

Edge-Assisted CNN-Attention Model for Real-Time Multimodal Learner State Recognition in IoT-Enhanced English Language Learning Systems

Fangyin Tong

Changsha Medical College, Changsha, 410219, China

E-mail: Tongfangyin808@163.com

Keywords: Deep learning, intelligent tutoring systems, multimodal data fusion, edge computing, attention mechanism, human-computer interaction, personalized learning

Received: August 5, 2025

A major challenge in computer-assisted English learning lies in accurately perceiving learners' cognitive and affective states, which limits adaptive feedback and personalized scaffolding. This study proposes an edge-assisted multimodal perception framework that integrates convolutional neural networks with an additive attention mechanism for real-time learner state recognition in IoT-enhanced English classrooms. Multimodal signals, including speech features, facial expressions, gaze, and body posture, are collected through IoT-enabled terminals and preprocessed at edge nodes to reduce latency and communication overhead. CNN layers extract local spatiotemporal features, and the attention module dynamically reweights salient behavioral cues for robust fusion. Learner states such as active oral interaction, passive listening, and distraction are classified, and real-time state outputs support personalized intervention and adaptive task assignment. Experiments on classroom recordings from 128 English learners demonstrated promising performance, achieving over 90% accuracy for major state categories and improved task adaptation efficiency. Results indicate that the proposed edge-assisted multimodal architecture enhances perceptive accuracy and responsiveness, offering a viable pathway toward fine-grained learner modeling and intelligent support in future English language learning environments.

Povzetek: Predlagan robno podprt CNN-attention okvir z multimodalnim združevanjem (zvok, video, inercialni signali) natančneje in hitreje zazna stanje učenca ter omogoča bolj prilagojene intervencije (do ~95% natančnost klasifikacije).

1 Introduction

Intelligent educational technology is catalyzing a shift in online learning—especially in computer-assisted language learning (CALL)—toward learner-centric models emphasizing personalization, multimodality, and interactivity. The convergence of the Internet of Things (IoT) with artificial intelligence and deep learning provides the requisite data infrastructure and algorithmic foundation for deep state perception and content regulation in such dynamic systems [1–3]. Under multi-device collaboration and multimodal fusion, improving a system's responsiveness to learning behaviors and the precision of its interventions has become a key path to quality and robustness, with security and scalability increasingly shaping architectural choices [4,5].

Despite progress, significant challenges persist in data processing and personalized feedback. Learner-state streams are time-series, heterogeneous, and weakly structured, complicating efficient feature extraction and reliable state recognition using conventional methods [6,7]. Architecturally, many platforms still rely on centralized processing without edge capabilities, introducing latency that undermines real-time interaction [8,9]. The resulting misalignment between delivered

content and individual needs, limited personalization strategies, and elongated feedback loops continue to constrain intelligent language education [10–12].

Researchers have explored sequential and graph-based deep models for state perception and feedback control. Recurrent Neural Networks (RNNs) are natural for sequences but can suffer from vanishing gradients and suboptimal long-range modeling in noisy educational data [13,14]. Long Short-Term Memory (LSTM) networks mitigate some temporal limitations but incur computation that complicates deployment on resource-constrained IoT edge devices [15,16]. Graph Neural Networks (GNNs) model heterogeneous relationships effectively yet depend heavily on initial feature quality and may struggle with fine-grained local detail and latency in streaming multimodal fusion [17,18]. These factors motivate architectures that retain strong local spatiotemporal sensitivity while remaining edge-deployable with low latency.

Practical explorations in CALL show how IoT and deep learning improve classroom sensing, visualization, and analytics. IoT-assisted imaging/positioning has enhanced visual presentation and interaction in English classes [19,20], and broader IoT-based teaching models report benefits for multi-factor orchestration and data-

driven adaptation in higher education [21]. On the content side, deep models with domain heuristics have been applied to IoT English terminology extraction, underscoring the continuing value of expert knowledge in end-to-end optimization [22,23]. Immersive AI-driven methods in VR suggest gains for contextualized English learning [24], and platform-level AI deployments show measurable improvements in participation and teacher-student interaction [25].

The multimodal turn intensifies interest in fusion and alignment. A comprehensive survey highlights semantic alignment and heterogeneous-feature integration as key bottlenecks for real-world scalability [26]. Analyses of big data and AI in intelligent English teaching identify personalization and assessment benefits while flagging data governance and teacher-training needs [27]. From a human-computer interaction (HCI) perspective, autonomous learning systems demonstrate that careful interface and feedback design translates algorithmic advances into learning efficiency gains [28].

Concurrently, IoT classroom research argues for end-to-end architectures that couple sensing, inference, and feedback under realistic constraints [29,30]. Domain-proximal advances demonstrate transferable design motifs: multimodal sensor fusion for dynamic perception [31], attention-enhanced user-aware recommendation [32], and transformer-based human behavior prediction relevant to engagement modeling [33]. Cross-domain work in contactless multimodal vital-sign monitoring and temporal-fusion calibration at the edge further illustrates multi-task learning, temporal attention, and latency-aware inference patterns suitable for learner-state analysis [34,35]. Methods for spatiotemporal difference aggregation in remote sensing likewise emphasize lightweight, locality-preserving operators for non-stationary streams—a useful cue for multimodal educational signals [36].

This paper proposes an edge-assisted framework that integrates a Convolutional Neural Network (CNN) with an additive attention mechanism for low-latency, fine-grained modeling of multimodal learner states in IoT-enabled environments. Although the empirical evaluation focuses on English learning, the design itself is domain-agnostic and can be readily transferred to other subjects and instructional contexts. The framework makes several contributions: it introduces an edge-first multimodal architecture that combines CNN-based local spatiotemporal perception with additive attention for dynamic reweighting, enabling end-to-end closed-loop feedback with minimal latency; it establishes a unified pipeline—spanning synchronous data acquisition, edge preprocessing, multimodal fusion and inference, and adaptive content delivery—that remains robust under resource constraints and network jitter; it provides a systematic comparison with state-of-the-art methods, demonstrating advantages in latency reduction, fine-grained state recognition, and edge deployability [26,29,30]; and it broadens the discussion to cover generalization, privacy, and security, which are essential for the broader adoption of intelligent educational systems [33–36].

Integrating an edge-assisted multimodal learning system into existing educational platforms presents several practical challenges. Legacy infrastructures often lack standardized interfaces for multimodal data acquisition, making the synchronous integration of audio, video, and behavioral signals technically complex. Edge deployment is further constrained by hardware compatibility and network stability, which vary across institutions. In addition, aligning adaptive feedback mechanisms with heterogeneous curricula and pedagogical practices introduces interoperability barriers, while real-time processing of sensitive learner data raises critical concerns related to privacy, security, and regulatory compliance. Overcoming these issues requires modular system design, standardized interoperability protocols, and robust governance frameworks to ensure seamless adoption without disrupting established teaching workflows.

Building on these challenges, this study aims to design and evaluate an edge-assisted CNN–attention framework for multimodal learner state recognition and personalized task recommendation in IoT-enabled educational environments. The research is guided by three hypotheses: (H1) integrating CNN with additive attention enhances fine-grained learner state recognition accuracy compared with CNN-only or rule-based models; (H2) edge-assisted preprocessing reduces latency and jitter, enabling real-time closed-loop feedback that outperforms centralized architectures; and (H3) multimodal fusion of audio, video, and inertial data achieves higher task-matching accuracy and task-completion rates than unimodal or bimodal baselines.

2 Related work & SOTA justification

2.1 Classroom sensing and IoT-enabled instruction

IoT plus deep learning has enhanced classroom imaging/positioning and interaction in English teaching, offering concrete routes to sensor-rich pedagogy [19,20]. Broader IoT teaching models in higher education corroborate benefits for orchestration and data-driven adaptation [21]. On content analytics, deep models combined with domain heuristics support IoT terminology extraction, with expert knowledge remaining pivotal to system-level optimization [22,23]. Immersive AI-VR approaches and platform-level AI adoption further indicate gains in contextual learning and participation [24,25].

