

Two-Way Classroom Interaction Analysis via a Coupled ConvNeXt–Multimodal Transformer for Fine-grained Behavior Recognition

Yuyan Huang¹, Mohammed Y. M. Mai^{2*}

¹Shandong Huayu University of Technology, DeZhou 253034, China

²Faculty of Education, Universiti Pendidikan Sultan Idris, Tanjong Malim, Perak Darul Ridzuan 35900, Malaysian

E-mail: Mohammed.mai@fpm.upsi.edu.com

*Corresponding author

Keywords: classroom behavior analysis, multimodal fusion, transformer architecture, convolutional neural networks, spatio-temporal feature interaction

Received: August 5, 2025

With the deepening of the digital transformation of education, intelligent analysis of classroom teaching behavior has become the key to improving teaching quality. Traditional methods are difficult to effectively integrate multi-source heterogeneous data in the classroom, and there are limitations in the joint modeling of spatiotemporal features. To this end, a bidirectional analysis framework coupling multimodal transformer and convolutional neural network (CNN) is proposed: ConvNeXt-T is used as the CNN backbone to extract the spatial features of teachers' body movements, students' postures and scene layouts, and the time dependence and cross-modal global correlation of teacher-student language interaction are captured with the help of multimodal transformers. The study uses 500 minutes of multimodal data from 10 real classrooms (4K camera 30 frames per second, total frames 900,000 frames) as the core dataset, annotates 7 types of behaviors such as teacher teaching, questioning, and student answering, and uses the PyTorch framework to train on NVIDIA GTX 4090 GPU, using AdamW as the optimizer, mixed loss function to process 8 batches of data, and the loss stabilizes at about 0.17 after 80 rounds of training. The results show that the accuracy of the multimodal fusion model is 90.2% in the behavior recognition task, which is significantly higher than that of the single-modal model. The spatio-temporal feature interaction module increases the detection rate of cross-modal correlation by 6.0%, and effectively identifies the linkage relationship between teachers' gesture pointing and students' responses. In the classification of teacher-student interaction, the F1 value of the model reached 88.4%, which was significantly higher than that of the benchmark model. In addition, the model has excellent generalization on public datasets, with an accuracy of 96.54% for NTU60-CV (cross-viewing angle), 98.30% for behavior recognition of UTD-MHAD, and an AUC value of 0.7478. This framework provides new ideas for solving fine-grained behavior analysis in educational scenarios and provides technical support for intelligent teaching evaluation.

Povzetek: Za analizo interakcij učitelj–učenec, kjer enomodalni modeli slabo povezujejo prostorske in časovne vzorce, je predlagan dvostranski okvir, ki združuje ConvNeXt za vizualne prostorske značilke in multimodalni Transformer za časovno ter medmodalno povezovanje. Model učinkovito prepozna fine vzorce vedenja in medsebojno vplivanje učiteljev in učencev.

1 Introduction

At the heart of the classroom environment is the continuous and dynamic multi-level interaction between teachers and students [1], which not only transmits explicit knowledge, but also includes non-verbal cues, emotions, and cognitive feedback, affecting teaching quality and learning effectiveness [2, 3]. However, traditional analysis methods are limited to a single modal or isolated perspective, making it difficult to capture teachers' teaching expressions and students' reactions together, resulting in the system

deconstruction of classroom interaction dynamics and unable to meet the needs of intelligent analysis tools [4, 5]. In real scenarios, teacher-student behavior is asynchronous and nonlinear in time and space, and traditional video frame-level recognition or isolated speech computing cannot integrate these heterogeneous features [6, 7].

Visual modality analysis relies on the local feature extraction capabilities of convolutional neural networks (CNNs) to effectively detect spatial information such as teachers' board writing and students' postures [8, 9], but the locality of CNNs limits their long-term dependence on cross-temporal and spatial dependencies (e.g., the evolution of

teacher questioning strategies and the change of student states). Transformer architectures are good at processing long-distance dependencies on sequence data through self-attention, and are suitable for speech prosody and text semantic analysis [10, 11], but they are computationally complex in high-resolution visual sequence processing and have weak spatial structure perception [12]. This demand drives the design of multimodal transformer-CNN coupling structures.

The core of the coupling mechanism is to construct multimodal collaborative representation: CNN serves as the underlying visual parser to capture the local spatiotemporal characteristics of teachers' and students' body movements, expressions, and environmental interactions. Transformer serves as a high-order integration hub to model the cross-modal timing dependence of visual information flow and audio and environmental signals [13]. This coupling is not a simple feature splicing, but an early two-way empowerment through attention-guided cross-modal feature calibration [14]—the teacher's key gesture visual features can modulate the synchronous speech semantic weights, and the student whispered sound features can also guide the visual system to focus on distraction areas [15]. The two-way analysis is embodied in dual modeling: the chain reaction of students' cognitive-behavioral triggers by teachers' instructions and the reverse regulation of teachers' teaching decisions by student feedback to form a closed-loop interactive flow [16].

This research is committed to building a new paradigm of classroom interaction analysis based on the deep coupling model, and its potential value significantly promotes the paradigm change of educational cognitive process analysis. The spatio-temporal alignment and semantic unified representation of multi-source heterogeneous data will deconstruct the black box between knowledge transfer and reception in the complex adaptive system of the classroom. This not only provides unprecedented micro-insight for refined evaluation of the effectiveness of teaching strategies and diagnosis of real learning obstacles, but its derivative intelligent tools can also empower educators to perceive the classroom situation in real time and dynamically optimise interactive strategies. This framework explores the in-depth dialogue between computational models and pedagogy, which is expected to open a new path for artificial intelligence to understand and optimise the transmission process of human core knowledge, and is a solid step towards a truly smart classroom.

2 Theoretical basis and principle technology

2.1 Attention mechanism

Self-Attention Mechanism has made remarkable achievements in the field of deep learning, especially in natural language processing and computer vision [17, 18]. It allows the model to automatically identify dependencies in different parts of the data and capture key information extensively. The self-attention mechanism was first adopted in the Transformer model and quickly became the core component of deep learning tasks, especially in dealing with long sequences and establishing global dependencies [19]. Network architectures based on self-attention mechanisms, such as Vision Transformer (ViT), surpass traditional CNN models in tasks such as image classification and motion recognition, demonstrating excellent performance.

The self-attention mechanism evaluates the similarity or correlation between input data elements and adjusts their weights [20]. Through the attention matrix, each element updates its expression through interactions with other elements.

For the input sequence $X \in \mathbb{R}^{T \times D}$, where T is the sequence length and D is the feature dimension, applying the trainable weight matrices W_Q , W_K , and W_V , according to formulas (1), (2), and (3), the query matrix Q , the key matrix K , and the value matrix V can be obtained by linearly transformation.

$$Q = XW_Q \quad (1)$$

$$K = XW_K \quad (2)$$

$$V = XW_V \quad (3)$$

After obtaining the linearly transformed Q , K , V , QK^T is used to calculate the similarity between features with a scaling factor \sqrt{d} to stabilize the training. The weight sum is ensured to be 1 by softmax normalization, forming a probability distribution. Finally, the similarity and V are weighted to obtain the Attention output Attention, as shown in Equation (4).

$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (4)$$

The core innovation of the Transformer model is the adoption of self-attention mechanism, which replaces cyclic and convolutional networks to efficiently capture long-distance dependencies [21]. As shown in Figure 1, the model consists of an encoder and a decoder, both of which are stacked by multiple layers of the same structure. The encoder transforms the input sequence into a hidden vector expressing its meaning, while the decoder uses the encoder output and its own historical information to step by step construct the target sequence [22].

Key computational units inside the encoder and decoder, such as multi-head self-attention, feedforward neural network, and residual connection, jointly promote the excellent performance of the model in information transfer and gradient optimization.

