

# GraphMASAL: A Graph-based Multi-Agent System for Adaptive Learning

Biqing Zeng\*  
South China Normal University  
Foshan, China  
zengbiqing137@163.com

Mengquan Liu  
South China Normal University  
Foshan, China  
2860251365@qq.com

Zongwei Zhen  
South China Normal University  
Foshan, China  
1169619484@qq.com

## ABSTRACT

The advent of Intelligent Tutoring Systems (ITSs) has marked a paradigm shift in education, enabling highly personalized learning pathways. However, true personalization requires adapting to learners' complex knowledge states (multi-source) and diverse goals (multi-sink); existing ITSs often lack the necessary structural-reasoning capability and knowledge dynamism to generate genuinely effective learning paths, and they lack scientifically rigorous validation paradigms. In this paper we propose GraphMASAL (A Graph-based Multi-Agent System for Adaptive Learning), which integrates (i) a dynamic knowledge graph for persistent, stateful learner modeling; (ii) a LangGraph-orchestrated trio of agents (Diagnoser, Planner, Tutor); (iii) a knowledge-graph-grounded two-stage neural IR component (dual-encoder dense retrieval with cross-encoder listwise re-ranking and calibrated score fusion); and (iv) a multi-source multi-sink (MSMS) planning engine with a cognitively grounded cost and an approximation guarantee via greedy set cover. Under blinded automated evaluations with matched inputs and inference settings across diverse student profiles, GraphMASAL consistently outperforms LLM prompting and structured ablations in planning—achieving stronger structural/sequence alignment of learning paths, higher coverage of weak concepts, and lower learning cost—while also surpassing prompt-based baselines in cognitive diagnosis. Agreement with expert/LLM-proxy ratings further supports the validity of our evaluation protocol. These findings indicate that grounding LLM agents in a dynamic knowledge graph, coupled with optimization under educational constraints, yields reliable, interpretable, and pedagogically plausible learning plans, advancing personalized and goal-oriented education.

## KEYWORDS

Intelligent Tutoring Systems, Large Language Model, Multi-agent

### ACM Reference Format:

Biqing Zeng, Mengquan Liu, and Zongwei Zhen. 2026. GraphMASAL: A Graph-based Multi-Agent System for Adaptive Learning. In *Proc. of the 25th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2026)*, Paphos, Cyprus, May 25 – 29, 2026, IFAAMAS, 9 pages.

## 1 INTRODUCTION

The advent of Intelligent Tutoring Systems (ITSs) has marked a paradigm shift in education, moving us closer to the long-held vision of

a truly personalized learning experience for every individual [2, 15]. The ultimate goal is an intelligent system that not only possesses a deep understanding of a dynamic knowledge domain but can also perform complex reasoning to execute multi-step, adaptive instructional tasks. Such a system would act as a personal mentor, guiding a learner from their current state of understanding towards their unique learning objectives.

However, the path to achieving this vision is fraught with systemic challenges that have prevented existing systems from reaching their full potential. We identify four interconnected bottlenecks that must be addressed in a holistic manner:

First is the Knowledge Layer challenge. In an era of rapid information growth, knowledge is not static. Most ITSs are built upon fixed curricula or static knowledge bases, which quickly become outdated and fail to capture the intricate, evolving relationships between concepts. A robust foundation for adaptive learning requires a dynamic knowledge graph that can continuously evolve, reflecting the current state of the domain.

Second, we face the Execution Layer challenge. The pedagogical process of "diagnose, plan, and tutor" is not a monolithic task but a complex workflow requiring distinct expertise. Entrusting this entire chain to a single, monolithic model or a simple series of large language model (LLM) calls leads to brittleness and a lack of specialized reasoning. True personalization requires a team of collaborative agents, each an expert in its sub-task, working in concert to create a seamless and effective learning experience.

Third, at the heart of this workflow lies the Decision Core challenge. Even with a perfect knowledge base and a team of expert agents, the fundamental task of charting an optimal learning path is a profound computational problem. Mapping a learner's initial knowledge state (a multi-source problem) to their diverse goals (a multi-sink problem) is an NP-hard optimization task. Existing agents lack the sophisticated "cognitive engine" needed to navigate this complexity and generate paths that are not just correct, but pedagogically efficient and sound.

Finally, towering over all these is the Scientific Validation challenge. If we build such a complex system, how do we scientifically measure the quality of its output? Evaluating a generated learning path is non-trivial. Simple accuracy metrics are insufficient; we need a new evaluation paradigm capable of assessing a path's logical coherence, pedagogical soundness, and structural integrity compared to an expert-defined ideal. Without this, the field cannot reliably measure progress.

To overcome these systemic barriers, we propose GraphMASAL, an integrated, graph-based multi-agent system for adaptive learning. Our framework is designed as a complete, end-to-end solution

\*Corresponding author

that systematically addresses each of the four challenges. The contributions of this paper are, therefore, four-fold:

We establish a dynamic knowledge graph as the system’s foundation, which not only ensures the domain knowledge is current and comprehensive but also allows for real-time updates to reflect a student’s evolving knowledge state.

We design a three-agent collaborative architecture (Diagnoser, Planner, Tutor) as the workflow engine, enabling modular, automated, and expert-driven execution of the tutoring process.

We introduce a novel multi-source, multi-sink path optimization algorithm as the core decision-making engine, empowering the Planner agent with the intelligence to generate highly effective learning routes.

We propose a new scientific evaluation paradigm, centered on our PathSim metric, to provide a robust and reliable methodology for measuring the quality of the system’s output. Beyond these, we introduce a KG-enhanced semantic retrieval pipeline with principled ranking fusion, and provide a formal treatment of the multi-source multi-sink (MSMS) planning with approximation guarantees under a cognitively grounded cost model.