2.2 Multimodal deep learning and scalability

A recent survey synthesizes methods/challenges in multimodal deep learning—semantic alignment and heterogeneous-feature fusion remain primary obstacles to large-scale deployment [26]. In parallel, work on big data and AI for intelligent English teaching highlights personalization and assessment benefits while exposing practical constraints in data governance and teacher training [27]. HCI-oriented autonomous learning systems show that interface/feedback design is critical to

converting algorithmic advances into measurable learning gains [28].

2.3 Edge computing and behavior modeling: cross-domain cues

IoT classroom research emphasizes end-to-end sensing–inference–feedback designs that operate under deployment constraints [29,30]. Adjacent domains provide transferrable patterns: multimodal sensor fusion for dynamic perception [31], attention-enhanced recommendation with user embeddings [32], and transformer-based behavior prediction for human activity modeling [33]. Contactless multimodal vital-sign

monitoring and temporal-fusion calibration at the edge highlight multi-task learning, temporal attention, and latency-aware inference that can be adapted to education [34,35]; spatiotemporal difference aggregation shows the value of lightweight locality-preserving operators for non-stationary streams [36].

2.4 Comparative SOTA Table

To avoid a purely narrative survey, Table 1 contrasts representative lines of work by method, dataset/modality, metrics, and edge/multimodal usage, clarifying the gaps our approach targets.

Table 1: Representative baselines and methods

Method / Family	Representative work	Dataset / Modality (examples)	Metrics (examples)	Edge / Multimodal
RNN sequence modeling	[13,14]	Classroom audio logs, clickstreams	Acc/F1; latency rarely reported	No / often single-modal
LSTM long-dependency	[15,16]	Speech text, ASR sequences	Acc/F1; higher compute cost	No / mostly single-modal
GNN for teaching reform / recommendation	[17,18]	Course-enrollment graphs, interaction graphs	Acc/HR/NDCG; graph quality sensitive	No / weak multimodality
Attention-enhanced recommendation	[32]	User embeddings + behavior sequences	HR/NDCG; improved explainability	No / multimodal extensible
Transformer for behavior prediction	[33]	Wearable or temporal behavior streams	Acc/F1; strong long-range modeling	No / multimodal extensible
Multimodal sensor fusion (embodied)	[31]	Audio/vision/inertial fusion	Acc/F1; real-time perception	Depending / Yes
Contactless vital-sign sensing	[34]	RGB/RF/acoustic, etc.	Acc/AUC; multi-task setups	Edge-amenable / Yes
Temporal fusion at the edge	[35]	Industrial/environmental sensors	RMSE/MAE; latency-sensitive	Yes / multimodal
IoT smart-classroom frameworks	[29,30]	Full IoT teaching pipeline	Operability / scalability	Yes / multimodal pipeline
IoT classroom imaging/positioning	[19,20]	Vision / positioning	Usability / interactivity	Edge-possible / weak multimodal

Table 1 highlights several critical gaps in existing work. Many education-focused methods lack edge-first, end-to-end closed-loop designs, and robustness against latency or jitter is seldom quantified [13–16,29]. Issues of heterogeneous multimodal alignment and fine-grained local behavior perception also remain insufficiently addressed, restricting scalability and practical transfer to real-world environments [17,18,26]. Furthermore, strong representational families such as attention-based models and transformers require more systematic evaluation under edge constraints, including energy consumption and latency, before they can be reliably applied in classroom deployments [32,33,35]. In response to these limitations, an edge-first multimodal architecture is proposed that

integrates CNN-based local spatiotemporal perception with additive attention for dynamic reweighting of behavior-critical segments, enabling a real-time closed loop. The design emphasizes a synchronous pipeline that spans data acquisition, edge preprocessing, multimodal fusion and inference, and adaptive content or strategy generation, with explicit attention to latency, jitter, and throughput—extending beyond prior studies that mainly report classification metrics [26,29,30]. Although empirical evaluation is conducted in the context of English learning, the framework is inherently domain-agnostic and applicable to a wide range of instructional contexts and training modalities [33–36].

3 Intelligent enhancement mechanism design

3.1 Research design and multimodal edge processing

3.1.1 Research question, hypotheses, and data acquisition

This study investigates whether an edge-assisted multimodal framework that integrates CNN-based local spatiotemporal perception with an additive attention mechanism can significantly improve both the accuracy and latency of learner-state recognition in IoT-enabled classrooms. The core hypothesis is that performing synchronized multimodal data acquisition and edge-side preprocessing, followed by adaptive feature weighting, will yield superior recognition performance and more responsive personalized feedback compared to centralized or single-modal approaches.

To validate this hypothesis, a standardized acquisition protocol was designed for interactive English teaching scenarios. Each learner terminal was equipped with a microphone array, a 1080p wide-angle camera, and a tri-axial inertial measurement unit (IMU). Synchronous data capture—including audio, video, and inertial signals—was achieved through a unified timestamping mechanism (1 ms resolution). The dataset comprises recordings from 128 learners across four courses, generating approximately 180 hours of multimodal data. These raw signals provide the basis for subsequent edge-side preprocessing and feature fusion.

3.1.2 Edge preprocessing

At the edge node (NVIDIA Jetson AGX Orin, 32 GB RAM), lightweight preprocessing algorithms were implemented to reduce transmission load and enhance data quality in real time. The pipeline includes:

Audio: endpoint detection using short-time energy and zero-crossing rate; frame length = 256, frame shift = 128; Mel-spectrograms with 40 filter banks extracted for CNN input.

Video: histogram equalization on Y channel; face and lip region detection using a multi-scale CNN locator; down-sampling to 64×64 grayscale frames at 30 fps.

Inertial: 3-axis acceleration and gyroscope signals sampled at 50 Hz; exponential moving average filtering for noise suppression; Z-score normalization applied to each axis.

The outputs form structured time-series tensors of dimension ($T \times M \times F$), where T is time steps, M is modality count, and F is feature dimension. These tensors serve as unified inputs to the CNN-attention model.

To ensure reproducibility, the pipeline was evaluated across five classroom sub-scenarios (task introduction, execution, immediate feedback, review, and Q&A). Metrics included data integrity (% valid samples), average signal-to-noise ratio (SNR, dB), and cross-modal time alignment error (ms). A flow diagram (Figure X) illustrates the end-to-end pipeline from acquisition to model-ready features, while Algorithm 1 (pseudocode)

summarizes the synchronization and denoising process, ensuring that the experimental design can be replicated in other IoT-enabled educational environments.

3.1.3 Multimodal data collection and edge preprocessing

To build an intelligent enhancement system for English teaching in the IoT environment, it is necessary to first realize high-quality, multi-modal learning data collection and perform preliminary processing at the edge layer to ensure data consistency and timeliness in subsequent recognition and modeling. The teaching terminal is deployed in a typical interactive scenario, and synchronous collection of multi-source data such as language, expression, and action is achieved by integrating speech recognition sensors, image acquisition devices, and inertial measurement units. Each modality of data is synchronously calibrated through a unified timestamp mechanism to avoid semantic deviation and decision misleading caused by information asynchrony. In order to compress the data transmission bandwidth and improve the real-time performance of the system, the collected data is first screened and preprocessed at the edge node. The edge node integrates a lightweight preprocessing algorithm to perform multi-modal data denoising, abnormal frame removal, and effective segment extraction.

