# GAT-CL: A Graph Attention and Contrastive Learning Framework for Multimodal Behavior Modeling in Intelligent Education Systems

Lei Zhu
Henan Vocational College of Tuina, Luoyang 471023, Henan, China
E-mail:Zhulei1123@outlook.com

*This study addresses the limitations of intelligent education systems in multimodal data fusion, scalability, and robustness by proposing a graph-based cognitive modeling framework enhanced with contrastive representation learning. Using interaction data from 186 students and 874,520 records over a semester, heterogeneous behavior graphs are constructed and encoded with a Multi-Head Graph Attention Network (GAT) to capture semantic and temporal dependencies. A contrastive learning module further strengthens embedding robustness, and the optimized representations drive a dynamic strategy engine for adaptive instructional resource allocation. Experimental results demonstrate 93.2% accuracy in learner behavior classification and 90.1% accuracy in clickstream prediction, with a 15.4% improvement in disengagement-signal retention compared to GCN, LSTM, Transformer, and GraphCL baselines. These findings validate the effectiveness and transferability of combining cognitive graph modeling with contrastive learning, advancing both theoretical foundations and practical capabilities of intelligent education systems to reduce dropout risk and enhance engagement.*

*Povzetek: Študija pokaže, da lahko sistem, ki učenje modelira kot "mrežo povezav" in se uči bolj robustnih predstavitev, bolje napove vedenje študentov ter pomaga prej zaznati upad motivacije in tveganje za odpad.*

## 1 Introduction

With the rapid development of information technology and the acceleration of digital education reforms, intelligent education systems have become key platforms for enhancing instructional efficiency and optimizing learner engagement. The integration of multimedia resources into teaching has created multimodal, interactive, and immersive learning environments [1–4]. These environments, however, generate fragmented, nonlinear, and high-frequency behavioral data that impose higher demands on adaptive content scheduling and intelligent responsiveness [5–6]. Conventional sequential models, such as long short-term memory networks (LSTM) and convolutional neural networks (CNN), have been widely adopted to capture temporal dependencies in learning behaviors [9–10]. Yet, they struggle to represent complex structural relationships, overlook latent graph-like patterns in behavioral sequences, and fail to adequately exploit multimodal synergies. As a result, personalized strategy generation remains unstable [11–12], particularly due to the lack of semantic alignment between behavior features and pedagogical content.

Recent efforts have attempted to improve adaptive instruction through hybrid management of behavioral data [17–23] and multimedia-based teaching integration [24–

28]. For example, Lee et al. [17] mapped learner behavior to the ICAP framework using deep learning; Zhao et al. [20] proposed a result-confirmation approach to interpret e-book reading patterns; and Cui [24] developed a multimedia teaching model for personalized language learning. While these advances enhanced interpretability and personalization, they still lack scalability across interdisciplinary, media-rich environments [25–26]. Moreover, static profile- or rule-based recommendation modules [15–16] are limited in dynamic adaptability, often resulting in poor content matching and ineffective feedback loops.

To overcome these limitations, graph-based methods have gained momentum in modeling the complex dependencies of learner behaviors. Graph neural networks (GNNs), particularly Graph Attention Networks (GAT), have demonstrated strong capabilities in capturing semantic proximities and structural relations [29–30]. In parallel, graph contrastive learning (GCL) has emerged as a powerful paradigm for enhancing embedding discriminability by leveraging subgraph alignment and perturbation strategies [31–34]. Recent surveys [32] and studies [33–34] highlight its ability to improve robustness in noisy, heterogeneous data environments. Similarly, the rise of Transformer-based multimodal models has provided promising tools for adaptive and inclusive education, integrating vision, text, and behavioral

modalities [35–37]. Applications range from multimodal attention modeling in educational intelligence [36] to domain-specific advising systems [37], underscoring the trend toward scalable multimodal fusion in educational AI.

Despite these advances, several challenges remain. Current models often focus on single-modality or low-dimensional behaviors, which limits their scalability and generalization across large-scale heterogeneous environments [42–44]. Existing frameworks also lack sufficient alignment between cognitive features and pedagogical strategies, thereby weakening interpretability and adaptability [41, 45]. Moreover, although graph contrastive learning and multimodal Transformers are rapidly evolving, their integration into dynamic, real-time educational systems has yet to be systematically explored [35, 42, 46].

To address these challenges, this paper proposes a graph-based multimodal behavior modeling and adaptive strategy optimization framework that integrates Graph Attention Networks with contrastive learning. Learner interaction data—including clickstreams, dwell times, access paths, and interaction frequencies—are encoded into heterogeneous behavior graphs. Multi-head GAT captures semantic and temporal correlations, while a contrastive learning module refines embeddings through positive–negative subgraph discrimination. The optimized representations feed a dynamic strategy engine that generates personalized instructional interventions in real time. In doing so, this study introduces a scalable graph–contrastive learning framework for multimodal learner modeling in media-rich education, provides empirical evidence on a large-scale dataset comprising 186 students and 874,520 interactions with significant performance gains over state-of-the-art baselines, and advances theoretical understanding of how cognitive graph modeling and contrastive learning jointly enhance the precision, interpretability, and adaptability of intelligent education systems.

Table 1: Comparative summary of prior studies

| Study & Year | Method | Dataset | Metrics | Key Limitation |
|---|---|---|---|---|
| Xuan (2022) [9] | DRN-LSTM | Classroom behaviors | Accuracy (85%) | Weak in structural modeling |
| Li et al. (2021) [10] | CNN for behavior recognition | Teaching videos | Precision/Recall | Ignores multimodal inputs |
| Zhao et al. (2021) [20] | ReCoLBA (result-confirmation) | E-book reading logs | Interpretability | Limited to single domain |
| Lee et al. (2023) [17] | DL + ICAP framework | STEM education | Accuracy (92%) | Focused on small-scale, domain-specific data |
| Liu et al. (2021) [23] | Hybrid learning management | Mgmt. courses | Engagement | Lacks scalability |
| GraphCL (2023) [34] | Graph Contrastive Learning | Benchmark graphs | Representation quality | Not applied to education |
| Wu et al. (2024) [33] | Cohesive subgraph GCL | Large graph datasets | Robustness | No education-specific validation |
| Bharathi et al. (2025) [35] | Multimodal Transformer | e-Learning | Engagement, Inclusiveness | Expensive, data-heavy |
| Xia & Niu (2025, Informatica) [38] | Transformer + Bi-LSTM | Vaccine sentiment tweets | Accuracy, F1 | Non-educational domain |
| Ji & Cao (2025, Informatica) [39] | Transformer fusion | Video forgery detection | Precision | Non-educational, but shows multimodal fusion potential |

As summarized in Table 1, prior research has advanced temporal modeling, interpretability, and multimodal integration in intelligent education systems. Nevertheless, sequential models often fail to capture the graph-like dependencies embedded in learner behaviors, interpretability-driven frameworks lack scalability across diverse contexts, and Transformer-based multimodal approaches remain computationally intensive while seldom linked to adaptive teaching strategies. Consequently, a critical research gap remains: few studies integrate graph-based modeling, contrastive learning, and adaptive strategy generation within large-scale, real-world educational settings. Addressing this gap constitutes the central contribution of the present work. Specifically, this study investigates how graph-based contrastive learning can enhance the robustness and scalability of multimodal learner behavior modeling, with the hypothesis that Graph Attention Networks combined with contrastive learning embeddings will outperform sequential and unimodal baselines in prediction accuracy, representation robustness, and learner engagement. Success is defined by achieving at least a 5% improvement over state-of-the-art baselines in behavior classification, demonstrating statistically significant gains ($p < 0.05$) in clickstream prediction and disengagement-signal retention, and validating adaptability in real-world, media-rich higher education datasets.