This paper details the architecture of GraphMASAL, elaborates on its core algorithmic and methodological innovations, and presents extensive experimental validation demonstrating its superiority over existing approaches.

## 2 RELATED WORKS

### 2.1 The Evolution of Personalized Learning and Its Systemic Challenges

Personalized learning in ITS seeks to scale one-to-one tutoring by constructing precise learner and content models [1, 12, 15], evolving from early rule-based systems to modern data-driven platforms that leverage machine learning [2]. Despite progress, deep personalization remains constrained by coupled challenges in knowledge dynamism, execution complexity, algorithmic optimization, and scientific validation. GraphMASAL is designed to address these challenges holistically.

### 2.2 Knowledge Representation for Adaptive Learning

A system’s ability to reason and personalize is fundamentally constrained by its method of knowledge representation. The field has evolved from static, schema-oriented ontologies—which, despite their rich semantics, are ill-suited for rapidly changing domains—to more flexible, instance-oriented Knowledge Graphs (KGs) [16]. Unlike concept or mind maps designed for human cognition, KGs are machine-readable structures that provide a robust foundation for advanced AI tasks.

### 2.3 Modern Agentic Workflows with LLMs and LangGraph

The emergence of Large Language Models (LLMs) has revolutionized the creation of intelligent agents [17]. However, orchestrating multiple LLM-based agents to perform complex tasks introduces significant challenges in state management, reliable communication, and workflow control [7]. To address this, a new generation of

agentic frameworks, such as LangGraph, has emerged. LangGraph enables the modeling of complex, stateful, multi-agent applications as graphs, where nodes represent agents or tools and edges define the control flow [4]. This provides a modular, transparent, and persistent structure that resolves the rigidity and state-management issues of earlier multi-agent systems.

GraphMASAL’s three-agent architecture (Diagnoser, Planner, Tutor) is an advanced application of this modern paradigm. By implementing this classic ITS role division within LangGraph’s robust orchestration framework, our system ensures a seamless and logically sound flow of information. This approach mirrors the evolution in software engineering from monolithic architectures to orchestrated microservices [8, 10], bringing a new level of scalability, maintainability, and extensibility to the design of Intelligent Tutoring Systems.

Beyond orchestration, our retrieval stack follows the dense-retrieval + cross-encoder reranking paradigm popular in neural IR and learning-to-rank [5]. Different from generic RAG, we ground dense retrieval on KG nodes and prerequisite subgraphs, enabling structure-aware reranking and personalization.

## 3 THE PROPOSED FRAMEWORK

To achieve truly personalized, dynamic, and effective adaptive learning, we designed and implemented GraphMASAL, a graph-based multi-agent system. The core principle of this framework is to unify the modeling of a student’s cognitive state, the domain knowledge structure, and pedagogical strategies within a dynamically evolving environment, driven by the collaborative reasoning of intelligent agents. This section elaborates on the framework’s four core system components—the semantic retrieval engine for knowledge grounding, the dynamic knowledge graph as the knowledge foundation, the multi-agent collaborative workflow, and the path optimization algorithm as the decision engine—as well as the PathSim methodology for scientific evaluation.

### 3.1 Semantic Retrieval Engine and Personalization

We operationalize a KG-enhanced retrieval pipeline with four stages: (1) **Embedding** texts (queries, concept names/descriptions, item stems) using `doubao-embedding-text-240715`; (2) **Vector Similarity** search over concept/item vectors (cosine TopK) backed by Neo4j indices; (3) **Cross-Encoder Reranking** to re-score candidates with context-aware matching; and (4) **Personalization** that annotates candidates with student-specific mastery and misconceptions. Given a query  $q$  and candidate  $d$ , we fuse bi-encoder similarity and cross-encoder scores:

$$\text{Score}(d \mid q) = \alpha \cos(\mathbf{e}(q), \mathbf{e}(d)) + (1 - \alpha) \sigma(f_{ce}(q, d)), \quad (1)$$

where  $\mathbf{e}(\cdot)$  is the embedding,  $f_{ce}$  is the cross-encoder,  $\sigma$  is a min-max normalization over the candidate set, and  $\alpha \in [0, 1]$  (default  $\alpha=0.5$ ). We return Top-K (default  $K=6$ ) to the Tutor agent or tools (e.g., QueryTool).

The system architecture of GraphMASAL, depicted in Figure 1, is a hierarchical and collaborative structure built upon a dynamic knowledge graph foundation that supports three specialized agents