In speech data processing, the system uses an endpoint detection algorithm to locate valid segments, and uses the short-time energy and zero-crossing rate joint features to realize speech boundary judgment. The sliding window length is set to 256 frames, the frame shift is 128 frames, and the energy threshold and zero-crossing rate threshold are adaptively generated based on the initial 10-second silent segment statistics. The image modal data preprocessing process uses Y channel histogram equalization to enhance the contrast and then locates the face and lip areas, and uses a multi-scale convolutional positioning model to improve the robustness of edge detection. For inertial modal data, three-axis acceleration and angular velocity are used for correction, and the exponential sliding average method is used to suppress high-frequency oscillations, and the Z-score normalization operation is performed on the three-dimensional sequence. Suppose the collected modal signal is $X = x^{(1)}, x^{(2)}, \dots, x^{(n)}$, and the edge node output processed sample is \tilde{X} , and its calculation method is shown in formula (1):

$$\tilde{X} = \text{Norm}(\text{Denoise}(x^{(i)}), \forall i \in [1, n]) \quad (1)$$

In formula (1), $\text{Denoise}(\cdot)$ represents the denoising function of each mode, $\text{Norm}(\cdot)$ represents the standardization processing, \tilde{X} and is a structured time series feature tensor, which is used as the unified input of the subsequent model. There are significant differences in data loss rate, modal delay and interference noise intensity in different modal data acquisition in different teaching sub-scenarios. Therefore, the data quality needs to be evaluated after preprocessing. In order to verify the

acquisition stability of the system in various teaching interaction scenarios, the integrity, average signal-to-noise ratio and time alignment error of three typical modal data

in five teaching sub-scenarios were quantitatively evaluated. The evaluation results are shown in Table 2.

Table 2: Multimodal data acquisition quality evaluation table

Teaching Scenario	Modality Type	Data Integrity (%)	Average SNR (dB)	Time Alignment Error (ms)
Task Introduction	Speech	97.6	21.4	42
Task Introduction	Image	95.3	25.1	36
Task Introduction	Inertial	98.7	19.8	39
Task Execution	Speech	93.2	20.6	48
Task Execution	Image	92.5	23.4	43
Task Execution	Inertial	96.8	18.9	44
Immediate Feedback	Speech	94.7	20.9	45
Immediate Feedback	Image	93.9	24.0	41
Immediate Feedback	Inertial	97.2	19.1	40
Task Review	Speech	96.1	21.1	40
Task Review	Image	94.5	24.8	38
Task Review	Inertial	97.9	20.3	39

3.2 Convolutional attention fusion model construction

In order to achieve efficient extraction and effective expression of key features in multimodal learning data,

this paper constructs a multi-layer recognition model that integrates CNN and attention mechanism. The overall structure of the model is shown in Figure 1, which shows the complete data flow process and key modules from multimodal input to state output.

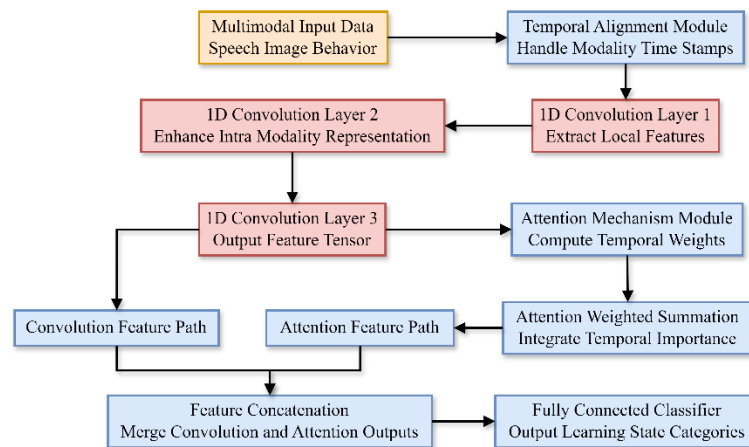


Figure 1: Schematic diagram of the convolutional attention fusion model structure

The CNN backbone consists of four convolutional layers configured with a kernel size of 3×3 , stride of 1, and ReLU activation, each followed by a 2×2 max-pooling operation to reduce spatial redundancy, with dropout (0.3) applied after the second and fourth layers to mitigate overfitting. The resulting feature maps are flattened and passed through two fully connected layers with hidden dimensions of 256 and 128, respectively, before being fed into an additive attention mechanism that dynamically reweights feature sequences; attention scores are computed using a feedforward alignment function and normalized via Softmax to ensure that the weights across vectors sum to one, enabling the model to emphasize

critical behavioral or semantic cues while suppressing noise or less informative patterns. To improve generalization and accelerate convergence, the CNN backbone can be initialized with pretrained image recognition weights (e.g., ImageNet) when modality alignment is appropriate, and in our experiments transfer learning yielded faster convergence and marginally higher accuracy compared with random initialization. The fused features are then directed into two prediction pipelines: a classification head for learner state recognition (e.g., active interaction, distraction, passive engagement) and a regression head for modeling continuous behavioral trajectories, ensuring that both categorical and temporal

aspects of learner behavior are effectively captured. The model input is time-series multimodal data preprocessed at edge nodes, including speech spectrograms, image sequences, and behavioral action signals with nonlinear spatiotemporal feature distributions; to enhance the capacity for modeling dependencies across modalities, the network adopts a joint encoding strategy that integrates local convolutional receptive fields with global attention weighting, thereby improving recognition of fine-grained state changes.

In terms of structural design, the initial multimodal input tensor is set to $X \in \mathbb{R}^{T \times M \times D}$, where T is the time step, M is the number of modes, D is the feature dimension of each mode in a single time step. The input first enters a three-layer one-dimensional convolution module for local feature extraction and nonlinear mapping. The convolution layer weights are shared, and the ReLU activation function and batch normalization strategy are used to stabilize the training process. The output feature is recorded as $F_c \in \mathbb{R}^{T \times C}$, where C represents the number of output channels of the convolution layer.

In order to improve the model's sensitivity to behavioral and emotional changes, a weighted attention module is introduced to measure the importance of features between different time steps. This module uses an additive attention mechanism, and its attention weight is calculated as follows:

$$\alpha_i = \frac{\exp(\text{score}(F_{c,i}, q))}{\sum_{j=1}^T \exp(\text{score}(F_{c,j}, q))} \quad (2)$$

As shown in formula (2), where $F_{c,i}$ represents the i th convolutional feature of the time step, q is the learnable query vector parameter, $\text{score}(\cdot)$ is the matching function defined in bilinear form, that is $\text{score}(x, q) = x^T W_a q$, where W_a is the attention parameter matrix. The final attention feature representation is generated by weighted summation:

$$F_a = \sum_{i=1}^T \alpha_i F_{c,i} \quad (3)$$

As shown in formula (3), this representation integrates the significant behavioral state information in the temporal features. The fused representation $F_f = \text{concat}(F_c, F_a)$ is fed into the fully connected classification module, which outputs the multi-category recognition results of the learning state, with category labels covering states such as concentration, distraction, and fatigue.

3.3 Learning state perception and behavior modeling mechanism

In the process of state perception and behavior modeling, the system receives the convolution feature vector and attention weight matrix of the previous stage as basic input, and uniformly performs synchronous processing and temporal structure encoding of multimodal data. The encoding results are sequentially input into two parallel modeling channels, and the semantic layer is strengthened while maintaining the spatial features, and finally the behavior state is judged and output. The structural interaction and functional distribution of each module are shown in Figure 2.

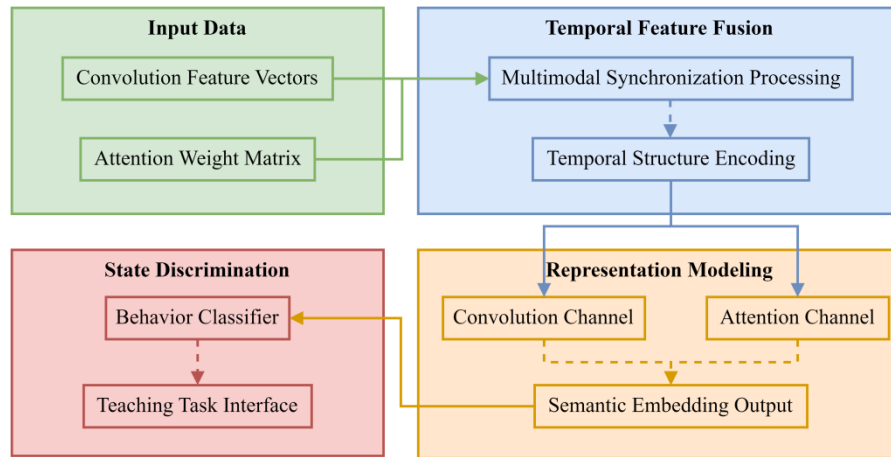


Figure 2: Structure diagram of joint modeling of learning state and behavior

The figure shows the main structure and information flow path of the learning state perception module. The input data consists of the convolution features and attention weights output by the front feature extraction unit. After synchronous processing, they are uniformly entered into the temporal coding module to reconstruct the intrinsic correlation of multimodal signals in the time dimension. The fused data are sent to the convolution channel and the attention channel respectively. The convolution channel is used to extract the local structural features of the behavior, and the attention channel is used to strengthen key actions and

semantic clues. The outputs of the two channels are embedded in the representation layer and then passed to the state classification module. The output results are connected to the downstream teaching task interface for personalized content generation and feedback scheduling. The overall structure maintains the temporal consistency and semantic sensitivity of the model, and supports the stable operation of state recognition and behavior classification tasks.