# 2 Design of behavior modeling and strategy generation method

## 2.1 Extraction of learning behavior features and construction of behavior graph

In the stage of extracting learning behavior features and building behavior graph, the behavior log data generated by the media teaching system is set as the original input, and the behavior event sequence is set as $S = e_1, e_2, \ldots, e_T$, where $e_t$ represents the behavior event of the learner at time t. Each behavior event $e_t$ is represented as a triple $e_t = (a_t, r_t, \tau_t)$, where $a_t$ represents the behavior action type, $r_t$ represents the resource identifier corresponding to the behavior, and $\tau_t$ is the timestamp of the behavior. According to the semantic normalization dictionary and the time density distribution function, $a_t$ and $r_t$ are discretized and mapped to define a unified behavior category space $\mathcal{A}$ and resource space $\mathcal{R}$. After mapping, the behavior events are uniformly embedded in a fixed-dimensional vector form.

A directed graph G = (V, E) can be constructed as the expression of the behavior graph structure. The node set V consists of all the behavior events of the learner in a certain time window. Assume that the sliding window size is $\omega$, and a sliding mechanism with a step size of $\delta$ is used in the behavior sequence to construct the graph structure for the continuous event segments, satisfying $|V| \leq \omega$ and ensuring that the graph structure has temporal integrity under the constraints of space complexity. In the figure, each edge $e_{ij} \in E$ connects event nodes $v_i$ and $v_j$. The strength of the edge is defined by the edge weight function $w_{ij}$. The weight calculation adopts the joint temporal-semantic mechanism, as shown in formula (1):

$$w_{ij} = \lambda_1 \cdot \text{sim}(a_i, a_j) + \lambda_2 \cdot \exp(-\gamma_1 |\tau_i - \tau_j|) \tag{1}$$

$\text{sim}(a_i, a_j)$ represents the semantic similarity function between action types, which is calculated using the cosine similarity of the embedding vector. $\tau_i$ and $\tau_j$ represent the timestamps of the corresponding actions, $\lambda_1$ and $\lambda_2$ are weighting coefficients, and $\gamma_1$ is the time decay factor, which controls the sensitivity of the edge weight to the change of time interval.

In order to suppress the risk of noise propagation caused by excessive edge connection density, structural filtering rules are introduced. Define the edge threshold $\theta_w$, if $w_{ij} < \theta_w$, discard the corresponding edge connection; at the same time, set the node degree upper limit $D_{max}$, if a node degree exceeds the upper limit, retain the $D_{max}$ connection with the highest edge weight, and set other edges to invalid, further limiting the complexity of the behavior graph and ensuring the convergence and stability of subsequent graph neural network calculations. The node representation initialization is achieved by jointly embedding the behavior action type, resource category and time information. Assume that the embedding vectors of $a_i \in \mathcal{A}$ and $r_i \in \mathcal{R}$ are $\mathbf{a}_i$ and $\mathbf{r}_i$ respectively, and the timestamp $\tau_i$ is normalized to the interval $[0,1]$ and embedded as the time vector $\mathbf{t}_i$. Then the initial representation $\mathbf{h}_i^0$ of the behavior node $v_i$ is defined as shown in formula (2):

$$\mathbf{h}_i^0 = \mathbf{W}_a \mathbf{a}_i + \mathbf{W}_r \mathbf{r}_i + \mathbf{W}_t \mathbf{t}_i + \mathbf{b} \tag{2}$$

Among them, $\mathbf{W}_a$, $\mathbf{W}_r$, and $\mathbf{W}_t$ are trainable weight matrices, and $\mathbf{b}$ is a bias term. This representation is passed to the subsequent GAT module as an input node feature vector, and is further used to learn the structural relationship and semantic coupling characteristics between behaviors. This method ensures that the temporal evolution trajectory and semantic association pattern of the learner's behavior are fully preserved during the construction of the behavior graph structure, laying the foundation for subsequent graph representation learning and teaching strategy generation.

## 2.2 GAT-driven behavior representation encoding

In the media teaching scenario, learners' behaviors have complex temporal structures and semantic dependencies. Traditional graph neural networks use average or static weight aggregation for adjacent nodes to hardly characterize the heterogeneous relationship characteristics between nodes. For this reason, GAT is introduced as the encoding mechanism of the behavior graph structure to achieve adaptive weighted learning of node semantic representation while maintaining the topological structure [29-30]. Figure 1 shows the overall composition of the behavior representation encoding module under the GAT structure and the interactive relationship between each functional unit.
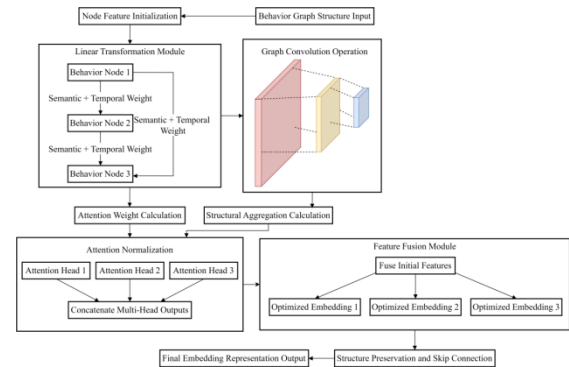


Figure 1: Graph attention encoding framework for learning behavior graph representation

GAT is built on the learning behavior graph structure to achieve deep encoding of behavior representation. The node feature matrix is set to $X \in \mathbb{R}^{N \times d}$, where $N$ is the number of nodes and $d$ is the original feature dimension. The original features are mapped using linear transformation to obtain the node representation $h_i = W x_i$, where $W \in \mathbb{R}^{d' \times d}$ is the trainable weight matrix and $d'$ is the mapped dimension. The attention weights between nodes are calculated based on the local adjacency structure. The attention relevance score $\psi_{ij}$ of node $j$ to node $i$ is obtained by formula (3):

$$\psi_{ij} = \text{LeakyReLU}(\vec{a}^\top [h_i \parallel h_j]) \tag{3}$$

$\vec{a} \in \mathbb{R}^{2d'}$ is a learnable parameter vector, and $\parallel$ represents a feature concatenation operation. To ensure information normalization, the softmax function shown in

formula (4) is introduced to normalize the scoring results in the node neighborhood:

$$\alpha_{ij} = \frac{\exp(\psi_{ij})}{\sum_{k \in \mathcal{N}(i)} \exp(\psi_{ik})} \quad (4)$$

$\alpha_{ij} \in [0,1]$ is the influence strength of node $v_j$ on node $v_i$ during the feature update process, and $\mathcal{N}(i)$ is the set of adjacent nodes of node $v_i$, satisfying $\sum_{j \in \mathcal{N}(i)} \alpha_{ij} = 1$. The final node embedding vector is the attention weighted aggregation result as shown in formula (5):

$$h_i' = f\left(\sum_{j \in \mathcal{N}(i)} \alpha_{ij} h_j\right) \quad (5)$$

$f$ is a nonlinear activation function, and $h_i' \in \mathbb{R}^{d'}$ is the node representation after a layer of GAT update. In order to enhance the model's ability to model multiple semantic channels, a multi-head attention mechanism is used to connect M independent attention subspaces in parallel, and the generated embedding representation $h_i^{\text{multi}}$ is as shown in formula (6):

$$h_i^{\text{multi}} = \|_{m=1}^M f\left(\sum_{j \in \mathcal{N}(i)} \alpha_{ij}^{(m)} h_j^{(m)}\right) \quad (6)$$

Here, $\alpha_{ij}^{(m)}$ and $h_j^{(m)}$ are the attention weight and node feature of the mth attention head, respectively. In the process of stacking multi-layer graph convolution, the semantic representation of the initial node is retained through the skip connection mechanism to alleviate the problem of feature over-smoothing. The formal expression is as follows:

$$\tilde{h}_i^{(l)} = h_i^{(l)} + h_i^{(0)} \quad (7)$$

Among them, $h_i^{(0)}$ is the initial embedding of the node, $h_i^{(l)}$ is the output of the lth layer, and $\tilde{h}_i^{(l)}$ is the final output of the fused residual. This structure not only ensures the local neighborhood expression ability, but also enhances the model's ability to retain and discriminate key nodes in the behavior path. After all nodes in the behavior graph are encoded by multi-layer GAT, a set of embedding representations $\mathbf{H} = \{h_1', h_2', \dots, h_n'\}$ with consistent dimensions is obtained, which serves as the input representation matrix of the subsequent contrastive learning optimization and strategy generation module.