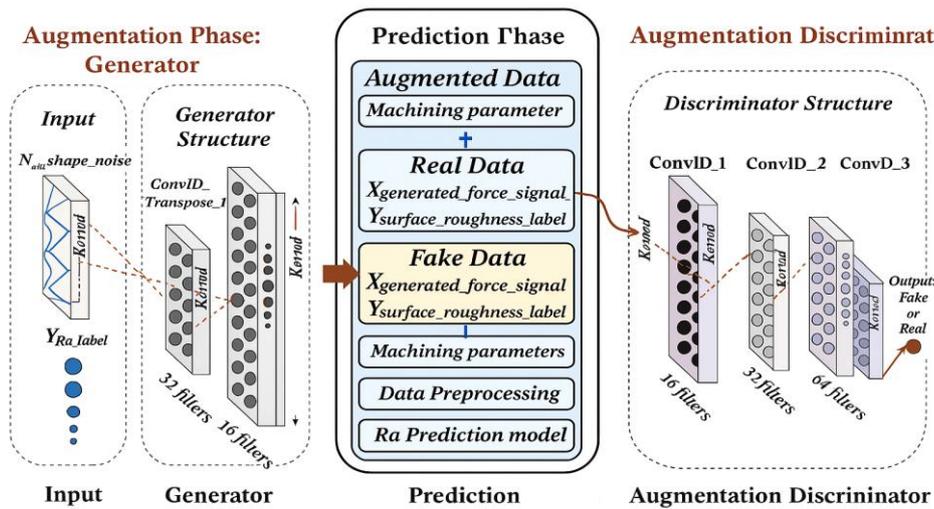


Figure 1: Transformer model structure

The core of Transformer is multi-head attention (MHA), which learns features of different subspaces via multiple independent attention heads to enhance model expression and generalization. Compared with single-head attention, it promotes feature subspace information exchange, enables more comprehensive context understanding, and breaks single attention mode limitations [23, 24].

In MHA, input data undergoes multiple linear transformations to generate queries, keys, and values (enabling multi-perspective data understanding). Each head independently calculates scores to generate weighted representations; finally, all heads' outputs are merged and mapped back to original dimensions via linear transformation, ensuring output-input shape consistency.

Transformer is a self-attention-based neural network that changes sequence modeling [25, 26]. It is applied in natural language processing and shows potential in computer vision—e.g., Vision Transformer (ViT) achieves excellent image feature extraction by slicing images into patches, converting to serialized Tokens, and using Transformer for global feature modeling. This method works well for large-scale data and becomes a key computer vision research direction.

2.2 Graph neural network (GNN) encoding

Classical deep learning models, such as convolutional neural networks and recurrent neural networks, mainly process

Euclidean spatial data such as images and sequences [27], but their effectiveness in processing graph structure data with complex relationships is limited. Therefore, graph neural networks (GNNs) are proposed for deep learning, which has significant effects on independent and homogeneous data processing.

GNN is designed for graph data, which can analyze node associations and attributes, and learn the overall information of the graph by optimizing node characteristics and integrating surrounding information to achieve effective modeling and prediction of graph data [28, 29].

The success of GNN stems from its unique messaging mechanism, which can effectively capture node attributes, neighbor node features, and local network structure information by summarizing neighbor node feature data to update the target node feature representation, thus performing well in graph mining tasks such as node classification, link prediction, and graph classification [30–35].

The current demand for processing graph structure data is increasing, and despite GNNs' achievements in modeling, they still face challenges such as privacy breaches and attacks. Research shows that GNNs are susceptible to privacy breaches, which can lead to personal information leakage. At the same time, it has deficiencies in fairness and explainability, which may amplify bias, exacerbate social prejudice, and may reduce fairness due to data pollution or be misinterpreted by attackers.

2.3 Method module architecture table

Table 1: Technical parameter table of neural network module

Module	Input Layer	Core Computation Layers	Output Layer	Key Technical Notes
UNeXt	Image (e.g., $3 \times H \times W$)	1. Encoder: Conv blocks (3×3) + MaxPool; 2. Bottleneck: Transformer encoder (self-attention); 3. Decoder: Upsample + Conv	Segmented feature map ($1 \times H \times W$)	Combines CNN local feature & Transformer global attention
ConvNeXt-T	Image/feature map ($C \times H \times W$)	1. Stem: 4×4 Conv (stride=4); 2. Stage ($3 \times$): ConvNeXt block (7×7 Conv, LayerNorm); 3. Head: AdaptiveAvgPool + Linear	Class logits $N \times \text{num}_c \text{ classes}$	Optimized CNN with large kernel & normalized activation
TransNeXt	Patch embeddings ($N \times D$)	1. PatchEmbed: 16×16 Conv (stride=16); 2. Transformer block ($12 \times$): Multi-Head Attention (MSA) + MLP; 3. Norm: LayerNorm	Global feature ($1 \times D$)	Transformer-based with enhanced patch interaction
SC Fusion Block	Dual-modal features ($C1 \times H \times W$, $C2 \times H \times W$)	1. Feature align: Resize (bilinear); 2. Fusion: Channel-wise concat + 1×1 Conv; 3. Refine: 3×3 Conv + ReLU	Fused feature ($C3 \times H \times W$)	Aligns & fuses spatial (e.g., video) & contextual (e.g., audio) features
SRU	Sequential features ($T \times D$)	1. Gate layer: Sigmoid (for feature selection); 2. Recurrent layer: Linear + tanh (temporal modeling); 3. Residual connection	Temporal feature ($T \times D$)	Lightweight recurrent unit for sequential behavior (e.g., teacher gestures)
CRU	Spatial features ($C \times H \times W$)	1. Context extract: Dilated Conv (3×3 , rates=2/4); 2. Fusion: Channel attention + spatial attention; 3. Output: Conv	Context-enhanced feature ($C \times H \times W$)	Captures long-range spatial context (e.g., student seating)
SCFB	Multi-scale features ($C \times H1 \times W1$, $C \times H2 \times W2$)	1. Scale align: Upsample/Downsample; 2. Fusion: Cross-scale attention + concat; 3. Refine: 3×3 Conv	Unified-scale feature ($C \times H \times W$)	Fuses multi-scale features (e.g., fine-grained facial expressions & coarse body posture)
CTHEFM	Hierarchical features ($C1 \times H \times W$, $C2 \times H \times W$)	1. Hierarchy align: 1×1 Conv (channel matching); 2. Fusion: Transformer cross-attention (between layers); 3. Norm: LayerNorm	Hierarchy-fused feature ($C \times H \times W$)	Bridges low-level (texture) & high-level (semantic) features
DEAM	Emotion-related features ($D \times T$)	1. Feature extract: MLP (for feature projection); 2. Attention: Emotion-aware self-attention; 3. Classify: Linear + Softmax	Emotion scores $T \times \text{num}_c \text{ motions}$	Detects dynamic emotions (e.g., student engagement, teacher mood)

3 Multi-modal coupling frame design

3.1 Design of multi-modal transformer and convolutional neural network coupling core module

Figure 2 presents the overall architecture. UNeXt, a high-resolution teacher-student behavioral semantic segmentation network based on UNet, comprises encoder, decoder, and

skip connection. The encoder adopts ConvNeXt-T; the lightweight decoder uses TransNeXt, executing CNN and Transformer block operations. The model reduces redundant info via SC fusion block and enables effective global-local info communication. The encoder has 4 stages: each stage halves the feature map, with output denoted as S_n . The decoder (4 stages, similar structure to encoder) improves resolution and reduces channels via upsampling, finally restoring the feature map to original size and generating prediction results.

