Research on the Construction of Personalized Learning State Recognition Model Integrating Multimodal Data

Ziqi Meng¹, Xiaoping Huang^{2,*}

¹Agricultural University Of Hunan, Hunan, China ²Nanyang Technological University, Singapore *Corresponding Author

Abstract: This study addresses the limitations of traditional learning state recognition methods in complex educational scenarios by proposing a multimodal deep learning model that integrates visual, acoustic, physiological, and behavioral data. The hierarchical feature fusion architecture combines heterogeneous data sources, while a dynamic attention mechanism enables weighted feature selection. lightweight design ensures real-time performance. **Experimental** demonstrate significant improvements over both baseline methods in recognition accuracy and F1 scores, with enhanced environmental noise resistance cross-scenario generalization capabilities. The research confirms that multimodal comprehensively integration captures learners' cognitive and emotional states, providing a reliable personalized analysis tool for intelligent education systems. Future work should focus on optimizing cross-cultural adaptability and addressing long-term concept drift issues.

Keywords: Multimodal Fusion; Learning State Recognition; Attention Mechanism; Personalized leaRning; Educational Artificial Intelligence

1. Introduction

of educational With the advancement informatization, learning state recognition has become a pivotal factor in optimizing personalized learning **Traditional** paths. single-modal data struggles to comprehensively capture learners' dvnamic cognitive emotional while multimodal data states, integration offers innovative approaches for precise modeling. This study aims to develop a multimodal personalized learning recognition model that integrates visual, acoustic, and interactive behavioral information.

overcoming the limitations of single modality to provide more reliable decision support for adaptive learning systems. The research will focus on three aspects: theoretical framework design, algorithm optimization, and model validation, driving technological innovation and practical applications in the field of intelligent education.

2. Feature Analysis and Preprocessing of Multimodal Learning Data

2.1 Analysis of Multimodal Data Sources and Characteristics

The core value of multimodal learning data lies in its ability to capture learners 'state information from multiple dimensions, compensating for the limitations of single-source data. In personalized learning scenarios, multimodal data primarily encompasses four categories: visual, auditory, physiological signals, and interactive behaviors. Visual data typically includes facial expressions, eye-tracking, and pose estimation, which reflect learners' focus levels, emotional states, and cognitive load. Auditory data reveals learners 'emotional tendencies and language comprehension abilities through tone, speech rate, and semantic analysis. For instance, speech characteristics during classroom Q&A or oral practice can indirectly indicate cognitive engagement. Physiological signals like electroencephalography (EEG), heart variability (HRV), and galvanic skin response (GSR) provide direct neural and physiological feedback, suitable for high-precision attention or stress monitoring. Interactive behavioral data originates from learning system logs, including clickstream patterns, answer duration, and page transition frequency, which quantify learners' behavioral patterns and strategic preferences. Different modalities of data exhibit significant differences. Visual and speech data, high-dimensional time-series signals, require

handling of interference factors such as lighting, occlusion, or background noise. physiological signals are highly objective, their collection is costly and susceptible to individual variations. Interactive behavioral data, though structured, lacks direct representation of implicit cognitive states. This heterogeneity demands researchers to thoroughly understand the physical significance and applicability of each modality before data fusion. For instance, facial expressions can assist in assessing learners' immediate emotions, but must be combined with eye-tracking data to distinguish between "superficial focus" and "deep engagement." Coordinated analysis of speech and interactive behaviors can more accurately frustration or cognitive bottlenecks. Therefore, characterizing the properties of multimodal data not only forms the basis for preprocessing but also serves as a critical foundation for feature selection and weight allocation during model design [1].

2.2 Alignment and Standardization of Cross-modal Data

In multimodal learning state recognition, the core challenge of data fusion lies in effectively aligning and standardizing data from different sources, scales, and temporal characteristics. Visual, acoustic, physiological signals, and behavioral interaction data often exhibit distinct acquisition frequencies and spatiotemporal distribution patterns. Direct fusion may lead to information redundancy or semantic conflicts. Temporal alignment stands as the primary task in multimodal data processing, particularly for continuous monitoring in dynamic learning scenarios. For instance, eye-tracking data updates at millisecond intervals while answer logs are recorded in seconds, requiring techniques like interpolation, sliding windows, or dynamic time warping (DTW) to achieve temporal synchronization. Spatial alignment involves coordinates unifying across multi-sensor data, such as spatial registration between facial keypoint detection results and eye-tracking heatmaps, to analyze correlations between visual attention distribution emotional expression changes.

Normalization and standardization are employed to eliminate dimensional discrepancies and distributional biases across modalities. Since data values from different modalities may vary dramatically, direct fusion can diminish the contribution of low-amplitude features. Common solutions include normalization (scaling data to fixed ranges) and standardization (transforming data into a zero-mean, unit-variance distribution through Z-scores). For unstructured data like audio spectra or image features, deep embedding techniques are required to map them into a unified semantic space. For instance, pre-trained models can extract emotional embedding vectors from speech data, which are then measured for similarity with visual features in latent space. Additionally, cross-modal attention mechanisms dynamically learn weights between modalities, enabling adaptive feature-level fusion. These methods not only enhance data comparability but also establish high-quality input foundations for subsequent modeling processes [2].

