corresponding to the low-ability group are listed as follows:

- {2, 5, 10, 11, 17, 19, 24, 30, 31, 39, 41, 43, 45, 60, 63, 66, 71, 72, 73, 74, 81, 86, 87, 90, 95, 98, 102, 105, 107, 108, 120, 122, 124, 126, 127, 131, 133, 135, 152, 156, 157, 172, 176, 185, 188, 198, 201, 205, 212, 218, 222, 223, 232, 236, 237, 241, 251, 255, 260, 265, 272, 275, 278, 279, 280, 281, 283, 287, 291, 293, 297, 299, 301, 302, 304, 307, 308, 311, 312, 313, 314, 315, 317, 323, 327, 329, 333, 334, 335, 336, 346, 347, 348, 350, 352, 354, 360, 369, 370, 375}

**Medium-Ability Students.** The student IDs corresponding to the medium-ability group are listed as follows:

- {12, 16, 21, 23, 25, 26, 29, 33, 35, 38, 40, 46, 49, 53, 54, 55, 59, 65, 67, 84, 92, 94, 99, 101, 111, 123, 132, 138, 140, 141, 144, 147, 154, 159, 161, 162, 163, 171, 177, 179, 180, 181, 183, 186, 187, 202, 206, 209, 210, 214, 216, 220, 221, 226, 228, 242, 245, 252, 253, 261, 262, 266, 267, 268, 274, 282, 288, 294, 295, 296, 303, 305, 309, 310, 316, 319, 322, 328, 331, 341, 342, 343, 359, 361, 363, 365, 366, 373, 374, 379, 380, 381, 383, 384, 388, 392, 394, 397, 399, 401}

**High-Ability Students.** The student IDs corresponding to the high-ability group are listed as follows:

- {1, 15, 20, 22, 37, 42, 50, 51, 56, 69, 75, 77, 79, 80, 93, 97, 100, 104, 109, 110, 115, 118, 119, 142, 150, 155, 164, 165, 169, 170, 173, 178, 190, 191, 197, 203, 233, 243, 246, 256, 259, 263, 269, 270, 271, 273, 276, 286, 289, 290, 298, 318, 321, 325, 338, 340, 344, 345, 349, 355, 364, 368, 371, 377, 390, 391, 398, 408, 426, 427, 428, 429, 430, 431, 432, 434, 438, 443, 444, 445, 447, 448, 451, 455, 458, 466, 468, 472, 479, 480, 481, 483, 486, 498, 500, 505, 506, 507, 508, 509}

In total, after the above selection process, we obtain 7,697 student–exercise interaction records, which are used in all subsequent experiments.

## D LLMs Pool

We list all large language models included in the LLM pool used in our experiments, together with their corresponding access addresses. Each model is indexed for ease of reference and reproducibility.

1. **Qwen1.5-0.5B-Chat:** <https://huggingface.co/Qwen/Qwen1.5-0.5B-Chat>
2. **Qwen1.5-1.8B-Chat:** <https://huggingface.co/Qwen/Qwen1.5-1.8B-Chat>
3. **Qwen1.5-7B-Chat:** <https://huggingface.co/Qwen/Qwen1.5-7B-Chat>
4. **Qwen1.5-32B-Chat:** <https://huggingface.co/Qwen/Qwen1.5-32B-Chat>