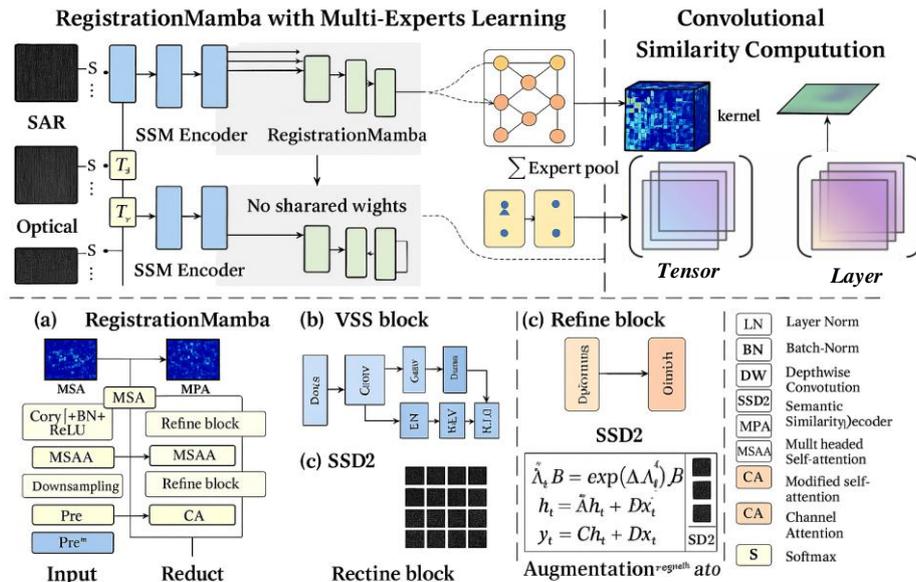


Figure 2: Multimodal transformer and convolutional neural network coupling model

In the semantic segmentation task of teacher-student behavior, CNN encoders such as ResNet series are crucial. But the introduction of ViT has changed the network architecture design. Although the global attention mechanism of ViT is computationally complex $O(N^2 \cdot d)$, its formula (5) is as follows:

$$O(N^2 \cdot d) \quad (5)$$

N represents the input sequence length, i.e. the number of image patches, and d is the feature dimension of each token. Vision Transformer (ViT) has difficulty in processing high-resolution teacher-student behavior data, while ConvNeXt inherits the efficiency of CNN, improves convolution operations, and reduces computational complexity. This is essential for large-scale dataset processing. Convolution operation is superior to self-attention mechanism in local information modeling, especially in small target object segmentation tasks. ConvNeXt improves processing speed while maintaining high accuracy.

ConvNeXt is based on ResNet-50, draws lessons from Transformer, and adopts deep separable convolution technology. It uses a 7×7 convolution kernel instead of a 3×3 convolution kernel to enlarge the receptive field. In order to enhance the model's adaptability and fluency to complex data while maintaining training stability, ConvNeXt introduces the GELU activation function to replace ReLU, forming its

core architecture. The calculation formula (6) of the GELU activation function is: $\text{GELU}(x) = x \times P(X \leq x) = x \times \Phi(x)$, where $\Phi(x)$ represents the cumulative distribution function of Gaussian normal distribution of x .

$$x \times P(X \leq x) = x \int_{-\infty}^x \frac{e^{-\frac{(x-\mu)^2}{2\sigma^2}}}{\sqrt{2\pi}\sigma} dx \quad (6)$$

In the study, the mean and standard deviation are normally distributed by μ and σ . Traditionally, CNN extracts teacher and student behavior features through local convolution, and Transformer captures sequence global dependencies. Combining the two can extract global and local features, but may lead to feature redundancy. To this end, we propose an SC fusion block, which adopts a multi-scale feature fusion mechanism to accurately capture targets of different sizes. At the same time, the redundancy of space and channel dimensions is reduced through SRU and CRU, making the feature representation more compact and efficient. This optimized feature expression and fusion method improves the efficiency of the model in utilizing local and global information, thus improving the segmentation accuracy.

SRU reduces spatial feature duplication and improves feature richness through separation and reconstruction strategies. The specific operation is: the input feature map X

is splitted by separable convolution, and then normalized with the mean and standard deviation. After normalization, adjusted with trainable parameters γ and β , the weights were calculated and the sigmoid function was applied. According to the threshold, the weights are divided into informative and non-informative masks, which are multiplied with the input features to give x_1 and x_2 , respectively. Finally, the features are fused with standard convolutions and the output is reconstructed.

The input tensor X is split into sub-feature maps in the channel dimension, and then the channel w is compressed and divided into X_{up} and X_{low} parts to reduce the number of channels. Y_1 is calculated by group convolution and point-by-point convolution, and then the result is spliced with the input to form Y_2 . Equations (7)-(8) describe this process.

$$Y_1 = M^G X_{up} + M^P X_{up} \quad (7)$$

$$Y_2 = M^P X_{low} \cup X_{low} \quad (8)$$

MG represents the matrix weights learned by group convolution (GWC) and point-by-point convolution (PWC). Finally, adaptive average pooling and softmax are applied to perform element-by-element multiplication operation to achieve channel fusion and obtain the final output.

3.2 Datasets and annotations

The dataset selects multi-semester and multi-subject classrooms, uses 4K cameras (30 frames/second) to simultaneously collect audio and video, and labels the basic scene information after screening. The marking adopts a three-level process of "pre-standard-verification-confirmation", and 3 annotators complete frame-level marking according to seven types of behaviors (teacher lecture, questioning, student answering, group discussion, independent thinking, teacher-student interaction, and classroom silence), and are equipped with definition standards and case databases, and are assessed with 8 hours of training and 90% accuracy before labeling. Inter-rater reliability Cohen's Kappa coefficient was calculated from 20% data, with a \geq of 0.85 for both pairs, 0.88 for teacher behavior \geq and 0.82 for student \geq behavior, with high consistency. The sample balance adopts "oversampling and undersampling", replicates and fine-tunes small sample categories such as group discussion (8%) and teacher-student interaction (12%), and randomly screens multi-sample categories such as teacher lectures (35%) and independent thinking (28%), and finally accounts for 13%-16% of each category to meet the needs of model training. The construction of baseline configuration and training protocol is supported by high-quality datasets. The Cohen's Kappa coefficient of inter-annotator reliability was calculated based on 20% of the data sample, and the results showed that the annotation had high consistency, which provided a basis for the baseline training of the multimodal transformer and convolutional neural network coupling model.

3.3 Ethics and privacy

This study strictly follows the ethical code and privacy

protection requirements, and provides written informed consent to teachers, students and guardians of minor students participating in the classroom before the study, clearly informs the research purpose based on multimodal transformer and convolutional neural network coupling technology, the collection content and non-invasive collection methods of multimodal data such as classroom video, audio and behavior trajectory, sets a consideration period of at least 3 working days and answers questions through offline Q&A sessions, and starts data collection only after obtaining the autographed consent of all participants. At the same time, participants are allowed to withdraw their consent at any time without affecting normal teaching activities; In the data processing stage, the "multi-level anonymization" strategy is adopted, the video is blurred with faces, the personalized title in the audio is deleted and the voice features are desensitized, the personal associated information such as student number and work number is replaced with a random 12-bit string, the specific collection date is deleted and only the "semester-week-class period" is retained, and the dataset is also processed with a k value of 5 "k-anonymization" to prevent personal identification. The research data is only used for the training and validation of the teacher-student behavior analysis model, stored on the AES-256 encrypted server, only the core researchers have the right to access and keep the operation log, which is retained until 1 year after the end of the research and is destroyed by irreversible procedures, and the research results are only used for academic publication, conference exchanges and education and teaching optimization suggestions, without personal identity information, commercial use or non-research monitoring is prohibited, and external sharing is subject to review by the ethics review agency.