## 2.3 Contrastive learning enhanced behavior embedding optimization mechanism

After the behavior graph representation is embedded by GAT, in order to improve the model's ability to aggregate similar structures in the learning behavior pattern and distinguish heterogeneous structures, a contrastive learning mechanism is introduced to construct an embedding optimization path. In the encoding stage, the training samples are expanded by constructing positive and negative behavior subgraph pairs, and the contrast loss function between graph embeddings is used to further constrain the spatial structure of the behavior representation. A set of positive sample graphs $G^+$ and negative sample graphs $G^-$ are generated through data perturbation, and the corresponding embedding vectors are $\mathbf{h}_v^+$ and $\mathbf{h}_v^-$ respectively. The positive sample graph is obtained by retaining the main nodes of the behavior path structure and perturbing the edge weights, while the

negative sample graph is generated by behavior path clipping and semantic perturbation.

When constructing the loss function, the Euclidean distance between behavior embeddings is used as the similarity metric, and the optimization goal is to minimize the embedding distance between positive samples and maximize the average distance between negative samples. The contrast loss function is defined as formula (8):

$$\mathcal{L}_{\text{contrast}} = \sum_{i=1}^N \left[ \| \mathbf{h}_i - \mathbf{h}_i^+ \|_2^2 - \frac{1}{K}\sum_{k=1}^K \| \mathbf{h}_i - \mathbf{h}_k^- \|_2^2 \right] \quad (8)$$

There are K negative samples in total. By maximizing the difference between the average distance of negative samples and the distance of positive samples, the model's ability to distinguish between aggregations of similar structures and heterogeneous structures is improved.

In order to ensure that the structural comparability between subgraphs can be maintained after the perturbation, the perturbation strategy is constrained to maintain structure. Let the perturbed subgraph be $G' = (V', E')$, and its adjacency matrix $A$ with the original graph $G$ is required to satisfy the maximum structural retention, that is, the control shown in formula (9) is performed during the perturbation process:

$$\| A - A' \|_F \leq \epsilon \quad (9)$$

$\|\cdot\|_F$ represents the Frobenius norm and $\epsilon$ is the upper limit of the perturbation amplitude. The final embedding vector $z_v$ is composed of the weighted combination of the original behavior embedding and the contrast optimized representation, and is defined by formula (10):

$$z_v = \alpha \cdot h_v + (1 - \alpha) \cdot h_v^{\text{contrast}} \quad (10)$$

Among them, $\alpha \in [0,1]$ is the weight parameter, and $h_v^{\text{contrast}}$ represents the optimization vector under the guidance of contrast loss. This embedding serves as the input basis for the subsequent generation of personalized teaching strategies, ensuring its dual robustness in semantic consistency and structural discriminability.

## 2.4 Personalized teaching strategy generation and resource scheduling

In the personalized teaching strategy generation and resource scheduling phase, the system first receives the learning behavior embedding vector optimized by GAT and contrastive learning as the input feature to build the teaching strategy matching model. Let the behavior embedding be denoted as $\mathbf{h}_u \in \mathbb{R}^d$, where $d$ represents the embedding dimension, and the historical behavior sequence embedding is denoted as $\mathbf{h}_u^{(1)}, \mathbf{h}_u^{(2)}, \dots, \mathbf{h}_u^{(t-1)}$. The dynamic feature state of the behavior sequence is extracted through the gated recurrent unit (GRU), and the state output is defined as formula (11):

$$\mathbf{s}_u^{(t)} = \text{GRU}(\mathbf{h}_u^{(t-1)}, \mathbf{s}_u^{(t-1)}) \quad (11)$$

Among them, $\mathbf{s}_u^{(t)}$ is the state vector at the moment of current strategy generation, $\mathbf{h}_u^{(t-1)}$ is the embedding input of the previous step behavior, and $\mathbf{s}_u^{(t-1)}$ is the state at the previous moment. After obtaining the current behavior state, the strategy matching function is designed to realize the personalized recommendation of teaching resources.

The embedding vector of teaching content is represented as $\mathbf{c}_j \in \mathbb{R}^d$, and the matching score calculation function is defined by the bidirectional attention fusion method commonly used in the dual-tower structure as formula (12):

$$\alpha_{u,j} = \sigma\big((\mathbf{W}_1 \mathbf{s}_u^{(t)})^\top (\mathbf{W}_2 \mathbf{c}_j)\big) \quad (12)$$

Among them, $\mathbf{W}_1, \mathbf{W}_2 \in \mathbb{R}^{d \times d}$ are trainable mapping matrices, and $\sigma(\cdot)$ represents the Sigmoid activation function, which is used to map the matching score to the interval [0,1]. All candidate teaching resources are sorted in descending order according to the score $\alpha_{u,j}$, and the top k resources are selected to form the recommendation set.

The resource scheduling module performs feedback path selection based on the above matching results combined with the teaching strategy graph model. The strategy graph structure is defined as $G_s = (V_s, E_s)$, where $V_s$ is the strategy node set and $E_s$ is the strategy transition edge set. Assuming the strategy node state is $\mathbf{v}_i$ and the transition relationship edge weight is $w_{ij}$, the current scheduling path of the system is calculated by the Bellman equation shown in formula (13) to calculate the path optimality function $Q(s,a)$:

$$Q(s,a) = r(s,a) + \gamma_2 \sum_{s'} P(s'|s,a) \max_{a'} Q(s',a') \quad (13)$$

Among them, $r(s,a)$ is the immediate feedback of taking action $a$ under the current state $s$, $P(s'|s,a)$ is the state transition probability, and $\gamma_2$ is the discount factor. The scheduling path is determined according to the principle of maximizing the $Q$ value, and the dynamic push process of teaching resources is finally controlled.

The system control layer maps the teaching content presentation strategy into task execution instructions based on the matching results and scheduling paths, and records the feedback data to update the strategy network. The whole process combines offline strategy pre-training with online fine-tuning to improve the system's response accuracy to changes in learning behavior, and achieves efficient adaptation and intelligent intervention control of teaching content while ensuring that system resource consumption is controllable.

# 3 Experimental setup and system deployment

## 3.1 Experimental platform and media teaching system construction environment

In constructing the experimental platform for the intelligent education system, it is essential to integrate multiple dimensions, including teaching function modules, algorithm deployment strategies, media resource processing, and front–end/back–end interaction design. Such integration ensures not only stable system operation but also flexible scalability in media-rich instructional scenarios. The system environment configuration directly influences both the inference performance of the deployed models and the responsiveness of user interactions, as well as the completeness of resource loading. To enhance reproducibility and transparency, this study reports the actual deployment structure of the proposed system, with the experimental platform configuration summarized in Table 2.

Table 2: Overview of the experimental platform configuration of the media teaching system

| Deployment Module | Hardware/Software Environment | Specifications | Description |
|---|---|---|---|
| Server Host | Windows 10 + WSL | Intel Xeon 2.4GHz×16 | Backend service deployment |
| Frontend Interface | Vue + Element UI | Resolution 1920×1080 | User behavior collection and display |
| Teaching Content Module | FFmpeg + OpenCV | Video encoding H.264 | Media resource loading and conversion |
| Model Service Container | Docker + PyTorch | CUDA 11.8 + cuDNN 8 | GAT model inference and strategy control |

Table 2 summarizes the deployment of the core modules of the intelligent education system at both software and hardware levels, including the server host environment, content processing framework, front-end configuration, and model service tools. Each component is optimized for media teaching tasks to ensure efficient multi-thread scheduling, video rendering, and behavioral data transmission. The system runs on a Windows 10 server with WSL support for deep learning models; the front-end is developed in Vue for interactive display; FFmpeg and OpenCV handle media transcoding and distribution; and Docker containers encapsulate GAT inference and strategy generation services. This deployment strategy enhances system stability, scalability, and resource scheduling efficiency, thereby supporting reproducible and practical evaluation of the proposed framework.