2.3 Noise Filtering and Feature Dimension Reduction Strategy

Multimodal learning datasets are frequently contaminated by significant noise, including visual data interference from ambient lighting variations, physiological signal drift caused by poor device contact, and abnormal click behaviors in system logs. Noise filtering serves as a critical step to ensure model robustness, requiring tailored strategies based on noise types. For high-frequency random noise, wavelet transforms or low-pass filtering can effectively smooth signals. For sudden disturbances, statistical outlier detection methods can identify eliminate abnormal segments. processing noise in time-series data, contextual consistency must be considered-such modeling learners' state transitions using hidden Markov models to correct misjudgments caused by transient interference. In speech data, independent component analysis (ICA) can separate environmental reverberation from target sound sources, while deep learning denoising algorithms can restore clear features from motion blur in visual data.

Feature dimensionality reduction addresses the challenges of high dimensionality and sparsity in multimodal data. Direct use of raw features not only incurs high computational costs but may also introduce redundant information. Linear methods like principal component analysis (PCA) extract key components through variance maximization, suitable for structured behavioral data; nonlinear approaches such as t-SNE and UMAP preserve complex topological relationships between modalities, ideal for

visualization clustering tasks. Deep or achieve end-to-end feature autoencoders through bottleneck structures, compression excelling in processing unstructured data like images and speech. Additionally, domain-specific feature selection crucial-methods like information gain or mutual information can quantify feature-learning correlations to eliminate weakly discriminative variables. The reduced feature set should balance compactness and representativeness, avoiding overfitting while supporting fine-grained classification requirements [3].

3. Design of Personalized Learning State Recognition Model

3.1 Hierarchical Architecture of Multimodal Fusion

In the construction of multimodal learning state recognition models, hierarchical architecture design serves as the core strategy for achieving efficient feature fusion. This architecture typically consists of three main modules: the data layer, feature layer, and decision layer, each responsible for specific functions and working in synergy. The data layer processes raw multimodal inputs such as visual images, speech waveforms, physiological signals, behavioral logs. Through sensor synchronization technology, it ensures temporal consistency while performing preliminary noise reduction and alignment. The feature layer employs modality-specific extraction networks-convolutional neural networks process visual data, recurrent neural networks analyze time-series speech signals, and graph neural networks model interaction sequences-enabling representation high-order semantic independent modal spaces. The decision layer achieves cross-modal interaction through gating cross-attention, mechanisms or mapping heterogeneous features into a unified learning state space. This hierarchical design not only preserves modality-specific characteristics but also avoids information loss through progressive fusion. For instance, lower-level fusion captures fine-grained correlations, while higher-level fusion integrates abstract semantics. The flexibility of this hierarchical architecture allows adaptation to various educational scenarios, whether real-time monitoring in classroom environments or asynchronous analysis in online learning platforms. Optimal performance can be

achieved by adjusting fusion depth and granularity [4].

3.2 Dynamic Feature Weighting based on Attention Mechanism

The attention mechanism provides dynamic feature selection capabilities for multimodal learning state recognition, addressing the limitations of traditional static weighting methods in adapting to individual differences changes. scenario Bycalculating intra-modal and cross-modal attention weights, this mechanism automatically focuses on the most discriminative feature segments. Within a single modality, temporal attention identifies critical frames or speech segments-such as highlighting frequent blinking or frowning moments in facial videos, or marking emotional turning points with intonation shifts in speech streams. Cross-modal attention evaluates the contribution of different data sources, for example, reducing the weight of response durations in behavioral logs when physiological signals indicate high cognitive load to avoid information interference. redundant dynamic weighting is particularly effective for learning state transitions, such as from "focused" to "fatigued," allowing models to capture subtle signs through evolving attention distributions. Additionally, hierarchical attention design distinguishes between global states and local features-for instance, group learning employs global attention to identify common patterns, individual-level attention pinpoints specific learners 'abnormal behaviors. The explainability of the attention mechanism also enhances model credibility. Through visualizing weight distributions, educators can intuitively understand the model's decision-making rationale, thereby optimizing instructional intervention strategies [5].

3.3 Lightweight Model Deployment and Real-time Optimization

The practical application of personalized learning state recognition models must meet the requirements of lightweight and low latency, presenting dual challenges for algorithm design and engineering implementation. In terms of model architecture, replacing standard convolution with deep separable convolutions, applying distillation techniques to compress large pre-training models, and adopting sparse connections to reduce parameter size can

significantly lower computational costs while maintaining accuracy. For instance, replacing ResNet with MobileNet architecture for visual data processing or using TinyBERT for speech feature extraction can multiply inference speeds on edge devices by several times. Real-time optimization requires combining pipeline parallelism with edge computing strategies, distributing feature extraction and fusion tasks across different hardware units: GPUs handle image processing, CPUs run lightweight temporal models, and FPGAs accelerate attention computations. Additionally, dynamic sampling technology can adjust input resolution or frame rate based on system load, such as analyzing only key facial features during low-power scenarios or performing segmented non-uniform sampling of speech signals. These optimizations not only ensure smooth operation of models on resource-constrained devices like tablets and embedded sensors but also support large-scale concurrent processing, providing a technical foundation for real-time feedback in smart classrooms or online education platforms. The ultimate goal of lightweight deployment is to achieve "seamless monitoring," integrating learning state recognition naturally into teaching processes rather than becoming an additional burden [6].