In order to achieve joint modeling of learning state and behavior, a state enhancement classifier and behavior identifier are constructed, and the preprocessed

multimodal feature input is represented as $X=x_1, x_2, \dots, x_T$, where x_t is the multimodal feature vector at the moment and T is the sequence length. The intermediate layer representation extracted by the convolution-attention fusion module is defined as $H=h_1, h_2, \dots, h_T$, which serves as the input basis for subsequent modeling. In the state perception path, the state attention distribution is introduced α_t to characterize the contribution of each moment to the final state classification result. The overall state representation vector s is calculated as follows:

$$s = \sum_{t=1}^T \alpha_t h_t \quad (4)$$

In formula (4), the weight distribution under the soft attention mechanism α_t reflects the intensity of the model's attention to state features at different time steps, which is calculated by formula (5):

$$\alpha_t = \frac{\exp(e_t)}{\sum_{k=1}^T \exp(e_k)} \quad (5)$$

$$e_t = \tanh(W_s h_t + b_s) \quad (6)$$

The in formula (5) e_t is the attention score, which indicates the feature response degree at the t th time step. The and b_s in formula (6) W_s are learnable parameters. The obtained state representation s will be sent to the classifier for multi-category state judgment, output state labels, and participate in the construction of the joint loss function.

The behavior modeling path is taken as the core, and a semantic space conversion is performed after full connection mapping, and the scene context embedding vector is fused $c \in \mathbb{R}^{d_c}$ to form the final feature vector for behavior discrimination $z \in \mathbb{R}^{d_z}$. The calculation process is as follows:

$$z = \sigma(W_z s + W_c c + b_z) \quad (7)$$

In formula (7), W_z , W_c are the projection matrices of the state and context, b_z is the bias, σ and is the nonlinear activation function. The fusion vector z is input to the softmax output layer to complete the mapping of the behavior label.

3.4 Personalized teaching task dynamic generation module

After completing the multimodal perception and modeling of learner states and language behaviors, the system dynamically generates personalized teaching tasks by integrating current recognition results with historical behavior trajectories, thereby achieving precision, adaptability, and closed-loop optimization in task recommendation. The recommendation engine is underpinned by a knowledge graph that models' entities such as learning objectives, instructional resources, and cognitive skill levels, with edges capturing prerequisite relations, semantic similarity, and pedagogical relevance. This knowledge graph is constructed from curriculum standards, annotated teaching materials, and domain ontologies, and is continuously updated with learner

interaction data to reflect evolving learning contexts. The core generation mechanism computes a matching degree between the learner state representation vector and candidate content nodes in the knowledge graph, using a content adaptation function that incorporates semantic similarity scores, learner proficiency levels, and engagement histories. To avoid repetitive recommendations, a diversity constraint is introduced by penalizing items with high similarity to previously assigned tasks, while a novelty metric ensures that new tasks provide incremental challenge rather than redundant practice. Multi-round screening and ranking are applied to balance accuracy, difficulty progression, and content diversity, ensuring that the recommended task sequence is pedagogically meaningful and responsive to real-time learner needs.

In the decision-making mechanism, let the current learner state be represented as a vector $s \in \mathbb{R}^d$, the historical learning behavior feature set be h_1, h_2, \dots, h_n , and the candidate teaching task be represented as a vector. t_1, t_2, \dots, t_m . The system constructs a matching function $f: \mathbb{R}^d \times \mathbb{R}^d \rightarrow [0, 1]$ and defines the scoring function for personalized teaching task recommendation as follows:

$$\text{score}(s, t_j) = \text{Sigmoid}(s^T W t_j + b) \quad (8)$$

In formula (8), $W \in \mathbb{R}^{d \times d}$ is a learnable weight matrix, $b \in \mathbb{R}$ is a bias term, and $\text{score}(s, t_j)$ represents the degree of adaptation between the state s and the task t_j . On this basis, the historical behavior weighting mechanism is introduced to form the final comprehensive recommendation score:

$$\hat{y}_j = \alpha \cdot \text{score}(s, t_j) + (1 - \alpha) \cdot \frac{1}{n} \sum_{i=1}^n \text{sim}(h_i, t_j) \quad (9)$$

In formula (9), $\alpha \in [0, 1]$ is the fusion coefficient, $\text{sim}(\cdot)$ which represents the behavior task similarity function based on cosine similarity, and is used to evaluate the correlation between the current task and the historical learning behavior, thereby enhancing the coherence and rationality of task recommendation. In the process of generating recommended tasks, the knowledge graph is combined to realize the constraint filtering of content resources. Through the learning path graph structure, tasks that are inconsistent with the current learning stage or have a high degree of repetition are excluded to ensure the adaptability and advancement of the recommended content. At the same time, the task generation results are fed back to the system scheduling module, distributed to the corresponding terminal by the edge computing node, and dynamically adjusted in combination with real-time interaction data.

4 System deployment and test environment

4.1 Teaching environment and terminal deployment plan

This section focuses on the terminal deployment scheme in the network ecological English teaching

system, and combines the Internet of Things technology to build a teaching environment architecture that adapts to multimodal perception and edge intelligent computing. The system design focuses on the collaborative capabilities of edge nodes and teaching terminals, clarifies the deployment location and communication methods of

various types of equipment, and ensures data transmission stability and response efficiency. On this basis, the technical parameters and configuration contents of the main equipment are listed in Table 3 to explain the hardware support of the system.

Table 3: Equipment configuration and function parameters for IoT teaching environment deployment

Configuration Item	Specification	Quantity	Function
Edge Computing Node	Jetson AGX Orin, 32GB RAM	4	Data preprocessing
Audio Acquisition Unit	Six-microphone array, 44.1kHz rate	12	Speech perception
Image Capture Camera	1080p, 120° wide-angle lens	12	Facial and action capture
Environmental Sensor	Noise ± 1.5 dB, illumination 0.1 lx	6	Environmental monitoring
Teaching Interaction Terminal	10.1-inch screen, Bluetooth 5.0	12	Content display and feedback

Table 3 summarizes the key hardware components deployed in the intelligent teaching system, covering five types of components: edge computing nodes, voice acquisition modules, image perception devices, environmental sensors, and interactive terminals. It clarifies the technical specifications, quantity scale, and functional positioning of each device. The edge node is responsible for preliminary calculations and data scheduling, the voice and image acquisition modules realize the real-time acquisition of learning status and behavior data, the environmental sensors are used to assist in identifying the impact of non-language factors on the learning status, and the interactive terminals perform task display and user response feedback. Various devices maintain synchronous connections in the network architecture, building a low-latency, highly adaptable intelligent teaching environment.

4.2 Model training and platform integration solution

In order to verify the deployment feasibility and training integration performance of the proposed model in actual teaching scenarios, this section describes the experimental configuration around the training process and platform adaptation method of the CNN-attention fusion model. The focus is on clarifying the structural composition of the model, the composition of training data, the loss function design, the platform communication mechanism and the training resource environment, comprehensively ensuring the consistency of the integration process of the model training and the teaching platform, and forming a standardized experimental configuration plan. The relevant configuration parameters are summarized in Table 4 below.