1114	5. <b>Qwen1.5-110B-Chat:</b>	<a href="https://huggingface.co/Qwen/Qwen1.5-110B-Chat">https://huggingface.co/Qwen/Qwen1.5-110B-Chat</a>	21. <b>Qwen3-32B:</b>	<a href="https://huggingface.co/Qwen/Qwen3-32B">https://huggingface.co/Qwen/Qwen3-32B</a>	1157
1115					1158
1116					
1117	6. <b>Qwen2-0.5B-Instruct:</b>	<a href="https://huggingface.co/Qwen/Qwen2-0.5B-Instruct">https://huggingface.co/Qwen/Qwen2-0.5B-Instruct</a>	22. <b>Qwen-Turbo (2025-07-15):</b>	<a href="https://bailian.console.aliyun.com">https://bailian.console.aliyun.com</a>	1159
1118					1160
1119					
1120	7. <b>Qwen2-1.5B-Instruct:</b>	<a href="https://huggingface.co/Qwen/Qwen2-1.5B-Instruct">https://huggingface.co/Qwen/Qwen2-1.5B-Instruct</a>	23. <b>Qwen-Plus (2025-09-11):</b>	<a href="https://bailian.console.aliyun.com">https://bailian.console.aliyun.com</a>	1161
1121					1162
1122					
1123	8. <b>Qwen2-7B-Instruct:</b>	<a href="https://huggingface.co/Qwen/Qwen2-7B-Instruct">https://huggingface.co/Qwen/Qwen2-7B-Instruct</a>	24. <b>Qwen-Max (2025-01-25):</b>	<a href="https://bailian.console.aliyun.com">https://bailian.console.aliyun.com</a>	1163
1124					1164
1125					
1126	9. <b>Qwen2-57B-A14B-Instruct:</b>	<a href="https://huggingface.co/Qwen/Qwen2-57B-A14B-Instruct">https://huggingface.co/Qwen/Qwen2-57B-A14B-Instruct</a>	25. <b>LLaMA-3.2-1B-Instruct:</b>	<a href="https://huggingface.co/meta-llama/Llama-3.2-1B-Instruct">https://huggingface.co/meta-llama/Llama-3.2-1B-Instruct</a>	1165
1127					1166
1128					1167
1129	10. <b>Qwen2-72B-Instruct:</b>	<a href="https://huggingface.co/Qwen/Qwen2-72B-Instruct">https://huggingface.co/Qwen/Qwen2-72B-Instruct</a>	26. <b>LLaMA-3.2-3B-Instruct:</b>	<a href="https://huggingface.co/meta-llama/Llama-3.2-3B-Instruct">https://huggingface.co/meta-llama/Llama-3.2-3B-Instruct</a>	1168
1130					1169
1131					1170
1132	11. <b>Qwen2.5-3B-Instruct:</b>	<a href="https://huggingface.co/Qwen/Qwen2.5-3B-Instruct">https://huggingface.co/Qwen/Qwen2.5-3B-Instruct</a>	27. <b>LLaMA-3.1-8B-Instruct:</b>	<a href="https://huggingface.co/meta-llama/Llama-3.1-8B-Instruct">https://huggingface.co/meta-llama/Llama-3.1-8B-Instruct</a>	1171
1133					1172
1134					1173
1135	12. <b>Qwen2.5-7B-Instruct:</b>	<a href="https://huggingface.co/Qwen/Qwen2.5-7B-Instruct">https://huggingface.co/Qwen/Qwen2.5-7B-Instruct</a>	28. <b>LLaMA-3.1-70B-Instruct:</b>	<a href="https://huggingface.co/meta-llama/Llama-3.1-70B-Instruct">https://huggingface.co/meta-llama/Llama-3.1-70B-Instruct</a>	1174