## 3.2 Dataset source and preprocessing process

The learning behavior data comes from the real use environment of a multimedia teaching system deployed in a middle school. The system covers nine teaching classes in three grades of junior high school. The teaching cycle

is a full semester, a total of eighteen weeks. The system is used as a teaching assistance platform for teachers and a self-learning support tool for students in daily teaching. The deployed terminals include teacher control terminal, student interaction terminal and resource service terminal. The data collection module is designed based on log tracking and behavior trigger recording mechanism. The system writes behavior events into the server log database in real time through the back-end interface. At the same time, the compensation synchronization of high-frequency behaviors is guaranteed through local cache to ensure data integrity and stability.

A total of 186 students' learning behavior data were collected, with a total of 874,520 records, covering various interactive behaviors of students in the media teaching process, forming a behavioral sequence set with users as the main index. The collected field types include behavior type code, event trigger timestamp, interaction position coordinate vector, teaching resource unique identifier, system response status code, task completion flag and user identity index. Each piece of data is uniformly constructed into a five-tuple form $(a_i, t_i, l_i, r_i, s_i)$, where $a_i$ represents the behavior event category, $t_i$ represents the trigger time, $l_i$ represents the interface space position vector where the behavior occurs, $r_i$ represents the associated resource identifier, and $s_i$ represents the behavior state code. The size of the system behavior type set is |A|=23, which constitutes a discrete event space. The behavior of each student is sorted by time-based index sequence to form the original sequence set $S = s_1, s_2, \ldots, s_N$, and each sequence $s_j = [(a_1^j, t_1^j), (a_2^j, t_2^j), \ldots, (a_m^j, t_m^j)]$ satisfies the monotonic time-increasing constraint $t_k^j < t_{k+1}^j$.

The behavior feature preprocessing process includes behavior type encoding conversion, behavior frequency normalization, time standardization and position coordinate transformation. The behavior type is converted into a $d$-dimensional vector representation by the mapping function $f_A: A \rightarrow \mathbb{R}^d$, and the initial behavior vector is constructed by one-hot vector embedding. The behavior frequency is normalized by the mean variance normalization method shown in formula (14):

$$x_i' = \frac{x_i - \mu}{\sigma} \quad (14)$$

Among them, $x_i$ is the original frequency statistics, $\mu$ is the mean of all behavior samples of this type, $\sigma$ is the standard deviation, and $x_i'$ is the normalization result. Time standardization adopts the maximum and minimum normalization strategy to transform the timestamp $t_i$ into $t_i' \in [0,1]$. The behavior location vector $l_i$ is encoded according to the spatial area divided by the interface

module and then embedded and transformed to form a fixed-dimensional position representation vector. Resource identifier $r_i$ is unified as a hash index, and behavior status $s_i$ is processed in a discrete classification manner, indicating whether the behavior is completed, whether it is responded to by the system, and whether it triggers an exception.

The data cleaning process strictly follows the three standards of behavior legitimacy, sequence integrity, and structural discriminability. All records with missing behavior status, timestamp conflicts, invalid resource identifiers, or non-teaching behaviors are removed. The behavior sequences whose interval between consecutive behaviors exceeds the upper limit of the maximum response cycle of the system is regarded as an incoherent behavior flow and processed in segments. After filtering, only the sequences whose behavior length is not less than the set threshold $L_{min} = 12$ are retained to ensure the expression density and topological connectivity of the input graph structure. The sliding window strategy is introduced in the construction of the behavior graph. The window size is set to $w = 8$. Only behavior pairs are constructed within the window range to reduce the density of the graph structure. The edge weight is set to a threshold of $\delta = 0.35$, and only the edges of $w_{uv} > \delta$ are retained in the final graph structure to control the size of the edge set and enhance the significance of semantic relationships. Finally, the training sample set and the test sample set are constructed. The sample division is non-overlapping based on the learner identity, with a ratio of 8:2. All samples are saved in the form of graph structure input, and their adjacency matrix $A \in \mathbb{R}^{n \times n}$ and node feature matrix $X \in \mathbb{R}^{n \times d}$ are stored respectively, where $n$ is the number of nodes in the graph, for the graph neural network module to perform behavior representation learning and strategy generation tasks.

## 3.3 Experimental parameter configuration and model training details

In order to verify the stability and effectiveness of the learning behavior modeling method under multiple training conditions, this paper systematically sets and adjusts the core hyperparameters of the model, covering key modules such as graph attention structure, contrastive learning mechanism and training convergence strategy, forming a set of representative parameter combination configuration schemes, as shown in Table 3.

Table 3: Model training parameter setting table

| Parameter | Configuration 1 | Configuration 2 | Configuration 3 | Configuration 4 | Configuration 5 |
|---|---|---|---|---|---|
| Learning Rate | 0.001 | 0.0005 | 0.0001 | 0.0005 | 0.0003 |
| Batch Size | 64 | 128 | 64 | 256 | 128 |
| Number of Attention | 4 | 8 | 4 | 8 | 6 |

| | | | | | |
|---|---|---|---|---|---|
| Heads Number of GAT Layers | 2 | 3 | 2 | 3 | 4 |
| Contrastive Loss Temperature | 0.5 | 0.3 | 0.7 | 0.5 | 0.4 |

Table 3 summarizes the key training settings of the proposed model, including batch size, learning rate, graph depth, and embedding temperature, each adjusted within controlled ranges to analyze their influence on behavior graph representation quality and strategy generation effectiveness. Combined experiments evaluate convergence speed, loss stability, and embedding discriminability under different configurations, establishing a reliable basis for subsequent performance comparison. A unified test set is then used to quantitatively assess behavior recognition, strategy recommendation accuracy, and scheduling efficiency, enabling a comprehensive evaluation of model operation under varying parameter combinations. The resulting performance trends and their implications for parameter selection in practical deployment are reported in Table 4.

Table 4: Parameter configuration performance comparison table

| Configuration | Learning Behavior Recognition Accuracy (%) | Teaching Content Matching Score (/1.0) | System Response Time (seconds) |
|---|---|---|---|
| Configuration 1 | 89.6 | 0.874 | 2.31 |
| Configuration 2 | 91.2 | 0.912 | 2.48 |
| Configuration 3 | 87.9 | 0.851 | 2.05 |
| Configuration 4 | 90.7 | 0.894 | 2.62 |
| Configuration 5 | 90.1 | 0.902 | 2.28 |

Table 4 compares the performance of different parameter configurations across recognition accuracy, content matching, and system response efficiency. Configuration 2 achieves the best overall balance, showing stable graph embedding performance that enhances semantic separation of behaviors and yields the highest matching score in strategy generation. Although Configuration 3 improves scheduling efficiency, its behavior representation quality is weaker, confirming that accuracy and adaptability are more consistently supported under Configuration 2.

The effectiveness of Configuration 2 derives from a dual reinforcement mechanism of structural perception and semantic discrimination. Specifically, setting the number of attention heads to 8 strengthens the capture of complex semantic associations, while a three-layer GAT deepens the extraction of higher-order features to retain long-term behavior patterns. The contrastive loss temperature (0.3) compresses embedding distances between positive and negative samples, improving discriminability, and a batch size of 128 balances stability with diversity in graph pair construction. With a learning rate of 0.0005, gradient descent remains stable during updates, avoiding oscillations in training. Together, these settings optimize recognition accuracy, response efficiency, and system adaptability, providing empirical guidance for future system deployment.

## 4 Result analysis

### 4.1 Comparative analysis of learning behavior modeling accuracy

To comprehensively evaluate the applicability and performance of the proposed method across diverse behavior modeling tasks, a comparative experimental group was constructed, including five representative models: RNN, LSTM, GCN, GAT, and the proposed GAT+Contrastive Learning (CL). RNN provides basic sequence modeling for short-term dependencies, while LSTM enhances long-range memory through gating mechanisms. GCN aggregates global neighbor information for coarse-grained structural modeling, whereas GAT incorporates attention weighting to capture fine-grained local semantics. Building on this, our GAT+CL approach introduces positive and negative behavior subgraph pairs to refine embedding boundaries, thereby improving behavioral discriminability.