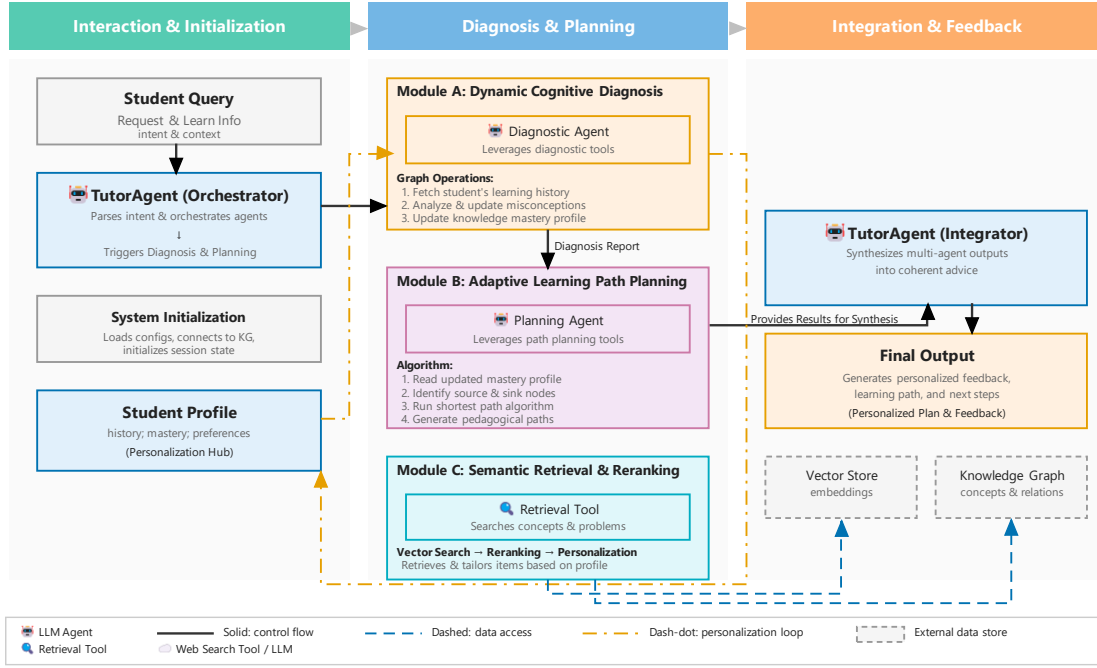


Figure 1: System architecture of GraphMASAL showing three phases with modular agent components and data flows.

(Diagnoser, Planner, Tutor) orchestrated via LangGraph, with planning driven by the MSMS algorithm and validated through the PathSim evaluation paradigm.

### 3.2 The Knowledge Foundation: Dynamic Knowledge Graph

The fundamental prerequisite for adaptive learning is the precise and real-time capture of a learner’s dynamic cognitive state. Static knowledge bases cannot meet this requirement. Therefore, we implement a dynamic knowledge graph mechanism that supports real-time updates and structural evolution. Its "dynamic" nature is manifested on two levels:

**Structural Dynamics:** The knowledge graph’s schema can be flexibly extended with the introduction of new courses and concepts.

**State Dynamics:** More critically, the attributes and relationships representing student states are updated in real-time with every student interaction (e.g., correctly or incorrectly answering a question).

#### Algorithm 1 KG-Enhanced Semantic Retrieval with Reranking

- 1: Input: query  $q$ , candidate set  $\mathcal{D}$ , weights  $\alpha$ , cutoffs  $K$
- 2: Embed  $q \rightarrow \mathbf{e}(q)$ ; for each  $d \in \mathcal{D}$ , get  $\mathbf{e}(d)$
- 3: Compute  $\text{sim}(q, d) = \cos(\mathbf{e}(q), \mathbf{e}(d))$ ; keep Top- $K_v$
- 4: Re-rank  $\mathcal{D}_{K_v}$  using cross-encoder  $f_{ce}(q, d)$ ; min-max normalize
- 5: Fuse scores by Eq.(1); return Top- $K$  with personalization annotations

This design transforms the knowledge graph from a static "knowledge map" into a dynamic "cognitive model," providing a solid and instantaneous data foundation for the decision-making of the higher-level agents, which is key to achieving truly personalized education.

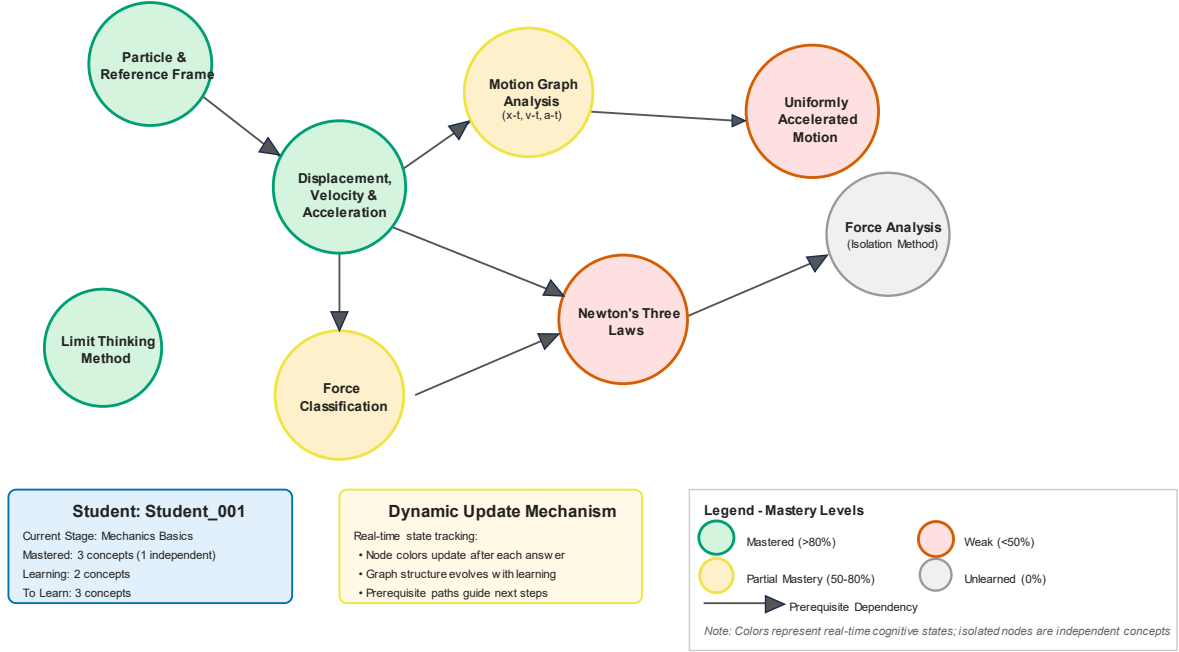
### 3.3 The Agentic Workflow: Multi-Agent Collaboration

We decompose the complex pedagogical task into three collaborative agents and orchestrate their interactions using the LangGraph framework to ensure the stability, orderliness, and state consistency of the information flow.