1136					1175
1137					1176
1138	13. <b>Qwen2.5-14B-Instruct:</b>	<a href="https://huggingface.co/Qwen/Qwen2.5-14B-Instruct">https://huggingface.co/Qwen/Qwen2.5-14B-Instruct</a>	29. <b>GLM-4-9B-Chat:</b>	<a href="https://huggingface.co/zai-org/glm-4-9b-chat">https://huggingface.co/zai-org/glm-4-9b-chat</a>	1177
1139					1178
1140					
1141	14. <b>Qwen2.5-32B-Instruct:</b>	<a href="https://huggingface.co/Qwen/Qwen2.5-32B-Instruct">https://huggingface.co/Qwen/Qwen2.5-32B-Instruct</a>	30. <b>GLM-4-32B-0414:</b>	<a href="https://huggingface.co/zai-org/GLM-4-32B-0414">https://huggingface.co/zai-org/GLM-4-32B-0414</a>	1179
1142					1180
1143					
1144	15. <b>Qwen2.5-72B-Instruct:</b>	<a href="https://huggingface.co/Qwen/Qwen2.5-72B-Instruct">https://huggingface.co/Qwen/Qwen2.5-72B-Instruct</a>	31. <b>GLM-4.5-Air:</b>	<a href="https://huggingface.co/zai-org/GLM-4.5-Air">https://huggingface.co/zai-org/GLM-4.5-Air</a>	1181
1145					1182
1146					
1147	16. <b>Qwen3-0.6B:</b>	<a href="https://huggingface.co/Qwen/Qwen3-0.6B">https://huggingface.co/Qwen/Qwen3-0.6B</a>	32. <b>GPT-3.5-Turbo (0125):</b>	<a href="https://platform.openai.com/docs/models/gpt-3.5-turbo">https://platform.openai.com/docs/models/gpt-3.5-turbo</a>	1183
1148					1184
1149	17. <b>Qwen3-1.7B:</b>	<a href="https://huggingface.co/Qwen/Qwen3-1.7B">https://huggingface.co/Qwen/Qwen3-1.7B</a>	33. <b>GPT-4.1-Mini (2025-04-14):</b>	<a href="https://platform.openai.com/docs/models/gpt-4.1-mini">https://platform.openai.com/docs/models/gpt-4.1-mini</a>	1185
1150					1186
1151	18. <b>Qwen3-4B:</b>	<a href="https://huggingface.co/Qwen/Qwen3-4B">https://huggingface.co/Qwen/Qwen3-4B</a>	34. <b>GPT-5-Mini (2025-08-07):</b>	<a href="https://platform.openai.com/docs/models/gpt-5-mini">https://platform.openai.com/docs/models/gpt-5-mini</a>	1187
1152					1188
1153	19. <b>Qwen3-8B:</b>	<a href="https://huggingface.co/Qwen/Qwen3-8B">https://huggingface.co/Qwen/Qwen3-8B</a>			1189
1154					1190
1155	20. <b>Qwen3-14B:</b>	<a href="https://huggingface.co/Qwen/Qwen3-14B">https://huggingface.co/Qwen/Qwen3-14B</a>	35. <b>DeepSeek-V3.2:</b>	<a href="https://platform.deepseek.com">https://platform.deepseek.com</a>	1191
1156					1192