The evaluation considers four behavior types—clicks (interest dynamics), video viewing (deep engagement), resource downloading (content value judgment), and bounce behavior (short-term exit, most challenging to model). Performance is assessed through five indicators: recognition accuracy, structure retention ratio, path consistency, semantic separation, and transfer prediction accuracy. This multidimensional comparison highlights differences in recognition effectiveness, structural preservation, sequence stability, embedding clarity, and predictive adaptability. The results are illustrated in Figure 2.
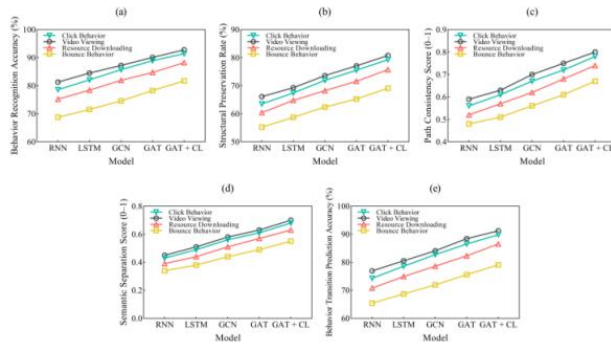
Figure 2: Comparison of five behavior modeling indicators of different models under multiple behavior types
Figure 2 (a): Comparison of behavior recognition accuracy
Figure 2 (b): Comparison of graph structure retention ratio
Figure 2 (c): Comparison of path consistency
Figure 2 (d): Comparison of semantic separation
Figure 2 (e): Comparison of behavior transfer prediction accuracy

As shown in Figure 2, the proposed GAT with contrastive learning achieves optimal performance across all four behavior types. In recognition accuracy, video-viewing behavior improves to 92.8%, benefiting from the weighted adjacency modeling of GAT and the enhanced category boundaries provided by contrastive learning. For structure retention, RNN and LSTM perform poorly on bounce behavior, whereas GAT+CL maintains 69.1%, reflecting better adaptability to sparse graph structures. Path consistency also improves from 0.75 with GAT alone to 0.80 with GAT+CL, indicating more stable global representations. Semantic separation strengthens progressively, rising from 0.34 with RNN to 0.55 with GAT+CL, supported by the discriminative effect of negative sample pairs. Behavior transfer prediction further demonstrates consistent gains, with click-type behavior reaching 89.8%, driven by enhanced trajectory modeling and temporal edge stability. Overall, these results confirm that integrating graph structure with contrastive learning substantially improves behavioral modeling across multiple evaluation dimensions.

## 4.2 Separability of embedded feature space

The experiment constructed a comparative experiment with three processing stages, corresponding to the original embedding, GAT embedding and GAT combined with contrastive learning embedding, to further verify the impact of the structural optimization method on the separability of the behavior embedding space. In the experiment, the principal component analysis method is uniformly used to reduce the ten-dimensional embedding vector to a two-dimensional space to eliminate the visual bias caused by the dimensional difference and maintain the comparability of different stages in the same embedding space. The behavioral data is divided into three categories of labels. By comparing the changes in the two-

dimensional distribution after the three embeddings, the synergy of structural modeling and contrastive learning in enhancing the discriminative ability of behavioral representation is revealed. The visualization results are shown in Figure 3.
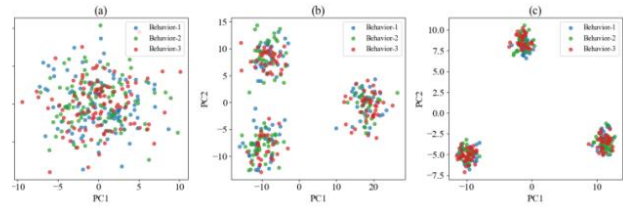


Figure 3: Comparison of two-dimensional visualization of behavior embedding space evolution
Figure 3 (a): Original embedding distribution
Figure 3 (b): GAT embedding distribution
Figure 3 (c): GAT+ contrast learning embedding distribution

As shown in Figure 3, the original embeddings display heavy overlap among the three behavior categories in two-dimensional space, with no clear boundaries and extensive mixing along the first two principal components. After introducing the graph attention mechanism, category boundaries become more distinct, and the overlap between behavior 1 and behavior 3 along principal component 1 is notably reduced, reflecting the enhanced structural perception of behavioral differences. With the further integration of contrastive learning, the three behavior categories are clustered into compact, well-separated sub-regions, demonstrating improved intra-class consistency and inter-class separability. These results confirm that combining structural modeling with contrastive learning substantially strengthens the discriminative power of behavior embeddings.

This improvement arises from the progressive enhancement of the embedding representation mechanism. While the original embeddings rely only on basic feature generation and lack semantic structural modeling, graph attention introduces multi-head correlation weights that reinforce valid connections among similar behaviors and suppress noisy links, thereby preserving local structural information. Contrastive learning further generates positive and negative subgraphs through structural perturbations, guiding the model to optimize intra-class similarity and inter-class distinction. This dual mechanism reduces redundancy, sharpens decision boundaries, and produces a high-density, low-overlap distribution in the embedding space—ultimately improving clustering quality and representation discriminability.

## 4.3 Effect of strategy recommendation path matching

In order to further evaluate the path fitting ability of the system in the strategy generation link, this experiment constructed a node similarity matching matrix between the target teaching task path and the recommended path generated by the system. The target path consists of six

standard teaching behaviors, which are entering the teaching task page, playing video explanation resources, browsing extended reading materials, participating in quizzes or small exercises, viewing system feedback reports, and returning to the task homepage and marking completion. This reflects the entire process of the idealized media teaching process from resource reception, content digestion to task closure. Correspondingly, the recommended path is dynamically generated by the system based on the behavior graph embedding, which contains five strategic behavior nodes, namely, entering the task homepage, clicking and playing video resources, reading recommended document materials, completing personalized recommendation tests, viewing recommendation feedback and jumping to the homepage. The above node sequence is mapped to a structural semantic path, and the behavior matching matrix is constructed by calculating the semantic similarity between the nodes in the recommended path and the target path. The results are shown in Figure 4.
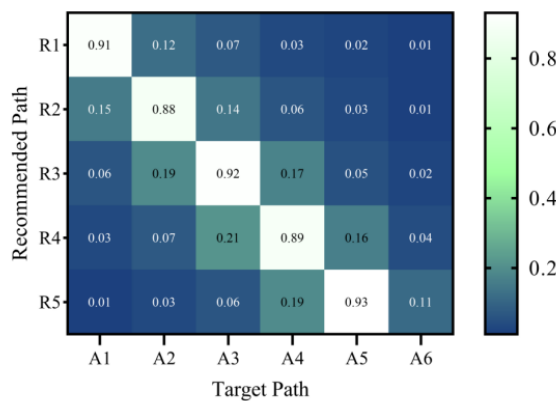


Figure 4: Node similarity heat map of recommended path and target task path

In Figure 4, the matching score between the fifth node of the recommended path and the target path reaches 0.93, representing the highest intensity region and indicating that the system achieves strong accuracy in strategy generation during the feedback presentation stage. This stage typically involves concentrated learner responses after task completion, where behaviors are relatively stable, allowing the system to form more consistent embeddings in graph representation learning. The third node shows a high matching score of 0.92, attributed to the stability of resource structures and the clear behavior patterns in the extended reading link, which

enhances the discriminative capacity of the GAT in capturing semantic representations. The fourth node records a score of 0.89, slightly lower due to the presence of multiple triggering modes in the test behaviors, which introduces local deviations in embedding. For the first node, although the score reaches 0.91, the high-heat zone is narrowly concentrated, suggesting that the initial behavior stage is structurally clear but semantically lightweight, making recognition rely more on structural rather than semantic similarity.Overall, the heat map demonstrates that the system achieves its highest strategy path fitting performance in the mid-to-late stages of the learning task, where behavioral patterns are richer and more stable.