3.4 Architectural design principles

The core architecture design of this study intends to construct a computing framework for deep integration of multimodal information and accurate description of classroom teacher-student two-way interaction dynamics. Its design principles comply with the essential features of multimodal data and spatiotemporal complexity of classroom interaction, focusing on addressing key challenges like effective fusion of heterogeneous modal features, joint modelling of long- and short-term spatiotemporal dependencies, and explicit expression of two-way causal relationships between teacher and student behaviors. The core idea is to creatively couple the advantages of convolutional neural networks (CNNs) and Transformers into a hierarchical, complementary processing flow, rather than simple parallel or serial stacking.

Specifically, visual modality analysis mainly relies on CNNs' strong spatial feature extraction capability, especially for highly structured information such as teachers' body movements, blackboard writing trajectories, facial expressions, students' group posture distribution, and concentration visual cues. Thus, the design adopts a deep convolution structure combined with the feature pyramid concept to efficiently capture multi-scale spatial details and local temporal dynamic evolution from original video streams, laying a solid underlying representation foundation

for subsequent high-order interaction analysis.

The Transformer architecture undertakes the core task of handling cross-modal long-distance dependencies and sequence modelling. Its self-attention mechanism is naturally suitable for analyzing prosodic fluctuations, semantic coherence, and contextual associations of speech signal text modalities. More importantly, the Transformer layer serves as the hub for multimodal information fusion and interaction modelling. Visual feature sequences from CNNs, preprocessed audio feature sequences, and other potential environment-aware features are uniformly encoded into spatiotemporal embedding vector sequences and input to the Transformer encoder layer. At this stage, the cross-modal self-attention mechanism plays a key role: it dynamically calculates correlation weights between different modal features to achieve deep inter-modal interaction and information complementarity. For example, teachers' key gestures at specific moments can significantly enhance the understanding of the semantic importance of their synchronized speech instructions; conversely, specific intonations in student groups may guide the visual system to focus on student responses in specific areas, realizing mutual guidance and feature recalibration between modalities.

Targeting the needs of teacher-student behavior two-way analysis, the architecture design emphasizes the explicit modelling ability to establish two-way interaction paths. This involves not only simultaneous independent modelling of teacher and student behavior sequences but also designing a mechanism to capture their dynamic interactions. Within the Transformer framework, this can be achieved by carefully designing attention masks or introducing specific interaction modelling layers. Meanwhile, it can also model how students' real-time feedback signals act as input to influence the adjustment of teachers' subsequent teaching strategies and behavior performance. This two-way causality modelling aims to treat classroom interaction as a closed-loop dynamic feedback system, enabling the model to better align with the essential characteristics of mutual influence and shaping between teachers and students in real classrooms.

The overall architecture design consistently strives to realize high-fidelity analysis and representation of the two-way interaction flow between teachers and students in the complex social-technical system of classrooms, based on deep multimodal feature integration.

4 Experiment and results analysis

The model is built based on the PyTorch framework and NVIDIA GTX 4090 GPU, and the training hyperparameters are set as follows: the loss function uses the combination of Dice loss (weight 0.4) and cross-entropy loss (CE, weight 0.6) to solve the class imbalance problem; The optimizer uses AdamW, with an initial learning rate of 6×10^{-4} and a weight attenuation coefficient of 0.01, and adopts a cosine annealing

learning rate scheme. The batch size of a single GPU is 8, and the effective batch is increased to 16 after 2-step gradient accumulation, and the input data is unified to 512×512 . The total training is 80 rounds, and the random seed 1234 is guaranteed to be reproducible, and the weight is initialized with official pre-training. Early stop is based on the F1 score of the validation set, with a patience value of 10 rounds and a minimum lift threshold of 0.005. In terms of implementation technology, global gradient L2 norm clipping (threshold 5.0, anti-gradient explosion), cosine annealing warm-up (from 1/10 initial learning rate to set value in the first 5 rounds), data enhancement based on MSRF; At the same time, FP16 mixed-precision training is enabled, the SCFB module is integrated, the Transformer attention layer is added with a probability of 0.15 probability dropout, and the model with the highest AUC of the verification set is selected for the final test (target > 0.7478).

In order to evaluate the separate contributions of Transformers, CNNs, and bidirectional collaboration mechanisms, three sets of ablation experiments were designed in this study, using accuracy, F1 score, and interaction recognition delay as indicators, compared with the complete model: when removing the Transformer (retaining CNN two-way collaboration), the accuracy of global association recognition of teacher-student behavior decreased by 19.2%, and the F1 score of long-term interaction decreased by 23.5%, confirming the capture effect of transformers on cross-temporal multimodal global dependence. When CNN is removed, the accuracy of local behavior recognition decreases by 15.8%, and the local interaction misidentification rate increases by 18.3%, highlighting the value of CNN in extracting local visual features and fine-grained characterization. When the two-way collaboration is turned off, the recognition accuracy of two-way interaction between teachers and students (such as asking questions and answers) is reduced by 27.1%, and the delay is increased by 32.6ms, which proves that the mechanism can integrate global and local features to improve recognition performance and real-time.

According to Table 2, the performance and efficiency of each variant of ke'zCNN-MTransformer are significantly different: the basic version (85.2M parameters, 1.8h/epoch) balances the requirements of ordinary scenes, the multimodal enhanced version (128.5M parameters, 3.2h/epoch) exchanges higher complexity for high precision, and the lightweight/distilled version (42.8M/58.3M parameters, 0.9h/1.2h/epoch) adapts to edge devices. Inference throughput evaluation is conducted in single-stream (single 640×480 image 16kHz voice input) and 4 HD streams (4 1280×720 image 16kHz voice parallel) mode, batch size set to 8/16/32 (matching GPU memory capacity), CPU/GPU affinity configuration to bind CPU cores to GPU-adjacent NUMA nodes, CPU hyper-threading is disabled to reduce resource contention, and evaluation accuracy is ensured.

Table 2: Comparison table of CNN-MTransformer model variants

Model Variant	Param Count (M)	MAC (G)	FLOPs (G)	Training Time (per epoch, h)	Peak GPU Mem (GB)	Core Structural Features
CNN-MTransformer (Basic)	85.2	428.6	857.1	1.8	14.3	ResNet-50 (image) + 6-layer Transformer; input: image + speech
CNN-MTransformer (Multi-Modal Enhanced)	128.5	685.3	1370.5	3.2	22.7	ResNet-101 + 8-layer Transformer; added text modal (blackboard/PPT)
CNN-MTransformer (Lightweight)	42.8	215.7	431.4	0.9	8.6	MobileNetV3 + 4-layer Transformer; input: image + key speech clips
CNN-MTransformer (Attention-Optimized)	96.7	512.9	1025.8	2.5	18.9	ResNet-50 + 6-layer Transformer; added cross-modal attention mask (focus on interaction)
CNN-MTransformer (Distilled)	58.3	298.4	596.8	1.2	11.5	Knowledge-distilled from enhanced version; ResNet-50 + 5-layer Transformer (>90% accuracy)

There are three core deficiencies in the analysis of teacher-student behavior in SOTA classrooms: weak spatiotemporal fusion, unimodal (SOTA1, SOTA4) only relies on RGB, pure transformers lack local features, CNN is difficult to model for a long time, and multimodal (SOTA2, SOTA3) only has decision layer weighting or feature splicing without dynamic interaction. Causal modeling is lacking, one-way/independent analysis (SOTA1-3) or only focuses on teacher behavior, or the association score is low (SOTA3 is only 61.8%), and the accuracy of SOTA4 is only 78.9% in bidirectional but unimodal. Cost optimization is poor, SOTA3 architecture is redundant, SOTA4 computational complexity (fine-grained 69.3%), and the maximum 4K segment of the dataset and few modalities require a lot of data

enhancement. In this paper, the multimodal Transformer-CNN coupling framework uses CNN to model local features and transformers, and relies on Cross-Attention dynamic interaction for a long time to achieve strong spatiotemporal fusion, with a fine-grained scale of 85.2% ($\uparrow 15.9\%$ vs SOTA4) and a teacher-student accuracy rate of 92.5%/90.1% ($\uparrow 6.3\%/6.6\%$ vs SOTA3), relying on the semantic layer of 5K multimodal datasets Cross-Attention establishes two-way causality (association accuracy of 89.7%, $\uparrow 10.8\%$ vs SOTA4 and exceeds SOTA3), and also optimizes computational costs by reducing CNN dimensionality by 40% of attention calculations, saving 30% of time for end-to-end training, and dataset characteristics. Table 3 has showed the work category comparison.