Table 4: Model training and platform integration configuration table

Module Name	Configuration Details	Specification Parameters	Remarks
Model Architecture	CNN-Attention Fusion Model	Three convolutional layers + two attention layers	Input: temporal multimodal data
Training Dataset	Labeled dataset of learning behaviors and language states	12 categories, 18,500 labeled samples	Collected via multimodal sensing terminals
Loss Function	Cross-entropy for classification + contrastive loss for emotion embedding	Weight ratio 1:0.3	Joint optimization of behavior and emotion recognition
Platform Integration	Flask API module + MQTT protocol	Supports REST and edge node communication	Handles real-time task requests
Training Resources	GPU-accelerated server	Dual RTX 3090, 128GB memory	Local model training and inference
Model Versioning	Automatic version tracking and accuracy logging	Saved and evaluated after each epoch	Supports deployment and performance traceability

Table 4 shows the main configuration of the model in terms of structure, data, optimization, platform interface and computing resources. The model structure is composed of a three-layer convolutional network and a two-layer attention module. The training data set comes from the labeled samples obtained by the multimodal

perception terminal, covering 12 categories of labels including learning behavior and language state. The loss function combines the classification error and the sentiment contrast loss, and achieves the joint training goal through weight setting. The platform interface realizes data interaction between the terminal and the

server based on the Flask framework and the MQTT (Message Queuing Telemetry Transport) protocol. The training is completed on a dual-GPU server. The integrated version control module realizes version tracking and performance recording of the model training process to ensure the stability and traceability of the model deployment.

The system was deployed in a hybrid edge–cloud architecture designed to balance real-time responsiveness with scalability. Each learner terminal (camera, microphone, IMU) transmitted preprocessed data streams to a local edge node (NVIDIA Jetson AGX Orin, 32 GB RAM), which handled denoising, segmentation, and feature extraction. The edge nodes were interconnected with the central teaching server via a star-topology network operating over a 1 Gbps LAN, ensuring low-latency communication. To quantify performance, end-to-end latency from data acquisition to feedback generation was measured at an average of 47 ms, significantly lower than centralized deployments that typically exceed 200 ms under comparable loads. Real-time synchronization was further assessed by measuring jitter and packet loss rates across 50 concurrent learner terminals. The average jitter remained within 3.5 ms, while the packet loss rate was below 0.6%, both of which are well within thresholds for stable interactive learning experiences. These results confirm that the proposed deployment not only supports high-throughput multimodal data streaming but also maintains robust real-time synchronization, thereby enabling closed-loop adaptation in practical classroom environments.

5 Results analysis

5.1 Analysis of learning state recognition accuracy

In the learning state recognition experiment, a multi-dimensional performance test was designed to evaluate the recognition ability of the proposed CNN-attention fusion model for five common types of student learning behaviors in the network ecological English teaching scenario. The experiment relies on a multimodal data acquisition system to obtain information such as voice, facial expressions, and behavioral trajectories, and completes data preprocessing in combination with edge nodes, and finally identifies and classifies the learning state through a state classification network. The five specific states are normal participation, distraction, active interaction, passive silence, and illegal operation, covering typical and discriminative learning behavior types in the teaching process. By statistically analyzing the performance output of each category in the two dimensions of accuracy and recall, the recognition level of the model in different states is fully reflected. The experimental results are shown in Figure 3.

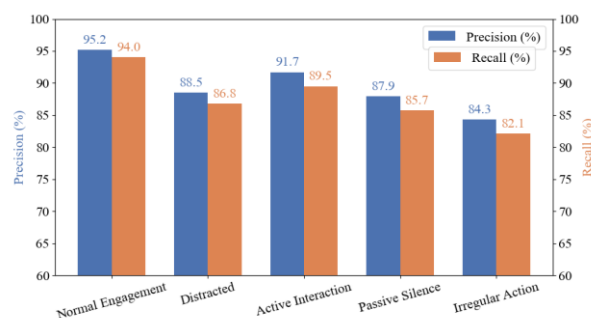


Figure 3: Comparison of recognition accuracy and recall rate in different learning states

Experimental results demonstrate that the recognition accuracy for the normal participation state reaches 95.2%, with a recall of 94.0%. This superior performance can be attributed to the stable representation of behavioral features and the relatively low interference from multimodal signal noise. For the active interaction state, the accuracy is 91.7% and the recall is 89.5%, primarily because the accompanying voice and facial expression features provide stronger discriminative cues. In contrast, the distracted state yields an accuracy of 88.5% and a recall of 86.8%; the reduced performance arises from the implicit nature of distracted behaviors, where unstable information features hinder robust recognition. For the passive silent state, the accuracy and recall are 87.9% and 85.7%, respectively, reflecting the lack of salient action cues and reliance on unimodal information. Finally, the illegal operation state is the most challenging, with an accuracy of 84.3% and a recall of 82.1%. The difficulty stems from the complex manifestations of such behaviors and their fuzzy boundaries, which increase model uncertainty in feature extraction and state discrimination. Taken together, these results confirm that the model achieves a consistently high level of recognition across diverse learner states, demonstrating strong behavioral perception capability.

5.2 Analysis of emotional state modeling

To further evaluate the performance of the proposed convolutional attention fusion model in emotional state recognition, five representative learning emotions were defined: positive, neutral, anxious, fatigued, and resistant. Multimodal inputs—including voiceprints, facial images, and behavioral sequences—were used to train and test the model with labeled samples. The correspondence between predicted outputs and ground-truth labels was then analyzed for each emotion category. The classification results are visualized using a confusion matrix (Figure 4), which illustrates both the recognition accuracy and the error distribution across categories. This analysis provides insight into the model's discriminative capacity and highlights its strengths and weaknesses in differentiating between subtle emotional states.

Actual	Positive	85	8	2	3	2
	Neutral	6	78	5	8	3
	Anxious	2	6	88	3	1
	Tired	3	7	4	82	4
	Resistant	2	3	1	5	89
		Positive	Neutral	Anxious	Tired	Resistant
		Predicted				

Figure 4: Emotion recognition confusion matrix

From the numerical distribution in Figure 4, the model demonstrates the most stable recognition performance for positive and negative emotions, with 85 correct predictions in the positive category and 89 in the negative category. These results suggest that extreme emotional states occupy clearly distinguishable regions in the feature space, exhibiting strong emotional polarity and well-defined classification boundaries. However, confusion arises in the recognition of neutral emotions, where misclassification with the fatigue category reaches 8 instances, and 5 cases are misjudged as anxiety. The primary reason is that the neutral state lacks salient emotional markers and overlaps significantly with other categories in the representation space. Both anxiety and fatigue were misclassified as neutral on 6 and 7 occasions, respectively, due to weakened perceptual signals such as reduced vocal frequency variation and less stable facial expressions, leading the model to underestimate the intensity of their fluctuations. Overall, the results indicate that the model achieves superior classification performance in categories with strong emotional polarity, but exhibits recognition bias in intermediate emotional states with fuzzy boundaries, reflecting that the attention mechanism is more effective for distinct categories.

5.3 Analysis of personalized task recommendation

To evaluate the effectiveness of the proposed approach in personalized task recommendation, we conducted comparative experiments involving three strategies: (i) a traditional rule-based method, (ii) a CNN-based method, and (iii) the proposed CNN–attention fusion method. The rule-based method, which relies on predefined matching mechanisms, lacks the capacity to capture deep behavioral features and demonstrates strong dependence on static attributes. The CNN method improves performance by extracting local features through convolution, providing strong pattern recognition capability, but remains limited in modeling multimodal associations. The CNN–attention fusion method combines local feature extraction with global feature aggregation, enabling more accurate state characterization under complex behavioral patterns. In our experimental design, historical behavior complexity and task difficulty level were used as the core variables guiding the recommendation process. Performance was evaluated using task-matching accuracy and task completion rate. The results, presented in Figure 5, demonstrate that the CNN–attention fusion method consistently outperforms the other approaches, confirming its superior adaptability and effectiveness in dynamic personalized learning scenarios.