**Diagnostic Agent:** Its core responsibility is "deep attribution" [9]. Upon receiving a student’s attempt record, it not only judges correctness but, more importantly, queries the knowledge graph to associate the student’s incorrect options with specific Misconception nodes and updates the mastery level of related Concept nodes accordingly. It answers the question, "Why did the student make this mistake?"

**Planning Agent:** Its core responsibility is "path optimization" [6, 19]. It receives the student’s current cognitive state (i.e., a set of weak concepts) from the Diagnostic Agent and uses this as input to invoke its core planning engine (see Section 2.4) to generate one or more optimal learning paths. It answers the question, "What should the student learn next, and how?"

**Tutor Agent:** Its core responsibility is "interaction and coordination." It acts as the master controller, managing multi-turn dialogues with the user, dispatching requests to the Diagnostic or Planning



**Figure 2: Dynamic Knowledge Graph Structure Physics Concept Network with Prerequisite Dependencies**

Agents, and integrating their structured outputs to generate insightful and easy-to-understand natural language feedback.

By constructing a StateGraph with LangGraph, we can precisely define the trigger conditions and execution logic for each agent, making the entire collaborative process resemble a sophisticated "cognitive processing pipeline" that guarantees the coherence and logicity of the tutoring service.

### 3.4 The Core Planning Engine: Multi-Source Multi-Sink Path Optimization Algorithm

The technical core of our system is the Multi-Source Multi-Sink (MSMS) Path Optimization Algorithm embedded within the Planning Agent. Traditional pathfinding algorithms (single-source, single-sink) are ill-equipped to handle the complex, realistic scenario where a student has multiple knowledge gaps (multi-sink) and a heterogeneous knowledge base (multi-source). The MSMS algorithm is specifically designed to address this challenge.

The algorithm models the learning path generation task as an optimization problem on a graph. The set of "sources" comprises the concepts the student has already mastered, while the set of "sinks" represents the weak concepts to be learned. The objective is to find a set of paths originating from the sources that covers all sinks while minimizing the total learning cost.

Crucially, the design preference of this algorithm is deeply inspired by theories from educational psychology and cognitive science. Specifically, its optimization objective—minimizing the total number of new concepts in the resulting paths—is intended to

proactively manage the learner's Cognitive Load [13]. By avoiding the introduction of excessive irrelevant or redundant concepts at once, we can reduce extraneous cognitive load, allowing the learner to allocate more cognitive resources to the internalization and construction of knowledge (germane load). Furthermore, the algorithm's "multi-source" nature, starting from what the student already knows, aligns with Constructivist Learning Theory [14]. It scaffolds new learning by activating and connecting to the student's existing cognitive schemata, which is posited to enhance learning motivation and self-efficacy.

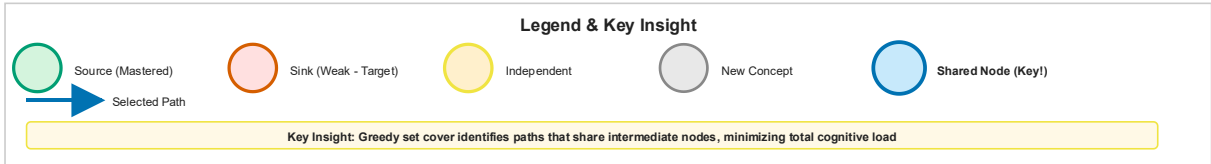
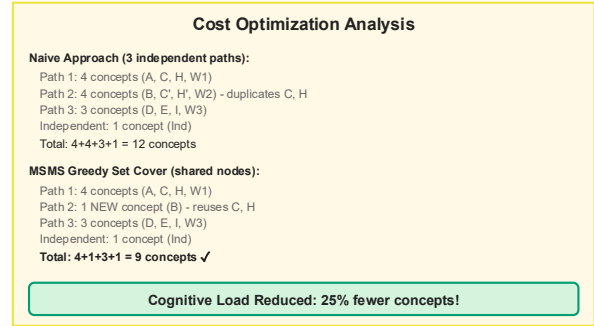
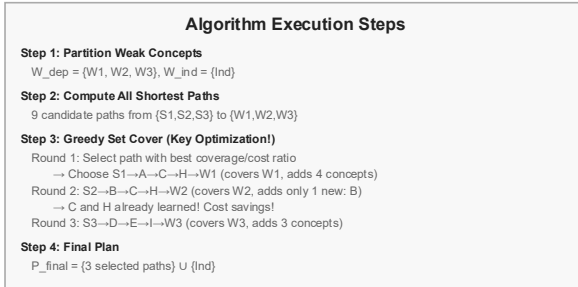
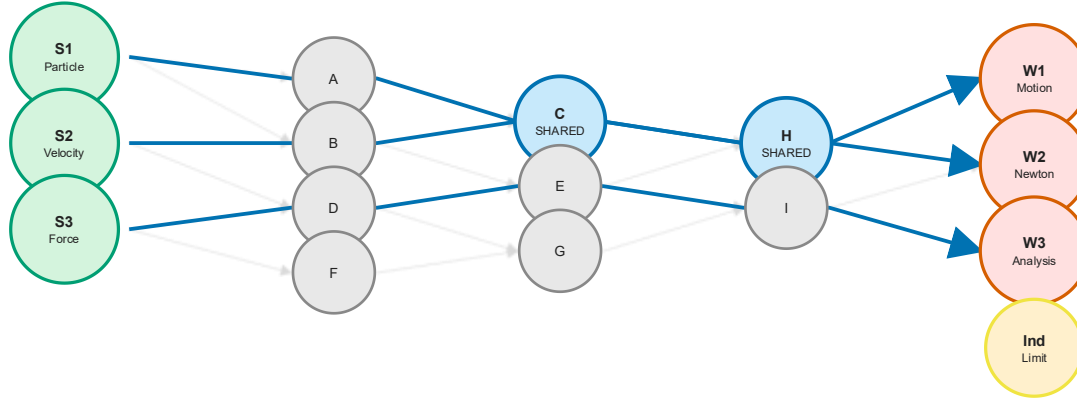
The detailed procedure is formalized in Algorithm 2.