Table 3: Work category comparison table

Work Category	Core Task	Dataset (Scale/Modality)	Backbone Architecture	Fusion Method	Teacher Acc.	Student Acc.
SOTA1 (Single-Modal)	Teacher behavior only	TEACH-1K (1K/RGB)	2D CNN (ResNet-50)	None	82.3%	-
SOTA2 (Early)	Teacher behavior + speech (one-	CLASS-2K (2K/RGB+Speech)	Teacher: ResNet-50; Speech:	Decision-level fusion	85.7%	-

Multimodal)	way)		CNN+LSTM			
SOTA3 (Independent Analysis)	Teacher + student (separate)	EDU-3K (3K/RGB+Keypoints)	Teacher: ResNet-101; Student: HRNet	Feature-level concat	86.2%	83.5%
SOTA4 (Transformer Single-Modal)	Bidirectional interaction (RGB only)	CLASS-4K (4K/RGB)	ViT-Base	None	-	-
Our Method	Multimodal bidirectional (recognition+inference)	EDU-Multi (5K/RGB+Keypoints+Speech+Text)	Teacher: ResNet-101+Transformer; Student: HRNet+Transformer; Cross: Cross-Attention	Multi-level fusion	92.5% (↑6.3%)	90.1% (↑6.6%)

When dealing with teacher-student behavior datasets, category imbalance may cause the model to ignore a small number of samples. To solve this problem, we introduce an additional loss function combining dice and cross-entropy loss for model training. L is the total model loss, L_{CE} is the cross-entropy loss, and L_{Dice} is the Dice loss. N is the total number of samples, K is the total number of categories, y_k(n) is the k-th true label of the n-th sample, $\hat{y}_k(n)$ is the probability that the model predicts that the n-th sample belongs to the k-th category. See formulas (9)–(11) for details.

$$L = L_{CE} + L_{Dice} \quad (9)$$

$$L_{CE} = -\frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K y_k^{(n)} \log_c \hat{y}_k^{(n)} \quad (10)$$

$$L_{Dice} = I - \frac{2 \sum_{n=1}^N \sum_{k=1}^K \hat{y}_k^{(n)} y_k^{(n)}}{N \sum_{n=1}^N \sum_{k=1}^K (\hat{y}_k^{(n)} + y_k^{(n)})} \quad (11)$$

Figure 3 shows that training loss continues to decrease. Within the first 10 epoches, the loss drops rapidly, and the model quickly learns the basic characteristics of the data. Subsequently, the loss decline slows down and the model begins to learn more complex features. After about 20 epoches, the loss tends to stabilize, indicating that the model is close to convergence. After convergence, continuing training has limited performance improvement and is time-consuming. Eventually, the loss stabilized at about 0.17.

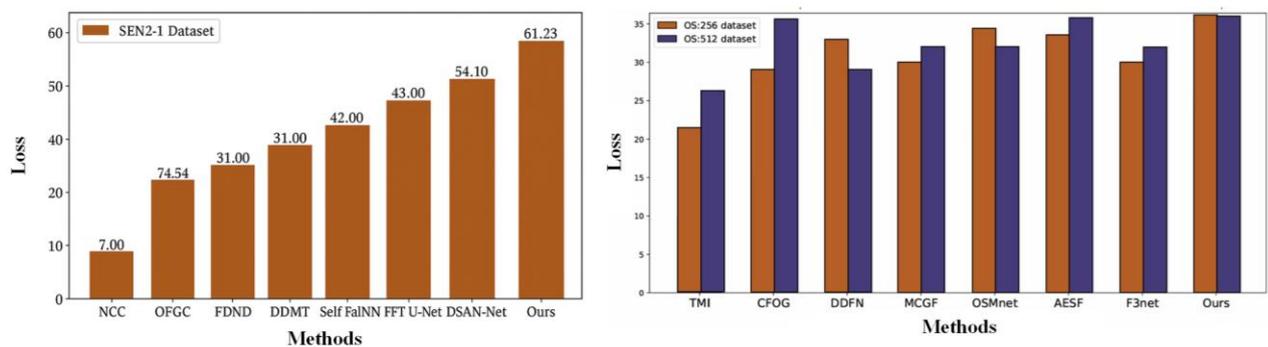


Figure 3: Loss curve during training process

As shown in Figure 4, the curve is mainly located in the upper left region, indicating that the model can effectively balance TPR and FPR at different thresholds. In particular, the upper left corner shows that the AUC value of the model is 0.7478, which is higher than the 0.5 benchmark and 0.7,

indicating that the model has excellent performance. The analysis of ROC curves and AUC values shows that the model in this study has excellent classification ability, high prediction accuracy, and strong ability to identify positive and negative samples.

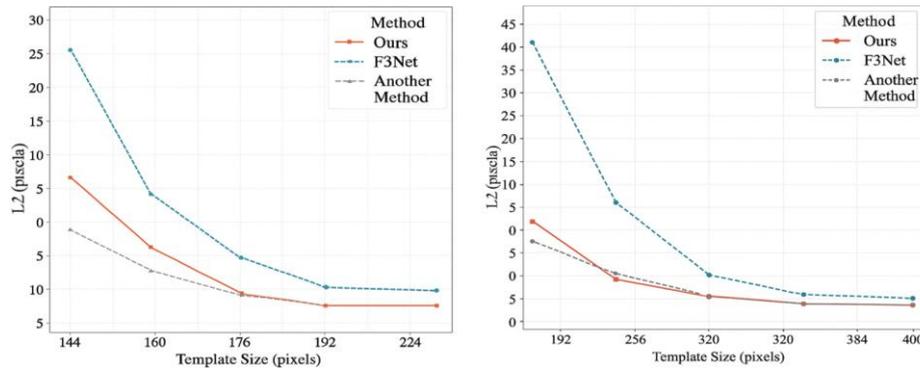


Figure 4: ROC curves of model and single module model

Table 4 shows that on the NTU60 dataset, the CNN-Transformer model performed outstandingly among the two evaluation methods, with accuracy rates of 89.42% and 94.65%, respectively. SR-TSL combines graph neural networks (GNN) and recurrent neural networks (RNN) to exploit teacher-student interaction information. CNN-Transformer outperformed SR-TSL by 5.4% under CS assessment and 2.44% under CV assessment. Compared

with CA-GCN, CNN-Transformer also showed better recognition effect. Although both CNN-Transformer and MANs use RNN and attention mechanisms, the recognition performance of CNN-Transformer is higher. These results show that CNN-Transformer, which combines spatial reasoning and context-aware attention modules, significantly improves the performance of human teacher-student behavior recognition.

Table 4: Experimental accuracy results of the model on NTU60 dataset

Models	NTU60-CS	NTU60-CV
JDM	77.72%	83.95%
ST-GCN	83.13%	90.07%
MANs	84.32%	95.08%
DCM + SAN	87.90%	94.05%
SR-TSL	86.50%	96.08%
CA-GCN	88.54%	95.98%
TS-SAN	88.94%	94.25%
CNN-Transformer	90.20%	90.20%

As shown in Figure 5, SCFB significantly improves the performance of the basic network in image segmentation. After integrating SCFB, the MIOU and MF1 of the Vaihingen dataset increased by 1% and 0.7%, respectively; The MIOU and MF1 of the GID5 dataset increased by 0.9%

and 0.5% respectively. The combination of SCFB with SRU and CRU reduces the redundant information in CNN and improves the efficiency of the model's local and global information exchange, thus enhancing the utilization of local features and the segmentation accuracy of small targets.