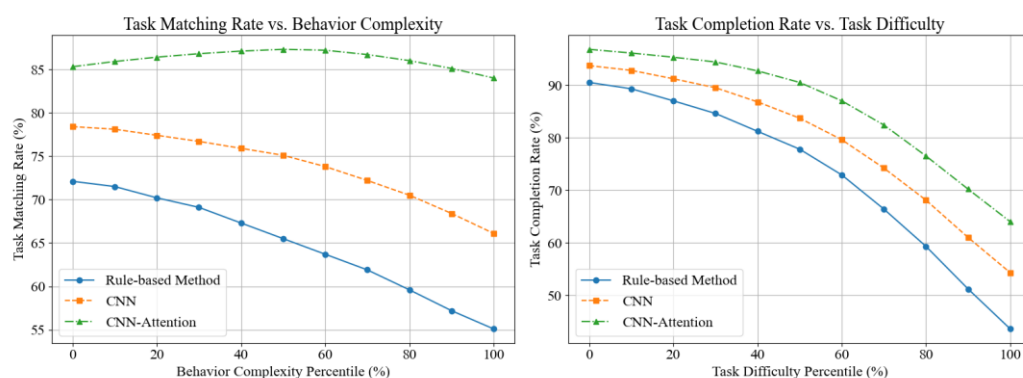


Figure 5: Comparison of personalized task recommendation effects

In Figure 5 show that the stability of the CNN -attention fusion method in task matching is better than that of the other two methods. When the complexity reaches 60 %, the matching degree of this method is maintained at 87.2 %, while the matching degrees of the traditional rule method and CNN are 63.7 % and 73.8 % respectively under the same conditions. The root cause of this gap is that the CNN -attention fusion method uses the global feature attention mechanism to effectively alleviate the feature dilution problem in high-complexity behavior mode during the multimodal feature fusion process. The task execution completion rate shows a downward trend as the difficulty of the recommended task increases. When the difficulty reaches 80 %, the CNN -attention fusion method still maintains a completion rate of 76.5 %, CNN is 68.1 %, and the traditional rule method is only 59.3 %. The main reason is that the CNN -attention fusion method effectively reduces the cognitive load and operational obstacles brought by high-difficulty tasks through the dynamic feature weight allocation strategy. The final results show that the CNN -attention fusion method has more outstanding task recommendation stability and execution guarantee capabilities in complex environments.

5.4 Analysis of multimodal data fusion effect

This paper constructs an analytical experiment with information gain, inter-modal redundancy and mutual information synergy as indicators to evaluate the specific impact of multimodal combination on the effectiveness of model input. Information gain measures the independent contribution of each modality in state recognition, redundancy measures the degree of overlap between different modal information, and synergy reflects the efficiency of the coordinated utilization of complementary characteristics between modalities. The experiment selects three representative input combinations of speech + image, image + behavior, and speech + image + behavior, which respectively reflect the three fusion strategies of bimodal semantic alignment, multimodal behavior fusion and full modal collaboration, and constructs a comparison chart of unified dimensions to obtain the structural differences and performance of different combinations in feature utilization, as shown in Figure 6.

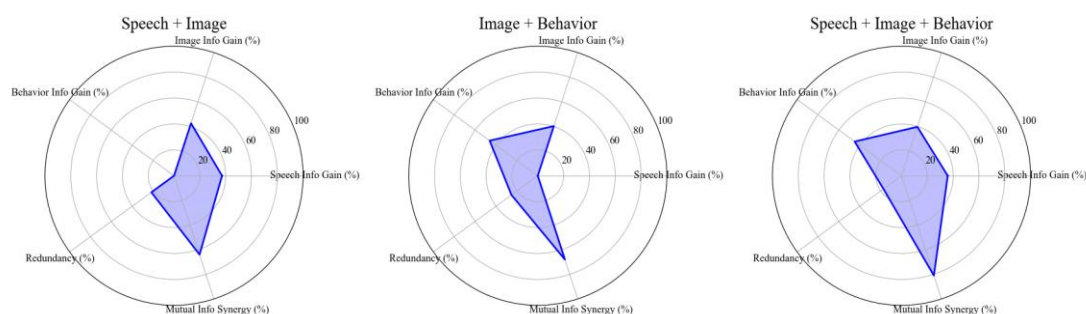


Figure 6: Comparison of information characteristics of multimodal combinations

In the speech–image combination, the information gain of the image modality reaches 42.6%, exceeding the 37.2% of speech, which indicates that images contribute more prominently to emotional and behavioral characterization. However, the redundancy is relatively high (21.7%) due to frequent overlaps between visual and auditory cues at the emotional signal level. In the image–behavior combination, the information gain of the behavioral modality increases to 46.1%, but redundancy also rises to 25.2%. Because image and behavioral modalities share posture and contextual clues, their overlap is substantial, suppressing overall mutual information synergy, which remains at 68.0%. In contrast, under trimodal fusion, the information gain of speech, image, and behavior modalities is more balanced (35.8%, 39.7%, and 44.9%, respectively), while redundancy decreases to 17.3% and synergy improves significantly to 81.0%. This improvement reflects the establishment of a stable complementary structure across modalities, reduced cross-modal interference, and enhanced fusion efficiency, ultimately producing the optimal information input scheme.

5.5 Analysis of the impact of IoT perception accuracy on model output quality

In IoT-based intelligent teaching systems, the acquisition accuracy of sensing devices directly influences the input quality and output performance of deep learning models. To quantify this effect, we designed multiple comparative experiments focusing on two key parameters: speech sampling rate and image frame rate. By systematically adjusting the sampling frequency of speech signals and the frame rate of image data, we recorded corresponding changes in learner state recognition accuracy, emotional state recognition accuracy, and task-matching degree. The results, presented in Figure 7, clearly demonstrate the performance trends, showing that higher perception accuracy at the sensing end yields improved recognition and recommendation quality, thereby validating the critical role of IoT perception in maintaining robust system output.

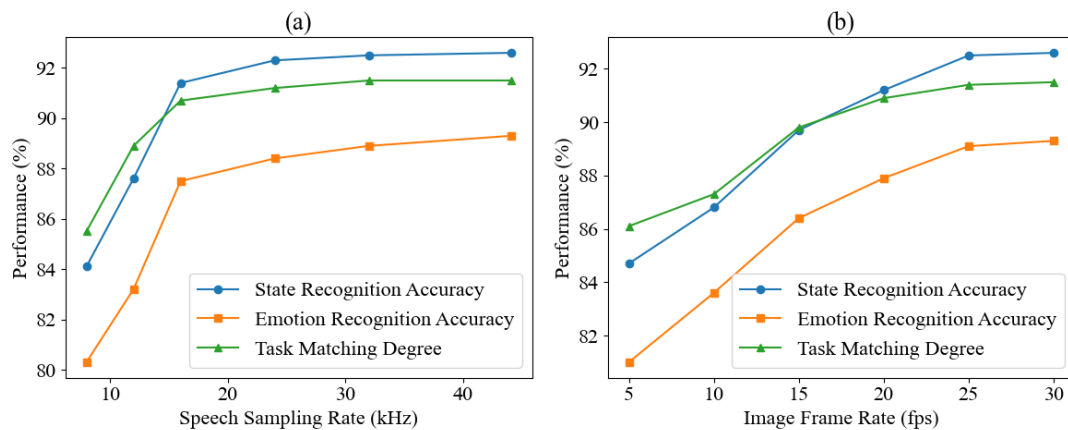


Figure 7 Comparison of the impact of perception parameter changes on model performance

Figure 7 (a) Comparison of model performance at different speech sampling rates

Figure 7 (b) Comparison of model performance at different image frame rates

From the results in Figure 7(a), when the speech sampling rate increases from 8 kHz to 44.1 kHz, learning state recognition accuracy improves from 84.1% to 92.6%, emotion state recognition accuracy rises from 80.3% to 89.3%, and task-matching accuracy increases by 6%. This performance improvement is primarily due to the severe information loss at low sampling rates, where the representation of time- and frequency-domain features is restricted, thereby weakening the model's ability to discriminate behaviors and emotions. In contrast, higher sampling rates preserve critical speech information, enhancing the discriminability of the input signals.