**3.4.1 Problem Formulation and Cognitive Cost Model.** Let  $G = (V, E)$  be the prerequisite graph over concepts. For a given student,  $S \subset V$  denotes mastered sources and  $W \subset V$  weak sinks. A feasible plan selects paths  $\mathcal{P}$  such that each  $w \in W$  is covered by at least one path  $P : s \rightarrow w$  with  $s \in S$ . We minimize the total *learning cost*

$$\min_{\mathcal{P}} \sum_{v \in \cup_{P \in \mathcal{P}} P} c(v) \quad \text{s.t. } \forall w \in W, \exists P \in \mathcal{P} : P : s \rightarrow w, \quad (2)$$

where  $c(v) = \lambda_1(1 - \text{mastery}(v)) + \lambda_2 \text{difficulty}(v) + \lambda_3 \text{fanout}(v)$ . The cost penalizes novelty and split-attention while favoring scaffolding from mastered concepts.

**Operational definitions and normalization.** We instantiate the cost terms on the prerequisite graph as follows, with all components min–max normalized to  $[0, 1]$  for scale comparability.



**Figure 3: MSMS Algorithm for Multi-Source Multi-Sink Path Optimization via Greedy Set-Cover: Minimizing the Total Number of Emerging Concepts**

$mastery(v)$  is the student-specific proficiency stored in edge attribute  $proficiency \in [0, 1]$  of relation HAS\_MASTERY\_OF.  $difficulty(v)$  is a concept-level proxy derived as the normalized mean difficulty of problems linked to  $v$  via TESTS\_CONCEPT; when unavailable, we use normalized Concept.level.  $fanout(v)$  quantifies local branching computed as normalized out-degree  $deg^{out}(v)$  under IS\_PREREQUISITE\_FOR. This design penalizes detours through highly branching regions while encouraging coherent, scaffolded progress.

**3.4.2 Complexity and Approximation Guarantee.** MSMS planning generalizes Set Cover: each candidate path covers a subset of sinks. Minimizing total node cost to cover all sinks is NP-hard by reduction. Our planner enumerates shortest source  $\rightarrow$  sink paths (depth-limited) and applies greedy set-cover selection over path-sets. Let  $|W|$  be the number of sinks; the greedy procedure achieves a  $(1 + \ln |W|)$  approximation to the optimal cover cost [3, 11]. Bounding

path length and pruning by prerequisite depth control enumeration cost while preserving high-quality candidates.

**3.4.3 Pedagogical Validity.** The cost  $c(\cdot)$  operationalizes cognitive load theory: novelty and difficulty encode intrinsic/extraneous load; fan-out mitigates split attention. Multi-source anchoring aligns with constructivism by tying new learning to mastered schemata, yielding plans that are both computationally principled and pedagogically meaningful.

### 3.5 The Scientific Evaluation Paradigm: PathSim Methodology

To scientifically and objectively evaluate the performance of our planning algorithm, we propose a novel path similarity evaluation methodology, PathSim. Traditional graph or set similarity metrics

---

**Algorithm 2** MSMS Path Optimization

---

```
1: Input:
2:  $G = (V, E)$   $\triangleright$  Knowledge graph where  $V$  are concepts and  $E$  are prerequisite edges
3:  $S \subset V$   $\triangleright$  Set of mastered concepts (sources)
4:  $W \subset V$   $\triangleright$  Set of weak concepts (sinks)
5:
6: Output:
7:  $P_{\text{final}}$   $\triangleright$  A final learning plan, composed of paths and independent concepts
8:
9: function MSMSPLANNER( $G, S, W$ )
10:    $\triangleright$  Step 1: Partition weak concepts
11:    $W_{\text{ind}} \leftarrow \{w \in W \mid \neg \exists w' \in W, (w', w) \in E\}$   $\triangleright$  Independent concepts
12:    $W_{\text{dep}} \leftarrow W \setminus W_{\text{ind}}$   $\triangleright$  Dependent concepts
13:
14:    $\triangleright$  Step 2: Compute all-pairs shortest paths from sources to dependent sinks
15:    $AllPaths \leftarrow \emptyset$ 
16:   for all  $s \in S$  do
17:     for all  $w \in W_{\text{dep}}$  do
18:        $path \leftarrow \text{DIJKSTRA}(G, s, w)$ 
19:       if  $path$  exists then
20:         Add  $path$  to  $AllPaths$ 
21:       end if
22:     end for
23:   end for
24:
25:    $\triangleright$  Step 3: Greedily select paths to cover all dependent sinks
26:    $P_{\text{optimal}} \leftarrow \text{GREEDY\_SET\_COVER}(AllPaths, W_{\text{dep}})$ 
27:
28:    $\triangleright$  Step 4: Combine results into the final plan
29:    $P_{\text{final}} \leftarrow (P_{\text{optimal}}, W_{\text{ind}})$ 
30:   return  $P_{\text{final}}$ 
31: end function
```

---

(e.g., Jaccard similarity) are severely limited when assessing learning paths because they ignore two crucial pieces of information:

The Topological Structure of a Path: A learning path is more than a set of concepts; its internal node sequence and connectivity (i.e., dependency relations) are paramount.

The Composite Structure of a Plan: A complete learning plan is often a hybrid of multiple paths and several standalone concepts.