## 4.4 Media resource response delay and scheduling efficiency

In order to deeply explore the resource response characteristics and scheduling efficiency performance of the intelligent education system in different media teaching task scenarios, this experiment designed five representative teaching interaction scenarios. By simulating five typical operating states: single-user access, small-scale return visits, large-scale interaction, multi-modal switching, and platform-level high load, the system's response behavior under different behavior complexities and request densities was comprehensively analyzed. Scenario 1 corresponds to the teaching content request after a single user enters the platform for the first time. The system is in the initial loading state, with concentrated resource demand and delay sensitivity. Scenario 2 simulates the return visit process of students in small classes. The platform can rely on historical cache to achieve moderate resource reuse. Scenario 3 involves real-time interaction in large classes. The number of users surges, resource requests are frequent, and the system's concurrent processing capabilities face significant challenges. Scenario 4 is set as multi-modal task switching. Users need to frequently switch between video, graphics, and interactive modules to test the flexibility of the system's dynamic scheduling. Scenario 5 builds a platform-level stress test scenario to create extreme loads through concentrated high-frequency access to test the system's response elasticity and scheduling stability under resource limits. Figure 5 shows the characteristics of media resource response delay and system resource utilization under three types of request conditions in five scenarios.
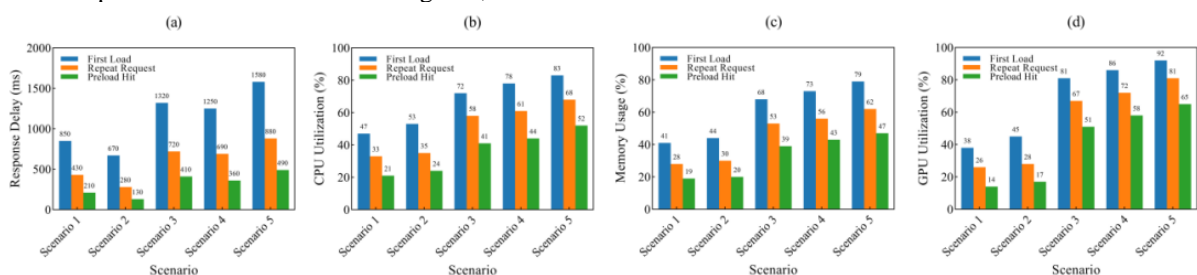


Figure 5: Combined analysis of media resource response and scheduling efficiency
Figure 5 (a): Analysis of media resource response delay

Figure 5 (b): Analysis of CPU utilization
Figure 5 (c): Analysis of memory utilization
Figure 5 (d): Analysis of GPU occupancy

As shown in Figure 5, the first-load delay of media resources increases from 850 ms in Scenario 1 to 1580 ms in Scenario 5 as scenario complexity and access intensity rise. This is primarily due to resource location, permission verification, and data distribution processes, which are more vulnerable to competition and delay under high concurrency. In contrast, delays under repeated requests are markedly lower—for example, 280 ms in Scenario 2 and 880 ms in Scenario 5—benefiting from cache scheduling and connection reuse, where higher cache hit rates yield faster responses. With preload hits, resource preparation is completed in advance, resulting in the lowest response times across all scenarios, demonstrating the effectiveness of preload optimization under high-frequency access.

From a resource perspective, CPU and GPU utilization peak in Scenario 5 at 83% and 92% respectively during the first load, reflecting the heavy computational and image processing demand in high-load conditions. Memory usage, however, remains relatively stable, indicating that pre-allocation and reuse mechanisms are effective. Overall, the results show that the preloading mechanism consistently improves response efficiency, while system scalability in multimodal, high-concurrency scenarios depends heavily on CPU and GPU resources.

## 4.5 System stability analysis under different behavior complexity scenarios

To further verify the response stability of the constructed model in the face of various learning behavior complexity scenarios, the experiment examines its operating performance during task execution through system-level monitoring experiments. The experimental design divides the behavior graph into levels according to the complexity of the structure. Ten levels of behavior complexity are set, from low to high, representing the gradual enhancement of the learning behavior graph in structural dimensions such as the number of nodes, edge density, path branching, and interaction frequency. By loading the corresponding level of behavior graph input in the simulation platform, the average response time, resource occupancy level and teaching task completion rate of the system at each level are recorded. A stability evaluation index system can be constructed to observe the real-time scheduling capability and robustness performance of the system when dealing with high-frequency and high-coupling behavior paths. The relevant data results are shown in Figure 6.
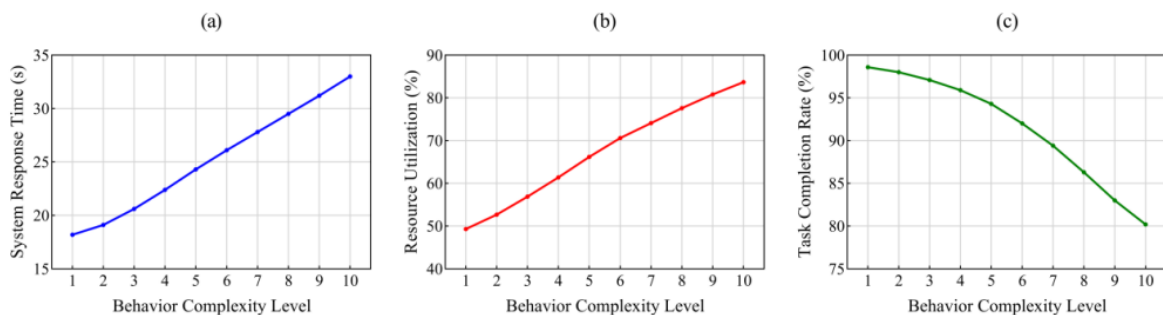


Figure 6: Analysis of the stability performance of the media teaching system at different levels of behavior complexity
Figure 6 (a): Changes in system response time
Figure 6 (b): Changes in resource usage
Figure 6 (c): Changes in task completion rate

As shown in Figure 6, system response time rises from 18.2 s at behavior complexity level 1 to 33.0 s at level 10, primarily due to the increased density of nodes and edges, which raises the computational cost of graph embeddings, and the higher concurrency of behavior paths, which intensifies scheduling pressure. Resource utilization grows from 49.3% at level 1 to 83.7% at level 10, reflecting those complex inputs trigger more concurrent requests in the multi-head attention mechanism and strategy generation module, thereby increasing thread utilization. Meanwhile, the task completion rate decreases from 98.6% to 80.2%, as higher complexity introduces greater ambiguity and interference in behavior paths, reducing the accuracy and timeliness of strategy recommendations.

Overall, these trends indicate that the system maintains strong adaptability across varying levels of behavioral complexity, but under extremely high input loads, computational bottlenecks and strategy deviations remain the primary constraints on stability.

## 4.6 Teaching content adaptability and task completion quality

This paper designs six representative media-based teaching tasks and conducts a comparative analysis of content recommendation and behavioral response effects for each. Task 1, basic knowledge explanation, emphasizes linear knowledge delivery with a clear structure but a single learning path. Task 2, multimedia case analysis, focuses on multimodal information

integration and reasoning, characterized by frequent content shifts. Task 3, interactive answering, requires learners to provide high-frequency, real-time feedback, resulting in high behavioral density and rapid interaction rhythms. Task 4, video demonstration learning, relies on visual information absorption, where behavior paths are primarily passive but demand sustained attention. Task 5, group collaboration, involves multi-user interaction with a complex and dynamic behavioral chain. Task 6, comprehensive skill assessment, integrates multiple knowledge points and operational steps, combining a relatively loose task structure with clear goal orientation.