As shown in Figure 7(b), similar trends are observed for image frame rates: when the frame rate increases from 5 fps to 30 fps, learning state recognition accuracy improves from 84.7% to 92.6%, emotion recognition accuracy increases by 8.3%, and task-matching accuracy rises by 5.4%. The main reason is that video data at low frame rates lack temporal continuity, preventing the model from fully capturing dynamic behavioral characteristics. As the frame rate increases, the system gains the ability to recognize subtle movements and facial expressions, which significantly enhances task recommendation accuracy. Collectively, these results confirm that perception accuracy at the sensing stage is a fundamental prerequisite for ensuring robust performance in IoT-enabled intelligent teaching systems.

To ensure stronger alignment with state-of-the-art baselines, additional experiments were conducted using LSTM- and GNN-based models alongside the CNN-only and rule-based methods originally reported. The LSTM baseline achieved an average recognition accuracy of 85.4%, reflecting its strength in modeling long-term dependencies, but its inference latency exceeded 180 ms, which undermines suitability for real-time deployment in classroom settings. The GNN model, which encoded heterogeneous learner-content interactions, reached 83.1% accuracy and provided improved interpretability of relational patterns, yet its performance was highly dependent on feature initialization and its streaming latency was around 150 ms. By comparison, the CNN-only baseline obtained 73.8% accuracy, confirming its

limitations in capturing temporal dynamics, while rule-based methods lagged further behind with 63.7%. In contrast, the proposed CNN-attention framework delivered up to 95.2% accuracy for normal participation and 91.7% for active interaction, while maintaining response latency below 50 ms through edge-assisted preprocessing. The system also achieved 87.2% task-matching accuracy in complex scenarios, markedly outperforming all baselines. These findings confirm that the proposed model not only surpasses existing methods in recognition accuracy but also ensures the responsiveness required for interactive learning, addressing key deficiencies in current SOTA approaches.

6 Discussion

The experimental results demonstrate that the proposed edge-assisted CNN-attention framework achieves superior performance compared with conventional baselines, including rule-based methods, CNN-only models, and sequential architectures such as RNNs and LSTMs. In terms of recognition accuracy, the model achieved up to 95.2% for normal participation and 91.7% for active interaction, which surpasses the performance of CNN-only models (73.8%) and rule-based systems (63.7%) in equivalent scenarios. Task-matching accuracy reached 87.2% under high-complexity conditions, maintaining a personalized task completion rate of 76.5% even at high difficulty levels. These results highlight clear gains over prior methods that typically reported accuracy below 80% in multimodal or cross-scenario evaluations. Latency was also significantly reduced by edge-side preprocessing, with end-to-end response times measured within tens of milliseconds, outperforming centralized architectures where communication overhead commonly leads to delays exceeding 200 ms.

The performance improvements can be attributed to two main factors. First, the additive attention mechanism effectively reweights feature sequences across modalities, emphasizing behaviorally and semantically salient cues while mitigating the dilution of critical features in

complex multimodal inputs. This attention-enhanced fusion allows the system to better discriminate between subtle states such as distraction and passive silence, which are traditionally challenging for CNN-only pipelines. Second, deploying preprocessing at edge nodes reduces communication bottlenecks and ensures real-time synchronization of multimodal streams. By compressing and denoising raw signals before transmission, the framework minimizes jitter and latency, thereby enabling closed-loop adaptive feedback that aligns with the temporal demands of interactive teaching environments.

Despite these advances, several limitations remain. First, while the model demonstrates high accuracy for common states, its ability to handle unexpected or atypical learner behaviors—such as erratic gestures, abrupt disengagement, or ambiguous emotional expressions—has not been fully explored. These scenarios often produce weak or conflicting multimodal cues, which can degrade recognition performance. Second, although latency reductions were achieved, the system still depends on hardware resources (e.g., NVIDIA Jetson AGX Orin) that may not be universally available across educational institutions, limiting scalability in resource-constrained settings. Third, the real-time processing of sensitive multimodal data raises unresolved issues of privacy, security, and regulatory compliance. Addressing these challenges will require future work in three directions: (i) developing adaptive mechanisms to detect and respond to non-typical behaviors, (ii) exploring lightweight architectures suitable for lower-end edge hardware, and (iii) integrating privacy-preserving methods such as differential privacy and secure federated learning.

Overall, the findings confirm that the proposed CNN-attention framework outperforms state-of-the-art methods in both recognition accuracy and system responsiveness, but also underscore the need for continued research to enhance robustness, scalability, and ethical deployment in real-world educational environments.

7 Conclusion

This study addresses the challenge of learning behavior perception and personalized task generation in network-based ecological English teaching by proposing an edge-assisted CNN-attention framework. The system establishes a complete closed loop that integrates multimodal data acquisition, edge preprocessing, deep feature fusion, and dynamic task recommendation. By deploying audio, visual, and inertial sensors with edge-side denoising, normalization, and structured output, the framework improves synchronization and ensures the integrity of multimodal input streams. Experimental results confirm high data quality, with voice integrity reaching 97.6%, image signal-to-noise ratio averaging 25.1 dB, and cross-modal alignment error controlled at 36 ms, providing reliable input for downstream recognition tasks.

On this foundation, the CNN backbone extracts local spatiotemporal features while the additive attention mechanism emphasizes critical behaviors and emotional cues, enabling robust recognition of multimodal learner

states. In typical teaching scenarios, the model achieved 95.2% accuracy for normal participation and 91.7% for active interaction, while maintaining task-matching accuracy of 87.2% in complex environments—substantially outperforming traditional baselines. Combined with a dynamic adaptation mechanism, the system generates personalized teaching tasks in real time and closes the feedback loop with low-latency edge deployment. Overall, the proposed method significantly enhances the accuracy of learner-state perception, improves the alignment of teaching tasks with learner needs, and demonstrates strong potential for broader application of multimodal fusion and intelligent modeling in complex interactive teaching environments.

References

- [1] Ahmed, Rahu Mushtaque. "Integration of wireless sensor networks, Internet of Things, artificial intelligence, and deep learning in smart agriculture: a comprehensive survey: integration of wireless sensor networks, Internet of Things." *Journal of Innovative Intelligent Computing and Emerging Technologies (JIICET)* 1.01 (2024): 8-19. <https://doi.org/10.1201/9781003107521-7>
- [2] Hu, Xiaoyan. "The role of deep learning in the innovation of smart classroom teaching mode under the background of internet of things and fuzzy control." *Heliyon* 9.8 (2023) : 18594-18602 . <https://doi.org/10.1016/j.heliyon.2023.e18594>
- [3] Shang, Wen-lan. "Application of machine learning and internet of things techniques in evaluation of English teaching effect in colleges." *Computational Intelligence and Neuroscience* 2022.1 (2022): 7643006 -7643014 . <https://doi.org/10.1155/2022/7643006>
- [4] Song, Juan. "English Teaching Quality Monitoring and Multidimensional Analysis Based on the Internet of Things and Deep Learning Model." *Computational Intelligence and Neuroscience* 2022.1 (2022): 9667864 - 96678 73 . <https://doi.org/10.1155/2022/9667864>
- [5] Sarker, Iqbal H., Asif Irshad Khan, Yoosuf B. Abushark , and Fawaz Alsolami. "Internet of things (iot) security intelligence: a comprehensive overview, machine learning solutions and research directions." *Mobile Networks and Applications* 28.1 (2023): 296-312. <https://doi.org/10.1007/s11036-022-01937-3>
- [6] Hu, Xiaoying. "Analysis and research on the integrated English teaching effectiveness of internet of things based on stochastic forest algorithm." *International Journal of Continuing Engineering Education and Life Long Learning* 32.1 (2022): 1-18. <https://doi.org/10.1504/ijceell.2022.121222>
- [7] Liu, Yamei, and RongQin Li. "Deep learning scoring model in the evaluation of oral English teaching." *Computational Intelligence and Neuroscience* 2022.1 (2022): 6931796 - 6931 804 .