The PathSim method, a Topology-Aware Hybrid Path Matching Algorithm, is designed to address this gap. It employs a multi-level comparison framework that separately computes the similarity of the independent concept sets (using Jaccard similarity) and the path sets. When calculating path similarity, it holistically considers the overlap of nodes and edges, as well as the sequence similarity (calculated via normalized Levenshtein distance), thus enabling a more accurate measurement of the structural and sequential alignment between two paths. The PathSim formalized in Algorithm 3.

The introduction of PathSim provides an effective and interpretable scientific paradigm for the evaluation and comparison of

---

**Algorithm 3** PathSim Similarity Calculation

---

```
1: Input:
2:  $Plan_A = (P_A, I_A)$   $\triangleright$  Learning plan A
3:  $Plan_B = (P_B, I_B)$   $\triangleright$  Learning plan B
4:  $w_p, w_i$   $\triangleright$  Weights for path and independent concept similarity
5:
6: Output:
7:  $TotalSim$   $\triangleright$  The total similarity score between  $Plan\_A$  and  $Plan\_B$ 
8:
9: function PATHSIM( $Plan_A, Plan_B$ )
10:    $\triangleright$  Step 1: Calculate similarity for independent concepts
11:    $Sim_i \leftarrow \text{JACCARD}(I_A, I_B)$ 
12:
13:    $\triangleright$  Step 2: Calculate similarity for path sets
14:   if  $|P_A| = 0$  and  $|P_B| = 0$  then
15:      $Sim_p \leftarrow 1$ 
16:   else if  $|P_A| = 0$  or  $|P_B| = 0$  then
17:      $Sim_p \leftarrow 0$ 
18:   else
19:      $\triangleright$  2a: Compute cross-similarity matrix  $M$ 
20:     for  $i \leftarrow 1$  to  $|P_A|$  do
21:       for  $j \leftarrow 1$  to  $|P_B|$  do
22:          $M[i, j] \leftarrow \text{COMPUTE\_PATH\_SIM}(P_A[i], P_B[j])$ 
23:        $\triangleright$  Considers nodes, edges, sequence
24:     end for
25:
26:      $\triangleright$  2b: Calculate symmetric average best match
27:      $Sim_{A \rightarrow B} \leftarrow \frac{1}{|P_A|} \sum_{i=1}^{|P_A|} \max_j (M[i, j])$ 
28:      $Sim_{B \rightarrow A} \leftarrow \frac{1}{|P_B|} \sum_{j=1}^{|P_B|} \max_i (M[i, j])$ 
29:      $Sim_p \leftarrow (Sim_{A \rightarrow B} + Sim_{B \rightarrow A}) / 2$ 
30:   end if
31:
32:    $\triangleright$  Step 3: Combine scores for the final similarity
33:    $TotalSim \leftarrow w_p \cdot Sim_p + w_i \cdot Sim_i$ 
34:   return  $TotalSim$ 
35: end function
```

---

complex, structured learning plans, filling a void in existing evaluation methodologies.

## 4 EXPERIMENTS

To evaluate the effectiveness of GraphMASAL's output items, we conduct both LLM-based automated and human evaluations [20].

### 4.1 Implementation and Experimental Setup

**System.** We implement GraphMASAL in Python. Neo4j serves as the persistent graph store, supporting dynamic graph operations including node insertion, edge pruning, and relation rewiring under agent control. Workflow orchestration is realized with LangGraph [4]. All inference is powered by a commercial large language model, with doubao-seed-1.6-250615 as the reasoning backbone. For retrieval and entity grounding over the knowledge graph [5],

we use doubao-embedding-text-240715 to index concept nodes, prerequisite relations, and item texts, enabling vector similarity search and semantic-level knowledge retrieval. We parameterize  $\alpha=0.5$ ,  $K=6$ , and enumerate shortest paths up to length 10 before greedy selection; all hyper-parameters are released with our code.

**Data.** Our primary experimental assets consist of a curated high-school physics knowledge graph derived from publicly available textbooks: 101 concept nodes, 168 prerequisite edges, and 200 exercises drawn from the open-source multi-modal, multi-subject problem set `mllm_multi_subject_data`. Each exercise is linked to one or more concept nodes and is accompanied by a misconception map enumerating common distractors and their underlying misunderstandings.

**Knowledge graph construction.** The knowledge graph is constructed from publicly available textbooks through an automated extraction pipeline. Concept nodes and prerequisite dependencies are extracted using doubao-seed-1.6-250615 with Chain-of-Thought prompting (CoTPrompt) [18], which guides the model to identify domain concepts and their logical prerequisites through structured reasoning. To ensure graph quality and efficiency, we apply deduplication over semantically similar concepts based on embedding cosine similarity (threshold  $> 0.9$ ), yielding a concise, non-redundant knowledge structure that supports efficient path planning and retrieval.

**Simulated students.** To obtain controlled, repeatable evaluations, we synthesize 15 distinct student profiles. Each profile specifies an initial mastery vector over the 101 concepts. For every profile, we simulate a 15-step answering sequence. At each step, the correctness probability is sampled as a monotone function of the student’s current mastery over the exercise-linked concepts, thereby inducing individualized learning trajectories with realistic variance.

**Hyper-parameters and infrastructure.** temperature is set to 0.2 and top\_p to 0.9 for deterministic reasoning. The context window accommodates the retrieved top-k knowledge snippets ( $k = 6$ ), the agent’s hidden chain-of-thought summary, and the tool outputs. All experiments are executed on a single workstation with standard CPU and GPU resources; retrieval and graph operations are I/O-bound and do not require specialized accelerators.

## 4.2 Automated Evaluation

Evaluating tutoring systems is challenging due to the open-ended nature of the tasks. To address this, we designed a scalable, automated evaluation pipeline based on quantitative metrics, which is then validated against an expert proxy.