Based on these task structures, two evaluation indicators—content adaptation and task completion quality—are employed to assess the system's overall collaborative performance across modules such as behavior modeling, strategy generation, and resource scheduling. The corresponding results are presented in Figure 7.
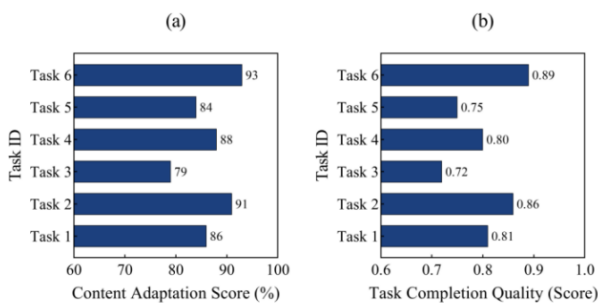


Figure 7: Comparison of content adaptation and task completion quality under different types of teaching tasks
Figure 7 (a): Teaching content adaptation
Figure 7 (b): Task completion quality analysis

As shown in Figure 7, the comprehensive skill assessment task achieved the best overall results, with a content adaptation rate of 93% and a task completion quality score of 0.89. This indicates that the system can more effectively capture learner behavior patterns in integrated tasks and provide accurate strategy matching. Multimedia case analysis (91%, 0.86) and video demonstration (88%, 0.80) also exhibit strong adaptability, though their completion quality differs due to varying demands on graphic recognition and short-term reasoning. Basic knowledge explanation (86%, 0.81) and group collaboration (84%, 0.75) present stable adaptability but lower completion quality, reflecting the constraints of linear knowledge delivery and the uncertainty of collaborative behaviors. By contrast, interactive question-answering tasks show the lowest performance (79%, 0.72), revealing the system's current limitations in modeling high-frequency, instant feedback behaviors. These results highlight that task characteristics significantly influence content recommendation accuracy and behavioral strategy adaptation, suggesting that further refinement of path modeling is needed for highly interactive and weakly structured tasks.

## 5 Discussion

The experimental results demonstrate that the proposed GAT+Contrastive Learning framework consistently outperforms sequential models (RNN, LSTM) and baseline GNNs (GCN, GAT) across multiple tasks. Its advantages stem from the combination of graph attention for fine-grained structural modeling, contrastive learning for embedding optimization, and representation fusion that reduces redundancy while enhancing intra-class consistency and inter-class separability. These mechanisms collectively improve recognition accuracy, structural retention, and semantic separation, thereby enabling more precise and adaptive teaching strategy generation. Nevertheless, performance in interactive question-answering and group collaboration tasks remains weaker, reflecting the difficulty of modeling high-frequency real-time interactions and multi-user dependencies, where dynamic feedback and irregular trajectories increase system complexity.

Despite the strong empirical results, this study has several limitations. First, the dataset includes only 186 students from a single institution, constraining generalizability to other educational settings. Second, the evaluation is restricted to media-rich higher education scenarios, leaving applicability in K–12 or vocational contexts unexplored. Third, the framework requires substantial computational resources, as multi-head GAT and contrastive learning introduce overhead that may hinder deployment in large-scale or resource-limited environments. Addressing these limitations in future work will require expanding datasets, applying domain adaptation techniques, and adopting model compression or graph pruning strategies to enhance scalability.

Beyond quantitative results, a preliminary qualitative survey with teachers and students revealed improved interpretability of behavioral feedback and stronger engagement through personalized recommendations, supporting the framework's practical value. To ensure reproducibility, this study explicitly defines its evaluation metrics: structure retention ratio (preservation of temporal–semantic links), path consistency (alignment between predicted and actual trajectories), semantic separation (inter- vs. intra-class embedding distance), and task completion quality (a composite score combining accuracy, timeliness, and performance outcomes). Together, these findings highlight the theoretical significance of graph-based contrastive modeling in advancing behavior analysis and its practical promise for adaptive intelligent education systems.

## 6 Conclusions

This study proposes a scalable framework for adaptive system design by modeling learner–media interactions as heterogeneous behavior graphs and combining Graph Attention Networks (GAT) with contrastive learning to optimize subgraph representations. Experimental results on a real-world multimedia course dataset demonstrate clear advantages: video engagement recognition reached 92.8%, disengagement signal

retention 69.1%, and clickstream prediction 89.8%, outperforming sequential and baseline graph models. These quantitative improvements confirm the framework's effectiveness in enhancing semantic discrimination, structural preservation, and adaptive strategy generation.

At the same time, the system shows limitations in high-frequency interactive tasks (e.g., real-time Q&A) and multi-user collaboration scenarios, where modeling dynamic, irregular behaviors introduce computational overhead and reduces efficiency. These shortcomings highlight the need for further optimization of graph scalability and behavioral path modeling.

Future research will extend this work through large-scale, cross-institutional experiments to verify generalizability, teacher-in-the-loop evaluations to assess the pedagogical value of generated strategies, and online deployment trials to examine scalability and real-time responsiveness. Additional directions include integrating richer multimodal signals such as eye-tracking and speech, and exploring graph compression and pruning to balance accuracy with computational efficiency. Together, these efforts aim to strengthen both the theoretical foundations of multimodal learner modeling and the practical feasibility of deploying intelligent education systems in diverse instructional contexts.

# References

[1] Multimedia Technologies in Education." International Journal of Computer Science and Network Security 22.6 (2022): 727–732.

[2] Tuhuteru, Laros, Desy Misnawati, Aslan Aslan, Zakiyatut Taufiqoh, and Imelda Imelda. "The Effectiveness of Multimedia-Based Learning to Accelerate Learning after the Pandemic at the Basic Education Level." Tafkir: Interdisciplinary Journal of Islamic Education 4.1 (2023): 128–141. https://doi.org/10.31538/tijie.v4i1.311

[3] Liu, Dongyang. "The Effects of Segmentation on Cognitive Load, Vocabulary Learning and Retention, and Reading Comprehension in a Multimedia Learning Environment." BMC Psychology 12.1 (2024): 4. https://doi.org/10.1186/s40359-023-01489-5

[4] Yorganci, Serpil. "The Interactive E-Book and Video Feedback in a Multimedia Learning Environment: Influence on Performance, Cognitive, and Motivational Outcomes." Journal of Computer Assisted Learning 38.4 (2022): 1005–1017. https://doi.org/10.1111/jcal.12658

[5] Zhao, Hui, and Lina Guo. "Design of Intelligent Computer-Aided Network Teaching System Based on Web." Computer-Aided Design and Applications 19 (2021): 12–23. https://doi.org/10.14733/cadaps.2022.s1.12-23

[6] Xiang, Meng, Fengfan Mao, and Li Xiao. "A Study on the Integration of Digital Language Teaching System into English Teaching." Journal of Computational Methods in Science and Engineering 23.2 (2023): 913–920. https://doi.org/10.3233/jcm-226491

[7] Wahyu, Riswan Tri, Nihta VF Liando, and Rinny Rorimpandey. "The Implementation of TikTok as Media Teaching to Improve Students' Speaking Ability." JoTELL 2.12 (2023): 1551–1564.

[8] Sumarsih, Parmin Parmin, and Lusi Rachmiazasi Masduki. "Effect of Problem-Based Learning Media Teaching Edugame to Improve Science Literacy and Collaboration on Solar System Material." Jurnal Elementaria Edukasia 7.4 (2024): 3354–3366. https://doi.org/10.31949/jee.v7i4.11619

[9] Xuan, Zhaozhen. "DRN-LSTM: A Deep Residual Network Based on Long Short-Term Memory Network for Students Behaviour Recognition in Education." Journal of Applied Science and Engineering 26.2 (2022): 245–252.

[10] Li, Guang, Fangfang Liu, Yuping Wang, Yongde Guo, Liang Xiao, and Linkai Zhu. "A Convolutional Neural Network (CNN)-Based Approach for the Recognition and Evaluation of Classroom Teaching Behavior." Scientific Programming 2021 (2021): 6336773. https://doi.org/10.1155/2021/6336773

[11] Song, Bo. "Multimodal Interactive Classroom Teaching Strategies Based on Social Network Analysis." International Journal of Networking and Virtual Organisations 30.1 (2024): 70–81. https://doi.org/10.1504/ijnvo.2024.136777

[12] Dewi, Ni Putu Dilia, Dadang Hermawan, and Dian Rahmani Putri. "Multimedia in Picture Series as Teaching Strategy in Encouraging English Learning Motivation to Bengkala Elementary Students." Yavana Bhasha 5.2 (2022): 165–175. https://doi.org/10.25078/yb.v5i2.1160

[13] Koutsoukos, Marios, Eleni Mavropoulou, and Serafeim Triantafyllou. "Teaching Specialty Courses Using Educational Applications: Five Teaching Scenarios from Technical Vocational Education." Journal of Education and Training Studies 12.4 (2024): 1–9. https://doi.org/10.11114/jets.v12i4.6926

[14] Rizou, Ourania, Aikaterini Klonari, and Dimitrios Kavroudakis. "Supporting Statistical Literacy with ICT-Based Teaching Scenario." International Journal of Education 9 (2021): 59–78. https://doi.org/10.5121/ije.2021.9405

[15] Kim, Seonghun, Sion Park, Jiwon Park, and Youjin Oh. "A Study on the Intention to Use the AI-Related Educational Content Recommendation System in the University Library." Journal of Korean Library and Information Science Society 53.1 (2022): 231–263.