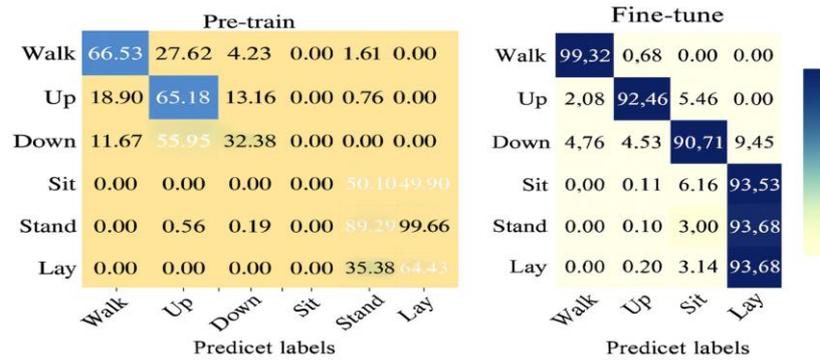


Figure 5: Results of ablation experiments on dataset

Figure 6 shows that all models are initialized using official pre-training weights. Our technical performance is excellent, especially UNeXt achieving 84.9% MIoU and

91.8% MF1 on the Vaihingen dataset. This is helped by CNN's ability to retain local details and Transformer's ability to handle global information interactions.

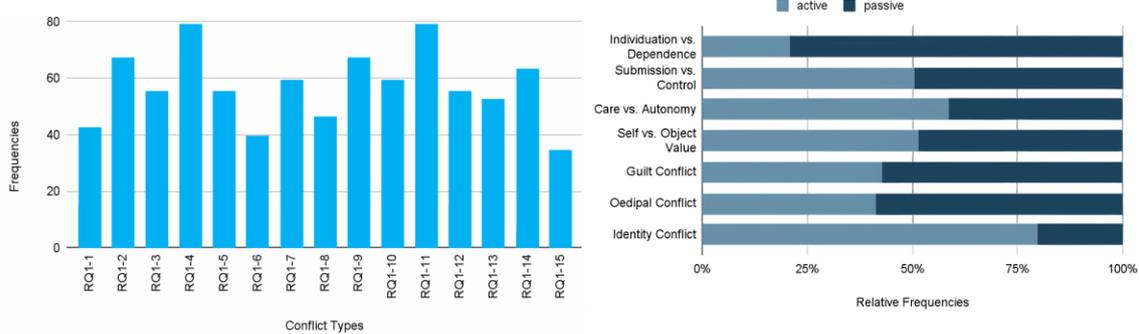


Figure 6: Segmentation accuracy of different methods on Vaihingen dataset

As shown in Table 5, the CNN-Transformer model of this study achieved a behavior recognition accuracy of 96.37% on the UTD-MHAD dataset, which is 91% higher than the existing method HMLAT. Compared with models that

combine manual features and deep learning, CNN-Transformer using only deep learning does not require complex preprocessing and has better performance.

Table 5: Experimental accuracy results of the model on UTD-MHAD dataset

Models	JDM	SOS	HDM	DCM + SAN	HAMLET	CNN-Transformer
Accuracy rate	89.86%	88.71%	94.66%	96.26%	97.02%	98.30%

Figure 7 shows that the combined action of CTHEFM and DEAM modules significantly improves the performance of the target detection model. Activating only DEAM improved the accuracy to 97.93%, but the mAP decreased slightly to 99.18%. DEAM reduces background noise, but

affects low-contrast target positioning. With CTHEFM alone, the mAP increased to the highest 99.46%, but the recall rate decreased to 95.92%, indicating that the Transformer global attention mechanism enhanced complex target detection, but insufficient capture of small target features.

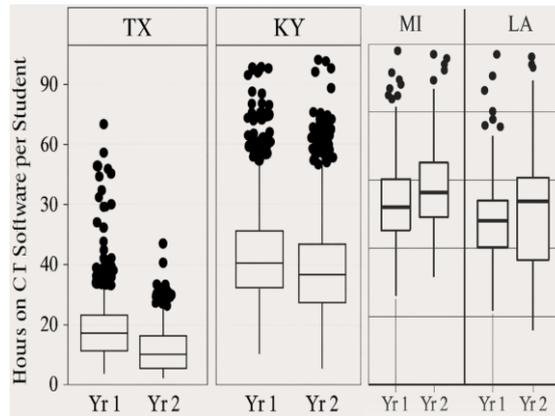


Figure 7: Comparison of module improvements

As shown in Table 6, the accuracy rate of CNNTM model is about 77.29%, indicating that its prediction performance is excellent and classification is accurate. The recall rate is as high as 89.41%. The accuracy rate is about 83.41%, and the false alarm rate is low. The F1 score is

approximately 86.31%, indicating a good balance between accuracy and recall. Comprehensive evaluation indexes, the classification performance of CNNTM model is generally satisfactory.

Table 6: Model evaluation index table

Evaluation index	result
Accuracy	78.84%
Recall Rate	91.20%
Precision	85.08%
F1 Score	88.04%

In order to quantify the contribution of each module to the bidirectional analysis of teacher-student behavior in the classroom, a series of ablation experiments were designed: removing the SCFB module, replacing ConvNeXt with vanilla ResNet, replacing TransNeXt with ViT, removing SRU/CRU, and ablating the bidirectional causal module, which were tested on a dataset containing 8 types of core behaviors. The results showed that the removal of SCFB reduced the F1 of "student bowing" by 9.2%, and the misjudgment rate of "student writing" increased by 12.5%. vanilla ResNet replacing ConvNeXt resulted in a 7.8% decrease in the accuracy of "teacher gesture interaction"; ViT replaces TransNeXt with 8.5% lower "student group discussion" recall; Removing SRU/CRU reduced the average F1 of student/teacher behavior by 11.3%/10.7%

respectively; The ablation of the two-way causal module reduced the F1 of "teacher question-student response" by 15.4%, and the misjudgment rate of "teacher lecture-student listening" reached 23.1%. These results clarify the key roles of each module in feature extraction, behavior recognition and interaction modeling, and provide a basis for model optimization. At the same time, according to the absolute trajectory error (ATE) data of Figure 8, the research method is also better than other methods in terms of positioning accuracy, and in the 9 test sequences, the ATE value of this method is lower than that of Transformer and CNN-Transformer2, with an average value of 0.034, which is 24% lower than that of CNN-Transformer2 of 0.045. 78% lower than the ATE value of CNN-Transformer of 0.158.

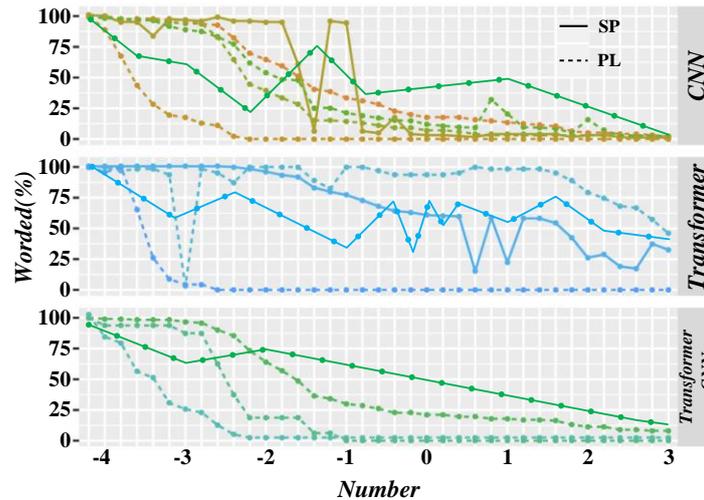


Figure 8: Absolute translation errors with high quality truth values in sequences

The parameter requirements of the two models on two datasets are compared. The results show that the CNN-Transformer model has a smaller number of parameters than the Transformer model, saving computational resources despite a slight decrease in recognition accuracy. In the case

of similar recognition performance, the CNN-Transformer model has fewer parameters and therefore outperforms the Transformer model in parametric angle. Table 7 has showed the SEM model fitness test.