**4.2.1 Overall Experimental Setup.** Our core evaluation methodology is a blind reasoning test. Crucially, we define the ground truth as an *oracle baseline*—not an external gold standard from human experts, but rather the output of our own MSMS planner running with complete metadata access. Specifically, for each simulated student attempt, we execute the "teacher" version of GraphMASAL that has access to all problem annotations (linked\_kp\_ids and misconception\_map), thereby producing the optimal diagnosis and learning path under perfect information. We then create a "blinded" version by removing these critical fields, forcing the deployed GraphMASAL to perform genuine reasoning and retrieval

**Table 1: Comparison of evaluation results for Cognitive Diagnosis.**

System	Precision	Recall	F1 Score
DirPrompt	0.46	0.57	0.51
CoTPrompt	0.54	0.60	0.57
<b>GraphMASAL</b>	<b>0.80</b>	<b>0.69</b>	<b>0.74</b>
w/o KG	0.62	0.64	0.63

based solely on problem content. This setup measures how well our system reconstructs the oracle decision when key information must be inferred rather than directly retrieved.

**4.2.2 Evaluating Cognitive Diagnosis.** To rigorously evaluate the cognitive diagnosis capability of GraphMASAL, we conducted a systematic comparison against two representative baseline approaches: (1) Direct Prompting (DirPrompt), wherein the language model is tasked with diagnosing student misconceptions through straightforward task instructions, and (2) Chain-of-Thought Prompting (CoTPrompt) [18], which incorporates explicit reasoning chains into the prompt to guide the model’s diagnostic process. Both baselines leverage the same underlying large language model but differ in their prompting strategies. The diagnostic task requires the system to identify which concepts a student has mastered based on their problem-solving behaviors, a fundamentally challenging classification problem due to the fine-grained nature of concept-level attribution and the limited observable evidence from each student’s answer history. Our experimental results, as shown in Table 1, reveal a substantial performance gap between GraphMASAL and the prompt-based baselines. GraphMASAL significantly outperforms both DirPrompt and CoTPrompt across precision, recall, and F1-score metrics. This marked improvement demonstrates that GraphMASAL’s integration of a dynamic knowledge graph, specialized diagnostic agent, and structured graph operations provides a more robust and reliable foundation for cognitive diagnosis than purely prompt-engineering approaches, which lack persistent memory of student states and struggle to reason over complex prerequisite dependencies.

**4.2.3 Evaluating Learning Paths.** We benchmark path planning against LLM prompting and structured ablations using three metrics: PathSim ( $\uparrow$ ; structural/sequence alignment), Coverage ( $\uparrow$ ; proportion of weak concepts covered), and Total Cost ( $\downarrow$ ; unique concepts to learn). For each of 15 simulated student profiles, all methods are evaluated under matched inputs (identical weak-concept set/targets, retrieval Top-K, and context window) and matched inference settings (same LLM backbone, fixed temperature and top\_p). To mitigate stochasticity, every configuration is repeated with 5 random seeds; we report mean  $\pm$  95% CI over seeds  $\times$  profiles. Statistical significance is assessed via paired two-sided tests on per-profile scores. MSMS (ours) attains the highest PathSim and Coverage while minimizing Total Cost, with improvements over all baselines that are statistically significant under paired tests in most profiles. *Shortest-per-sink* achieves near-perfect Coverage but incurs higher cost and lower PathSim due to lack of global de-duplication. Removing the branching penalty (*No-fanout*) increases redundancy,

**Table 2: Path planning baselines: PathSim ( $\uparrow$ ), Coverage ( $\uparrow$ ), Total Cost ( $\downarrow$ ).**

Method	PathSim	Coverage	Total Cost
<i>LLM-based</i>			
CoT Prompt	$0.71 \pm 0.12$	0.82	$16.6 \pm 4.2$
<i>Structured / Ablations</i>			
w/o KG	$0.63 \pm 0.15$	0.78	$18.3 \pm 4.7$
Bi-encoder only	$0.74 \pm 0.10$	0.88	$15.5 \pm 3.9$
Random Path	$0.52 \pm 0.18$	0.76	$21.9 \pm 5.6$
Shortest-per-sink	$0.79 \pm 0.08$	0.96	$17.8 \pm 3.1$
No-fanout penalty	$0.81 \pm 0.09$	0.97	$18.6 \pm 3.3$
<i>Proposed</i>			
<b>MSMS (ours)</b>	<b><math>0.857 \pm 0.07</math></b>	<b>0.98</b>	<b><math>14.2 \pm 2.8</math></b>

raising cost and slightly lowering PathSim. *Bi-encoder only* outperforms *w/o KG* yet lags MSMS without reranking/global optimization. *w/o KG* and *Random Path* yield the lowest alignment and coverage, with the latter exhibiting the highest cost. *CoT Prompt* forms reasonable paths but underperforms MSMS on structure alignment and cost control.

### 4.3 Human Validation of the Automated Evaluation

To validate the reliability of our automated evaluation framework, we conducted a correlation analysis between automated metrics and human expert judgments. We randomly sampled 50 diagnosis-planning pairs from the experimental dataset. For each pair, three experienced physics educators independently rated the quality on a 5-point Likert scale, with ratings averaged as the human judgment baseline.

Table 3 presents the Pearson correlation coefficients. Both diagnostic and planning quality show moderate-to-strong positive correlations with human ratings, with statistical significance at  $p < 0.01$  and  $p < 0.001$  respectively. Critically, human validation confirms that the oracle baseline itself—the idealized plan under perfect information—possesses pedagogical validity. This forms a *logical closure*: our system reliably generates paths that closely match the oracle baseline, and human experts independently judge these oracle-aligned paths to be of high quality. Together, these findings validate not only the soundness of our automated metrics but also the pedagogical merit of the decision strategy our system consistently recovers.