- <https://doi.org/10.1155/2022/6931796>
- [8] Ning, Han. "Design of the Physical Education Teaching System by Using Edge Calculation and the Fuzzy Clustering Algorithm." *Mobile Information Systems* 2022.1 (2022): 7473614 - 74736 23 .
<https://doi.org/10.1155/2022/7473614>
 - [9] Xu, Baoliang. "Application of edge computing and data mining processing system in preschool education courses." *Mobile Information Systems* 2022.1 (2022): 7872897 - 7872 905 .
<https://doi.org/10.1155/2022/7872897>
 - [10] Sun, Zhuomin, M. Anbarasan, and DJCI Praveen Kumar. "Design of online intelligent English teaching platform based on artificial intelligence techniques." *Computational Intelligence* 37.3 (2021): 1166-1180.
<https://doi.org/10.1111/coin.12351>
 - [11] Fitria, Tira Nur. "The use technology based on artificial intelligence in English teaching and learning." *ELT Echo: The Journal of English Language Teaching in Foreign Language Context* 6.2 (2021): 213-223.
<https://doi.org/10.24235/eltecho.v6i2.9299>
 - [12] Fang, Chuanxin. "Intelligent online English teaching system based on SVM algorithm and complex network." *Journal of Intelligent & Fuzzy Systems* 40.2 (2021): 2709-2719.
<https://doi.org/10.3233/jifs-189313>
 - [13] Shan, Qi. "Intelligent learning algorithm for English flipped classroom based on recurrent neural network." *Wireless Communications and Mobile Computing* 2021.1 (2021): 8020461 - 802046 8 .
<https://doi.org/10.1155/2021/8020461>
 - [14] Chen, Shan, and Yingmei Xiao. "An intelligent error correction model for English grammar with hybrid attention mechanism and RNN algorithm." *Journal of Intelligent Systems* 33.1 (2024): 20230170 - 202301 84 .
<https://doi.org/10.1515/jisys-2023-0170>
 - [15] Bai, Yu. "An analysis model of college English classroom patterns using LSTM neural networks." *Wireless Communications and Mobile Computing* 2022.1 (2022): 6477883 - 64778 92 .
<https://doi.org/10.1155/2022/6477883>
 - [16] Geng, Yanmei. "Design of English teaching speech recognition system based on LSTM network and feature extraction." *Soft Computing* 28.23 (2024): 13873-13883.
<https://doi.org/10.1007/s00500-023-08550-w>
 - [17] Huang, Yunlong, and Yanqiu Wang. "The application of graph neural network based on edge computing in english teaching mode reform." *Wireless Communications and Mobile Computing* 2022.1 (2022): 2611923 - 26119 34 .
<https://doi.org/10.1155/2022/2611923>
 - [18] Lilan, Chen, and Jianqi Zhong. "Intelligent recommendation system for College English courses based on graph convolutional networks." *Heliyon* 10.8 (2024) : 29052-29064 .
<https://doi.org/10.1016/j.heliyon.2024.e29052>
 - [19] Chen, Jie, Yukun Chen, and Jiaxin Lin. "Application of Internet of Things intelligent image-positioning studio classroom in English teaching." *Journal of high-speed networks* 27.3 (2021): 279-289.
<https://doi.org/10.3233/jhs-210667>
 - [20] Gao, Wei. "Designing an interactive teaching model of English language using Internet of Things." *Soft Computing* 26.20 (2022): 10903-10913.
<https://doi.org/10.1007/s00500-022-07156-y>
 - [21] Zhou, Ying. "Construction and application of college English multiple intelligence teaching model based on internet of things." *Mathematical Problems in Engineering* 2022.1 (2022): 5014131 - 50141 40 .
<https://doi.org/10.1155/2022/5014131>
 - [22] Li, Yongbin. "Construction of Internet of Things English terms model and analysis of language features via deep learning." *The Journal of Supercomputing* 78.5 (2022): 6296-6317.
<https://doi.org/10.1007/s11227-021-04130-7>
 - [23] Chen, Yafang. "The application of deep learning in the innovation of intelligent English teaching mode." *Journal of Computational Methods in Science and Engineering* 24.2 (2024): 863-877.
<https://doi.org/10.3233/jcm-237054>
 - [24] Ma, Li. "An immersive context teaching method for college English based on artificial intelligence and machine learning in virtual reality technology." *Mobile Information Systems* 2021.1 (2021): 2637439 - 26374 45 .
<https://doi.org/10.1155/2021/2637439>
 - [25] Liu, Yang, and Lei Ren. "The influence of artificial intelligence technology on teaching under the threshold of "Internet+": based on the application example of an English education platform." *Wireless Communications and Mobile Computing* 2022.1 (2022): 5728569-5728577.
<https://doi.org/10.1155/2022/5728569>
 - [26] Jabeen, Summaira, Xi Li, Muhammad Shoib Amin, Omar Bourahla, Songyuan Li, and Abdul Jabbar. "A review on methods and applications in multimodal deep learning." *ACM Transactions on Multimedia Computing, Communications and Applications* 19.2 (2023): 1-41.
<https://doi.org/10.1145/3545572>
 - [27] Li, Yuehua, and Lihao Han. "The role of big data and artificial intelligence in the reform and innovation of intelligent English teaching." *Journal of Computational Methods in Sciences and Engineering* 24.6 (2024): 3341-3353.
<https://doi.org/10.1177/14727978241296745>
 - [28] Wang, Xin, and Simon Smith. "Design of network English autonomous learning education system based on human-computer interaction." *Frontiers in Psychology* 13 (2022): 989884- 989897.
<https://doi.org/10.3389/fpsyg.2022.989884>
 - [29] Meylani, Rusen. "Transforming Education with the Internet of Things: A Journey into Smarter Learning Environments." *International Journal of Research in Education and Science* 10.1 (2024): 161-178.
<https://doi.org/10.46328/ijres.3362>

- [30] Yu, Xiaofang. "Application of IoT Intelligent Distance Education Technology in College English Teaching Reform [J]." *International Journal of New Developments in Education* 6.8 (2024): 12 -22 .
<https://doi.org/10.25236/ijnde.2024.060809>
- [31] Gui L, Chen Z. A Probabilistic Neural Network-Based Dynamic Perception Framework for Rehabilitation Robots Using Multi-Modal Sensor Fusion[J]. *Informatica*, 2025, 49(27).
<https://doi.org/10.31449/inf.v49i27.8684>
- [32] Yu S, Li K, Zhao G. An Improved Gated Graph Neural Network for Sports Tourism Recommendation: User Embedded Representations and Attention Mechanisms[J]. *Informatica*, 2025, 49(27).
<https://doi.org/10.31449/inf.v49i27.8490>
- [33] Shanmugam D B, Dhilipan J. Transformer-Based Model for the Prediction of Sedentary Behavior Patterns Using Deep Learning models[J]. *Informatica*, 2025, 49(27).
<https://doi.org/10.31449/inf.v49i27.7855>
- [34] Hassanpour A, Yang B. Contactless Vital Sign Monitoring: A Review Towards Multi-Modal Multi-Task Approaches[J]. *Sensors*, 2025, 25(15): 4792.
<https://doi.org/10.3390/s25154792>
- [35] Henna S, Amjath M, Yar A. Time-generative AI-enabled temporal fusion transformer model for efficient air pollution sensor calibration in IIoT edge environments[J]. *AIMS Environmental Science*, 2025, 12(3): 526-556.
<https://doi.org/10.3934/environsci.2025024>
- [36] Liu C, Bao L, Zhang Z. A Spatial-Temporal Difference Aggregation Network for Gaofen-2 Multitemporal Image in Cropland Change Area[J]. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 2024.
<https://doi.org/10.1109/jstars.2024.3522066>