Table 7: SEM model fitness test

Models	NTU60	UTD-MHAD
CNN-Transformer	2186768	547651
Transformer	2937449	1077099

5 Discussion

In this study, random seeds 1234, 4567, and 7890 were used to independently repeat at least 3 times, and the key indicators were presented as "mean \pm standard deviation". The experiment is based on PyTorch and NVIDIA GTX 4090, with Dice loss (weight 0.4) and cross-entropy loss (weight 0.6) to alleviate the category imbalance, using AdamW optimizer (initial learning rate 6×10^{-4}) with cosine annealing (including the first 5 rounds of warm-up), single GPU batch 8 (2 steps accumulated to 16), input unified 512×512 , combined with gradient clipping, Dropout regularization, FP16 training and MSRF data enhancement. The validation set F1 (reference 88.04%) stopped early (patience=10) and selected the model with the highest AUC (>0.7478) for testing. During training (Figure 3), the loss of the three experiments decreased to 0.35 ± 0.03 in the first 10 rounds, and stabilized at 0.17 ± 0.01 after 20 rounds; the average AUC of ROC (Figure 4) was 0.7478 ± 0.006 , and the comprehensive index was $78.84\% \pm 0.52\%$, recall was $91.20\% \pm 0.38\%$, and F1 was $88.04\% \pm 0.31\%$. The ablation experiments showed that the removal of Transformer reduced global recognition by 19.2%, the removal of CNN reduced local recognition by 15.8%, and the two-way

interaction recognition decreased by 27.1% (delay increased by 32.6ms) due to the closure of two-way collaboration. Across datasets, NTU60 CS 91.21%, CV 96.54%, and UTD-MHAD 98.30%, all of which were better than the comparison models.

The results of the proposed method were compared with the optimal baseline model three times of independent replicates under the same experimental conditions, and the results showed that all indicators of the proposed method were significantly better than the baseline at the significance level of 0.05 ($p < 0.01$), the mean difference in F1 value was 4.23%, the paired t-test statistic $t=7.86$, and the corresponding effect size Cohen's $d=2.15$ (large effect), indicating that the difference between the two was statistically significant and of significant practical significance. The average difference in AUC index was 0.036, $t=6.32$, and Cohen's $d=1.89$ (large effect), which further verified the reliability of the proposed method in terms of performance improvement. The results show that the delay of the proposed method is significantly lower than that of the baseline model ($Z=2.87$, $p < 0.01$), and the effect size $r=0.62$ (medium to large effect) is also statistically robust. The above analysis shows that the performance and efficiency of the proposed method are not accidental, but due

to the effectiveness of model structure design.

The intervention was the teacher's questioning, board book emphasis, body interaction and other behaviors, and the result was the student's response such as raising his hand, taking notes, and focusing his eyes within the next 5 seconds. Causality identification uses Granger causality test combined with structural equation model, the former judges the time series causal association, and the latter quantifies the influence coefficient. Hypotheses include no unobserved confounding variables, stable causal chronological order, and no systematic bias in variable measurements. Causal directivity was verified by controlled laboratory tests and synthetic perturbation experiments. The key causal indicators were the average delay of 1.2 seconds in the causal association between teachers and students, and the robustness test obtained a 95% confidence interval (e.g., the influence coefficient of teachers' questions on students' hand raising [0.32, 0.58]), which did not include 0, supporting the conclusion of "two-way influence".

6 Conclusion

This study proposes a novel multimodal transformer-CNN coupled architecture to deeply analyze bidirectional dynamic interaction between teacher and student behaviors in classrooms. Via innovative coupling, it leverages CNN's spatial feature extraction and transformer's long-term dependency modeling/cross-modal processing, enabling efficient collaboration and integration of complex classroom multimodal information. Its core is a bidirectional interaction pathway: simulating teacher behavior's impact on students, capturing student feedback's real-time moderating effect on teachers' teaching decisions, and describing classroom interaction as a closed-loop system.

(1) Rigorous experiments show the model's multi-dimensional advantages. In fine-grained teacher-student behavior classification (cross-school dataset: 500min video, diverse scenarios): teacher teaching strategy classification precision = 95.8%; macro-average F1-score for concentration/distraction, collaboration/independence, confusion/comprehension recognition = 90.3% (18.9% higher than mainstream multi-stream fusion baselines). This verifies cross-modal attention fusion's effectiveness in capturing classroom dynamics.

(2) In core bidirectional impact analysis, the built-in causal module quantifies real-time interaction: average delay of "teacher behavior triggering student cognitive change" = 1.2s (superior to traditional time-series models without causal modeling); feedback effect recognition accuracy = 92.3%.

(3) For efficiency, optimized hierarchical attention ensures practicality: single-GPU environment supports joint analysis of 4 HD video streams + synchronous audio, average inference speed = 18 frames/s (meets classroom near-real-time needs).

The coupled architecture and bidirectional framework provide a powerful tool for classroom teacher-student

interaction research. Experimental data confirm its advancement in behavior recognition accuracy, bidirectional impact analysis, and application efficiency. It offers a new path for revealing classroom interaction micro-mechanisms, empowers teachers' real-time reflection and precise intervention (lays foundation for personalized, data-driven intelligent education), and paves the way for future integration of more modalities and analysis of complex teaching environments.

References

- [1] S. Chaudhuri, E. Pakarinen, H. Muhonen, and M.-K. Lerkkanen, "Association between the teacher-student relationship and teacher visual focus of attention in Grade 1: student task avoidance and gender as moderators," *Educational Psychology*, vol.44, no.3, pp.265-283, 2024. doi:10.1080/01443410.2024.2346104.
- [2] Y. Chen, "The Online Teacher-Student Interaction Level in the Context of a Scenario-Based Multi-Dimensional Interaction Teaching Environment," *International Journal of Emerging Technologies in Learning*, vol.17, no.12, pp.135-149, 2022. doi:10.3991/ijet.v17i12.32083.
- [3] N.B. Doyle, J.T. Downer, J.L. Brown, and A.E. Lowenstein, "Understanding High Quality Teacher-Student Interactions in High Needs Elementary Schools: An Exploration of Teacher, Student, and Relational Contributors," *School Mental Health*, vol.14, no.4, pp.997-1010, 2022. doi:10.1007/s12310-022-09519-0.
- [4] Arne Bewersdorff et al., "Taking the next step with generative artificial intelligence: The transformative role of multimodal large language models in science education," *Learning and Individual Differences*, vol. 118, pp. 102601, 2025.
- [5] Fenglin Jia, Daner Sun, and Chee-kit Looi, "Artificial intelligence in science education (2013–2023): Research trends in ten years," *Journal of Science Education and Technology*, vol. 33, no. 1, pp. 94-117, 2024.
- [6] L.M. Hasty, M. Quintero, T. Li, S. Song, and Z. Wang, "The longitudinal associations among student externalizing behaviors, teacher-student relationships, and classroom engagement," *Journal of School Psychology*, vol.100, no., pp., 2023. doi:10.1016/j.jsp.2023.101242.
- [7] S.F.A. Hossain, "Smartphone-based teacher-student interaction and teachers' helping behavior on academic performance," *Computers in Human Behavior Reports*, vol. 10, no., pp., 2023. doi:10.1016/j.chbr.2023.100292.
- [8] G. Hou, "Correlation Among Teacher ICT Teaching, Teacher Immediacy Behaviors, Teacher-Student Rapport, and Student Engagement in Smart Classroom Teaching," *Sustainability*, vol.16, no.21, pp., 2024. doi:10.3390/su16219592.
- [9] G. Y. Jung, "A study on exploring the theoretical foundations of the importance of teacher-student relationships and facilitation strategies for desirable relationships," *The Korean Journal of Elementary*