[16] Li, Xiu, Aron Henriksson, Martin Duneld, Jalal Nouri, and Yongchao Wu. "Evaluating Embeddings from Pre-Trained Language Models and Knowledge Graphs for Educational Content Recommendation."

Future Internet 16.1 (2023): 12.
https://doi.org/10.3390/fi16010012

[17] Lee, Hsin-Yu, et al. "Exploring the Learning Process and Effectiveness of STEM Education via Learning Behavior Analysis and the Interactive-Constructive-Active-Passive Framework." Journal of Educational Computing Research 61.5 (2023): 951–976.
https://doi.org/10.1177/07356331221136888

[18] Zhang, Chenhong, Mingli Gao, Wenyu Song, and Chong Liu. "Formative Evaluation of College Students' Online English Learning Based on Learning Behavior Analysis." iJET 17.10 (2022): 240–255.
https://doi.org/10.3991/ijet.v17i10.31543

[19] Li, Xiufang. "Learning Behavior Analysis and Learning Effect Evaluation in Open Online Courses." Creative Education 13.4 (2022): 1337–1352.
https://doi.org/10.4236/ce.2022.134081

[20] Zhao, Fuzheng, Gwo-Jen Hwang, and Chengjiu Yin. "A Result Confirmation-Based Learning Behavior Analysis Framework…" Educational Technology & Society 24.1 (2021): 138–151.

[21] Jia, Yufeng, and Qi Zhao. "The Learning Behavior Analysis of Online Vocational Education Students and Learning Resource Recommendation Based on Big Data." iJET 17.20 (2022): 261–276.
https://doi.org/10.3991/ijet.v17i20.34521

[22] Fan, Ju, Yuanchun Jiang, Yezheng Liu, and Yonghang Zho. "Interpretable MOOC Recommendation: A Multi-Attention Network for Personalized Learning Behavior Analysis." Internet Research 32.2 (2022): 588–605.
https://doi.org/10.1108/intr-08-2020-0477

[23] Liu, Yang. "Blended Learning of Management Courses Based on Learning Behavior Analysis." iJET 16.9 (2021): 150–165.
https://doi.org/10.3991/ijet.v16i09.22741

[24] Cui, Qichao. "Multimedia Teaching for Applied Linguistic Smart Education System." International Journal of Human–Computer Interaction 39.1 (2023): 272–281.
https://doi.org/10.1080/10447318.2022.2122111

[25] Shen, Yu-bao, and Thippa Reddy Gadekallu. "Resource Search Method of Mobile Intelligent Education System Based on Distributed Hash Table." Mobile Networks and Applications 27.3 (2022): 1199–1208.
https://doi.org/10.1007/s11036-022-01940-8

[26] Hu, Yan, QiangQiang Li, and Shih-wei Hsu. "Interactive Visual Computer Vision Analysis Based on Artificial Intelligence Technology in Intelligent Education." Neural Computing and Applications 34.12 (2022): 9315–9333.
https://doi.org/10.1007/s00521-021-06285-z

[27] Zhang, Kun. "Design and Application of Intelligent Teaching System for Network and New Media Major Driven by Artificial Intelligence

Technology." International Journal of Embedded Systems 17.1–2 (2024): 150–159.
https://doi.org/10.1504/ijes.2024.143759

[28] Mei, Chang. "Application of e-Learning and New Media Teaching Platform Based on Human–Computer Interaction Technology in Civil and Commercial Law Courses." Entertainment Computing 50 (2024): 100677.
https://doi.org/10.1016/j.entcom.2024.100677

[29] Wang, Xingfei, Ke Zhang, Muyuan Niu, and Xiaofen Wang. "SemSI-GAT…" IEEE TKDE (2025).
https://doi.org/10.1109/tkde.2025.3528496

[30] Chen, Yong, Xiao-Zhu Xie, Wei Weng, and Yi-Fan He. "Multi-Order-Content-Based Adaptive Graph Attention Network for Graph Node Classification." Symmetry 15.5 (2023): 1036.
https://doi.org/10.3390/sym15051036

[31] Xu, Dongkuan, et al. "InfoGCL: Information-Aware Graph Contrastive Learning." NeurIPS 34 (2021): 30414–30425.

[32] Ju, Wei, et al. "Towards Graph Contrastive Learning: A Survey and Beyond." arXiv:2405.11868 (2024).

[33] Wu, Yucheng, et al. "Graph Contrastive Learning with Cohesive Subgraph Awareness." WWW 2024.
https://doi.org/10.1145/3589334.3645470

[34] Meng, En, and Yong Liu. "Graph Contrastive Learning with Graph Info-Min." CIKM 2023.
https://doi.org/10.1145/3583780.3615162

[35] Bharathi, S. T., et al. "Multimodal Transformer-Based Approach for Building Adaptive, Interactive, and Inclusive e-Learning Systems." ICSSAS 2025.
https://doi.org/10.1109/icssas66150.2025.11081234

[36] Chen, Yanlin, et al. "Image Sensor-Supported Multimodal Attention Modeling for Educational Intelligence." Sensors 25.18 (2025): 5640.
https://doi.org/10.3390/s25185640

[37] Assayed, Suha Khalil, Manar Alkhatib, and Khaled Shaalan. "A Transformer-Based Generative AI Model in Education: Fine-Tuning BERT for Domain-Specific Student Advising." In Breaking Barriers with Generative Intelligence, Springer (2024).
https://doi.org/10.1007/978-3-031-65996-6_14

[38] Zilong, XIA, and Fanchao Niu. "Deep Learning-Based Sentiment Analysis of COVID-19 Pfizer Vaccine Tweets Using Transformer and Bi-LSTM Architectures." Informatica 49.30 (2025).
https://doi.org/10.31449/inf.v49i30.8479

[39] Ji, Zheng, and Luhao Cao. "Multi-Modal Video Forgery Detection via Improved EfficientNet with Attention and Transformer Fusion." Informatica 49.30 (2025).
https://doi.org/10.31449/inf.v49i30.8831

[40] Han, Yulin, and Ruihao Zeng. "EEG-Based Emotion Recognition…" Informatica 49.30 (2025).
https://doi.org/10.31449/inf.v49i30.8817

[41] Liu, Hai Chuan, et al. "Novel Multimodal Contrast Learning Framework Using Zero-Shot Prediction for Abnormal Behavior Recognition." Applied Intelligence 55.2 (2025): 110.
https://doi.org/10.1007/s10489-024-05994-x

[42] Li, Yuwei, and Botao Lu. "Intelligent Educational Systems Based on Adaptive Learning Algorithms and Multimodal Behavior Modeling." PeerJ Computer Science 11 (2025): e3157.
https://doi.org/10.7717/peerj-cs.3157

[43] Liu, Kang, et al. "Multimodal Graph Contrastive Learning for Multimedia-Based Recommendation." IEEE Transactions on Multimedia 25 (2023): 9343–9355.
（无 DOI）

[44] Guo, Feipeng, et al. "Dual-View Multi-Modal Contrastive Learning for Graph-Based Recommender Systems." Computers and Electrical Engineering 116 (2024): 109213.
https://doi.org/10.1016/j.compeleceng.2024.109213

[45] Liang, Bin, et al. "Fusion and Discrimination: A Multimodal Graph Contrastive Learning Framework for Multimodal Sarcasm Detection." IEEE Transactions on Affective Computing 15.4 (2024): 1874–1888.
https://doi.org/10.1109/taffc.2024.3380375

[46] Wang, Shan. "Multimodal Learning Data Intelligent Personalized Learning Resource Recommendation System for Web-Based Classrooms." ICAICA 2024.
https://doi.org/10.1109/icaica63239.2024.10823003