1194	<b>E Case Study</b>	
1195	<b>E.1 Case Study of Negative Bias Caused by Weak LLMs</b>	
1196		
1197	We present a concrete example illustrating the negative bias introduced by a low-capacity LLM when simulating a student who answered correctly. The student id is 190, the exercise id is 181.	1241
1198		1242
1199		1243
1200		1244
1201	<b>Exercise</b> <i>If we wish to grant a privilege to a user and also allow the user to pass the privilege on to other users, which of the following should we append to the appropriate grant command?</i>	1245
1202		1246
1203	<b>Options:</b> (A) With grant privilege (B) Allow grant privilege (C) Allow grant option (D) With grant option <b>Correct Answer:</b> D	
1204	<b>Knowledge Concepts:</b> Discretionary access control (DAC), GRANT	
1205		
1206		
1207		
1208		
1209		
1210	<b>Student Response</b> The student answered correctly by selecting <b>(D) With grant option.</b>	
1211		
1212	<b>Qwen2.5-7B Response (Weak Baseline)</b> GRANT. The task requires understanding the SQL GRANT command and its limitations. The GRANT command cannot be used to grant privileges on a specified tuple directly. Therefore, the correct answer is (B) A specified tuple. Yes	
1213		
1214		
1215		
1216		
1217		
1218	<b>Our Method (Ability-Aware Assignment) Response</b> Discretionary access control (DAC), GRANT. Since the question involves granting privileges and passing those privileges on to other users, the correct answer should allow for the propagation of granted privileges. The most suitable option is (D) With grant option, as it enables the user to grant the same privilege to other users. Yes.	
1219		
1220		
1221		
1222		
1223		
1224		
1225		
1226	<b>Analysis</b> Although the student demonstrated correct understanding of the GRANT command with the <i>with grant option</i> , Qwen2.5-7B produced an incorrect response. This indicates that the intrinsic capacity of Qwen2.5-7B is insufficient to faithfully simulate students who have already mastered this concept.	
1227		
1228		
1229		
1230		
1231		
1232		
1233	In contrast, our ability-aware framework assigns an LLM whose diagnosed capability better aligns with the student’s cognitive state. As a result, the simulated response correctly reflects the student’s mastery of discretionary access control and privilege delegation.	
1234		
1235		
1236		
1237		
1238		
	<b>E.2 Case Study of Positive Bias Caused by Strong LLMs</b>	1239
		1240
	We present a representative example illustrating positive bias introduced by a high-capacity LLM, where the model overestimates a student’s cognitive ability and produces a correct response even when the student answered incorrectly. The student id is 190, the exercise id is 177.	1241
		1242
		1243
		1244
		1245
		1246
	<b>Exercise</b> <i>Consider the following database schema:</i>	1247
		1248
	Movie(title, production year, country, run time, major genre)	1249
	primary key: {title, production year}	1250
	Person(id, first name, last name, year born)	1251
	primary key: {id}	1252
	Director(id, title, production year)	1253
	foreign keys: [title, production year] $\subseteq$ Movie[title, production year], [id] $\subseteq$ Person[id]	1254
	Writer(id, title, production year, credits)	1255
	foreign keys: [title, production year] $\subseteq$ Movie[title, production year], [id] $\subseteq$ Person[id]	1256
		1257
		1258
		1259
		1260
		1261
		1262
	<i>Which of the following relational algebra expressions output the titles and production years of movies written by writers born after 1980?</i>	1263
		1264
	(A): $\pi_{title, production\ year}(Movie \bowtie (Writer \bowtie \sigma_{year\ born > 1980}(Person)))$	1265
	(B) $\pi_{title, production\ year}(Writer \bowtie \sigma_{year\ born > 1980}(Person))$	1266
		1267
	<b>Options:</b> (A) only A (B) only B (C) both A and B (D) neither A nor B	1268
		1269
	<b>Correct Answer:</b> C	1270
	<b>Knowledge Concept:</b> Relational Algebra	1271
		1272
		1273
	<b>Student Response</b> The student answered <b>incorrectly.</b>	1274
		1275
	<b>GPT-5-mini Response (Strong Baseline)</b> Relational Algebra. Join Writer with Person filtered by year born > 1980, then project title and production year. The Movie join in A is redundant because Writer has a foreign key to Movie, so both expressions yield the same result. Final answer: (C) both A and B. Yes.	1276
		1277
		1278
		1279
		1280
		1281
		1282
	<b>Our Method (Ability-Aware Assignment) Response</b> Relational Algebra. The solution involves	1283
		1284

1285 understanding how to join tables in relational al-  
1286 gebra and apply selection conditions. Both alter-  
1287 natives involve joining the Writer table with the  
1288 Person table to filter out writers born after 1980.  
1289 However, alternative A also joins the result with  
1290 the Movie table and projects the title and produc-  
1291 tion year, which is necessary to get the required  
1292 information. Alternative B only joins Writer with  
1293 Person but does not project the movie details, mak-  
1294 ing it incomplete for the required output. Therefore,  
1295 the correct answer is A. Yes.

1296 **Analysis** In this example, the student failed to  
1297 correctly reason about equivalence in relational al-  
1298 gebra and answered the question incorrectly. How-  
1299 ever, the high-capacity model GPT-5-mini pro-  
1300 duced the correct answer by leveraging its strong  
1301 internal reasoning ability.

1302 In contrast, our ability-aware framework assigns  
1303 an LLM whose diagnosed capability better matches  
1304 the student’s cognitive state. As a result, the sim-  
1305 ulated response reproduces the student’s incorrect  
1306 reasoning pattern, providing a more faithful repre-  
1307 sentation of the learner’s actual mastery level.