**Table 3: Correlation between human expert ratings and automated metrics (N=50).**

Evaluation Task	Correlation (r)	p-value
Diagnostic Quality	0.65	$< 0.01$
Planning Quality	<b>0.68</b>	$< 0.001$

## 5 CONCLUSION

This work introduced GraphMASAL, a graph-centric multi-agent ITS that integrates: (i) a dynamic knowledge graph for persistent, stateful learner modeling; (ii) a LangGraph-based orchestration layer enabling interpretable, tool-augmented collaboration among Diagnoser, Planner, and Tutor agents; (iii) a KG-grounded two-stage neural IR component (dual-encoder dense retrieval with cross-encoder listwise re-ranking and calibrated score fusion) for concept and item grounding; and (iv) a multi-source multi-sink (MSMS) planning engine with a cognitively grounded cost and a  $(1 + \ln |W|)$  approximation via greedy set cover over source→sink paths. Across blinded evaluations and expert validations, GraphMASAL yielded substantially higher diagnostic fidelity ( $F1 = 0.74$ ) and stronger structural alignment of learning plans (mean PathSim = 0.857) with positive correlations to human judgments (diagnosis  $r=0.65$ , planning  $r=0.68$ ), indicating both computational effectiveness and pedagogical plausibility. Beyond empirical gains, our formulation clarifies how explicit graph structure, optimization under educational constraints, and modular agent tooling jointly contribute to reliability and interpretability in ITS. Future work includes principled incorporation of multimodal resources into the KG, tighter coupling of affect/motivation signals with planning objectives, and longitudinal field deployments to quantify learning gains and equity across populations.



## REFERENCES

- [1] Abderrahim Abyaa, Mohammed Khalidi Idrissi, and Samir Bennani. 2019. Learner Modelling: Systematic Review of the Literature from the Last 5 Years. *Educational Technology Research and Development* 67 (2019), 1105–1143.
- [2] Areej Alkhatlan and Jugal Kalita. 2018. Intelligent Tutoring Systems: A Comprehensive Historical Survey with Recent Developments. *arXiv preprint arXiv:1812.09628* (2018).
- [3] Václav Chvátal. 1979. A Greedy Heuristic for the Set-Covering Problem. *Mathematics of Operations Research* 4, 3 (1979), 233–235.
- [4] Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V. Chawla, Olaf Wiest, and Xiangliang Zhang. 2024. Large Language Model Based Multi-Agents: A Survey of Progress and Challenges. In *Proc. of the 33rd International Joint Conference on Artificial Intelligence (IJCAI)*.
- [5] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. In *Advances in Neural Information Processing Systems*, Vol. 33. 9459–9474.
- [6] Haoran Li, Rwitajit Majumdar, Morris Ricky Appelman Chen, and Hiroaki Ogata. 2021. Goal-Oriented Active Learning (GOAL) System to Promote Reading Engagement, Self-Directed Learning Behavior, and Motivation in Extensive Reading. *Computers & Education* 171 (2021), 104239.
- [7] Benjamin D. Nye, Darryl Mee, and Mark G. Core. 2023. Generative Large Language Models for Dialog-Based Tutoring: An Early Consideration of Opportunities and Concerns. In *LLM@AIED*. 78–88.
- [8] Joon Sung Park, Joseph C. O'Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. Generative Agents: Interactive Simulacra of Human Behavior. In *Proc. of the 36th Annual ACM Symposium on User Interface Software and Technology (UIST)*.
- [9] Chris Piech, Jonathan Bassen, Jonathan Huang, Surya Ganguli, Mehran Sahami, Leonidas J. Guibas, and Jascha Sohl-Dickstein. 2015. Deep Knowledge Tracing. In *Advances in Neural Information Processing Systems*, Vol. 28.
- [10] Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, Juyuan Xu, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2024. ChatDev: Communicative Agents for Software Development. In *Proc. of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL)*, Volume 1: Long Papers.
- [11] Petr Slavík. 1997. A Tight Analysis of the Greedy Algorithm for Set Cover. *Journal of Algorithms* 25, 2 (1997), 237–254.
- [12] Robert A. Sottolare, Arthur Graesser, Xiangen Hu, and Heather Holden. 2013. *Design Recommendations for Intelligent Tutoring Systems: Volume 1 - Learner Modeling*. Technical Report. US Army Research Laboratory.
- [13] John Sweller, Paul Ayres, and Slava Kalyuga. 2011. Measuring Cognitive Load. In *Cognitive Load Theory*. Springer, 71–85.
- [14] Jeroen J. G. Van Merriënboer and Paul A. Kirschner. 2017. *Ten Steps to Complex Learning: A Systematic Approach to Four-Component Instructional Design* (3rd ed.). Routledge.
- [15] Kurt VanLehn. 2006. The Behavior of Tutoring Systems. *International Journal of Artificial Intelligence in Education* 16, 3 (2006), 227–265.
- [16] Jindong Wang, Lei Zou, Wei Wang, and Zhengxu Wang. 2021. A Survey on Knowledge Graphs for Education. *Tsinghua Science and Technology* 26, 5 (2021), 710–724.
- [17] Sheng Wang, Tingkai Xu, Hongru Li, Chengwen Zhang, Jie Liang, Jiliang Tang, Philip S. Yu, and Qingsong Wen. 2024. Large Language Models for Education: A Survey and Outlook. *arXiv preprint arXiv:2403.18105* (2024).
- [18] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. In *Advances in Neural Information Processing Systems*, Vol. 35. 24824–24837.
- [19] Yanan Yun, Haoming Dai, Ran An, Yang Zhang, and Xuequn Shang. 2023. Doubly Constrained Offline Reinforcement Learning for Learning Path Recommendation. *Knowledge-Based Systems* (2023), 111242.
- [20] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhenhao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena. In *Advances in Neural Information Processing Systems*, Vol. 36.