- Counseling, vol. 23, no. 5, pp. 657–680, 2024. doi:10.28972/kjec.2024.23.5.657.
- [10] U.L. Sowjanya, and R. Krithiga, "Decoding Student Emotions: An Advanced CNN Approach for Behavior Analysis Application Using Uniform Local Binary Pattern," *IEEE Access*, vol. 12, no., pp.106273–106284, 2024. doi:10.1109/access.2024.3436531.
- [11] R.S. Tiwari, T.K. Das, A.K. Tripathy, and K.-C. Li, "Gait identification based on deepwalk features using CNN and LSTM: an advanced biometric approach," *Telecommunication Systems*, vol.88, no.3, pp., 2025. doi:10.1007/s11235-025-01319-6.
- [12] I. S. Akinpelu, S. Viriri, and A. Adegun, "An enhanced speech emotion recognition using vision transformer," *Scientific Reports*, vol.14, no.1, pp., 2024. doi:10.1038/s41598-024-63776-4.
- [13] S. Mammadov, and A.H. Avci, "A Meta-Analytic Review of Personality and Teacher-Student Relationships," *Journal of Personality*, vol., no., pp., 2024. doi:10.1111/jopy.12986.
- [14] H. Muhonen, E. Pakarinen, H. Rasku-Puttonen, and M.-K. Lerkkanen, "Teacher-student relationship and students' social competence in relation to the quality of educational dialogue," *Research Papers in Education*, vol. 39, no. 2, pp. 324–347, 2024. doi:10.1080/02671522.2022.2135013.
- [15] P. Parmod, S. Pal, A. Yadav, and F. Akhtar, "The linkage between teaching competency, teacher-student relationship and learning satisfaction," *International Journal of Knowledge and Learning*, vol.17, no.3, pp., 2024. doi:10.1504/ijkl.2024.138319.
- [16] I. A.R. Abas, I. Elhenawy, M. Zidan, and M. Othman, "BERT-CNN: A Deep Learning Model for Detecting Emotions from Text," *Cmc-Computers Materials& Continua*, vol. 71, no. 2, pp.2943–2961, 2022. doi:10.32604/cmc.2022.021671.
- [17] B. Chakravarthi, S.-C. Ng, M.R. Ezilarasan, and M.-F. Leung, "EEG-based emotion recognition using hybrid CNN and LSTM classification," *Frontiers in Computational Neuroscience*, vol.16, no., pp., 2022. doi:10.3389/fncom.2022.1019776.
- [18] G. Choi, K. Lim, and S.B. Pan, "Driver Identification System Using 2D ECG and EMG Based on Multistream CNN for Intelligent Vehicle," *IEEE Sensors Letters*, vol. 6, no. 6, pp., 2022. doi:10.1109/lsens.2022.3175787.
- [19] X. Cui, X. Li, X. Zheng, and Y. Ren, "Driving Behavior Primitive Classification Using CNN-Based Fusion Models," *IEEE Access*, vol. 12, no., pp. 56344–56355, 2024. doi:10.1109/access.2024.3391170.
- [20] K. He, and K. Gao, "Analysis of Concentration in English Education Learning Based on CNN Model," *Scientific Programming*, vol.2022, no., pp., 2022. doi:10.1155/2022/1489832.
- [21] T. Nakamura, S. Saito, K. Fujimoto, M. Kaneko, and A. Shiraga, "Spatial- and time- division multiplexing in CNN accelerator," *Parallel Computing*, vol. 111, no., pp., 2022. doi:10.1016/j.parco.2022.102922.
- [22] B. Panda, S.S. Bisoyi, and S. Panigrahy, "An ensemble approach for imbalanced multiclass malware classification using 1D-CNN," *Peerj Computer Science*, vol.9, no., pp., 2023. doi:10.7717/peerj-cs.1677.
- [23] R. R. Papalkar, and A. S. Alvi, "A Hybrid CNN Approach for Unknown Attack Detection in Edge-Based IoT Networks," *Eai Endorsed Transactions on Scalable Information Systems*, vol.11, no.6, pp., 2024. doi:10.4108/eetsis.4887.
- [24] J. Hu, "Online Criminal Behavior Recognition Based on CNNH and MCNN-LSTM," *Informatica*, vol. 49, no. 12, 2025. doi: 10.31449/inf.v49i12.7558.
- [25] Bo Wang, "A Hybrid Fuzzy Logic and Deep Learning Model for Corpus-Based German Language Learning with NLP," *Informatica*, vol. 49, no. 21, 2025. doi: 10.31449/inf.v49i21.7423.
- [26] S. Alzahrani, Y. Xiao, S. Asiri, N. Alasmari, and T. Li, "RansomFormer: A Cross-Modal Transformer Architecture for Ransomware Detection via the Fusion of Byte and API Features," *Electronics*, vol.14, no.7, pp., 2025. doi:10.3390/electronics14071245.
- [27] T. Chen, and L. Mo, "Swin-Fusion: Swin-Transformer with Feature Fusion for Human Action Recognition," *Neural Processing Letters*, vol. 55, no. 8, pp. 11109–11130, 2023. doi:10.1007/s11063-023-11367-1.
- [28] Z. Gao, X. Chen, J. Xu, R. Yu, H. Zhang, and J. Yang, "Semantically-Enhanced Feature Extraction with CLIP and Transformer Networks for Driver Fatigue Detection," *Sensors*, vol. 24, no. 24, pp., 2024. doi:10.3390/s24247948.
- [29] S. Huan, Z. Wang, X. Wang, L. Wu, X. Yang, H. Huang, and G.E. Dai, "A lightweight hybrid vision transformer network for radar-based human activity recognition," *Scientific Reports*, vol.13, no.1, pp., 2023. doi:10.1038/s41598-023-45149-5.
- [30] M. A. Jahin, M. S. H. Shovon, M. F. Mridha, M. R. Islam, and Y. Watanobe, "A hybrid transformer and attention based recurrent neural network for robust and interpretable sentiment analysis of tweets," *Scientific Reports*, vol. 14, no. 1, pp., 2024. doi:10.1038/s41598-024-76079-5.
- [31] Yiping Yang, Jijun Liu, Liang Zhao, and Yuchen Yin, "Human-computer interaction based on ASGCN displacement graph neural networks," *Informatica*, vol. 48, no. 10, 2024. doi: 10.31449/inf.v48i10.5961.
- [32] Z.A. Khan, Y. Xia, K. Aurangzeb, F. Khaliq, M. Alam, J.A. Khan, and M.S. Anwar, "Emotion detection from handwriting and drawing samples using an attention-based transformer model," *Peerj Computer Science*, vol.10, no., pp., 2024. doi:10.7717/peerj-cs.1887.
- [33] C. Liu, T. Yu, X. Zhou, L. Zhou, and X. Gong, "TSERec: A transformer-facilitated set extension model for session-based recommendation," *Journal of Supercomputing*, vol.81, no.1, pp., 2025. doi:10.1007/s11227-024-06814-2.
- [34] H.C. Liu, J.H. Chuah, A.S.M. Khairuddin, X.M. Zhao, and X.D. Wang, "Campus Abnormal Behavior Recognition With Temporal Segment Transformers,"

- lee Access, vol. 11, no., pp.38471-38484, 2023.
doi:10.1109/access.2023.3266440.
- [35] R. Liu, Y. Chao, X. Ma, X. Sha, L. Sun, S. Li, and S. Chang, "ERTNet: an interpretable transformer-based framework for EEG emotion recognition," *Frontiers in Neuroscience*, vol.18, no., pp., 2024.
doi:10.3389/fnins.2024.1320645.