4. Model Validation and Performance Evaluation

4.1 Experimental Design and Baseline Model Comparison

To validate the effectiveness of personalized learning state recognition models for multimodal data fusion, this study designed a systematic experimental protocol comprising three key components: dataset construction, baseline model selection, and evaluation framework design. The dataset incorporates multimodal data collected from authentic classroom environments, encompassing visual, acoustic,

physiological, and behavioral modalities. Through expert annotation, learning states were categorized into four types: focused, confused, fatigued, and distracted. cross-validation was employed to ensure result reliability, with training, validation, and test sets proportioned at 6:2:2. To comprehensively assess model performance, four representative baseline models were selected for comparison: an LSTM classifier based on single-modal data. an early fusion convolutional neural network, a late fusion random forest model, and a SVM classifier employing traditional engineering. These baseline models span from conventional machine learning to deep learning approaches, fully demonstrating the value of multimodal fusion. All experiments were conducted on a unified NVIDIA Tesla V100 GPU environment, with identical preprocessing procedures and hyperparameter tuning strategies applied across all models to ensure fairness [7]. Table 1 presents the comparative analysis between our model and the baseline model across four dimensions: accuracy, F1 score, inference latency, and parameter size. Our model achieves 89.7% accuracy and 88.3% F1 score, outperforming the best baseline models by 5.2 and 4.8 percentage points respectively, demonstrating the effectiveness of multimodal fusion. Notably, despite its large parameter size, the hierarchical architecture and lightweight design keep inference latency below 83ms, meeting real-time requirements. The comparison reveals limitations in single-modal models: the LSTM classifier using only visual data shows low recognition rates for "confused" states, while late fusion methods underperform in cross-modal association modeling. findings provide clear directions for future model optimization and validate our core innovation-achieving more precise learning state recognition through dynamic feature weighting and hierarchical fusion.

Table 1. Comparison of Model Performance

typesofmodels	precision(%)	F1score(%)	Inferencedelay(ms)	Numberofparameters(M)
single mode LSTM	76.5	74.2	45	3.2
Early fusion CNN	81.3	79.8	62	12.7
Late fusion random forest	79.6	77.5	28	-
Traditional feature SVM	72.8	70.1	15	-
The model	89.7	88.3	83	18.5

4.2 Construction of Multi-Dimensional Evaluation Indicators

In evaluating the performance of personalized learning state recognition models, a single metric often fails to fully reflect practical application

value. Therefore, multi-dimensional a assessment framework is essential. This system should conduct comprehensive evaluations four core dimensions: recognition accuracy, real-time responsiveness, robustness, and interpretability. The accuracy dimension encompasses traditional metrics like precision, recall, and F1 score, while incorporating differentiated weights for specific learning states-such as prioritizing "confusion" detection to ensure timely instructional interventions. The real-time responsiveness dimension focuses on engineering parameters including inference latency, memory consumption, and energy which efficiency, determine deployment feasibility in educational environments. The robustness dimension assesses stability under adversarial testing scenarios like lighting voice interference, and device variations. heterogeneity. The interpretability dimension employs methods such as attention visualization, feature contribution analysis, and decision path tracing to make model reasoning transparent and understandable to educators.

The development of this multidimensional evaluation framework must adhere to the unique requirements of educational scenarios. For instance, in accuracy assessment, it should not only evaluate overall performance but also analyze the model's adaptability to learners of different age groups and subject matter. Real-time testing requires simulating multi-device concurrent scenarios in actual classroom environments rather than ideal laboratory conditions. Robustness verification should focus on common teaching disruptions like projector flickering or background noise during group discussions. Explainability evaluation educational must integrate psychology theories to ensure the model's recognition of state characteristics aligns with human teachers 'experiential judgments. This comprehensive assessment strategy not only objectively highlights the model's technical

strengths but also reveals its potential limitations in practical education applications, providing clear directions for future optimization. The resulting evaluation system serves both the rigorous standards of academic research and the practical considerations of educational product implementation.

4.3 Analysis of Model Robustness and Generalization Ability

In practical applications of personalized learning state recognition models, robustness and generalization capabilities are key indicators of their engineering applicability. This study conducts systematic analysis from three dimensions: data distribution, environmental and interference. cross-scenario revealing the model's stable performance in complex educational environments. Regarding data distribution, test subsets with diverse learning styles, cognitive levels, and age groups were constructed to evaluate the model's adaptability individual differences. to particularly its recognition effectiveness for students with special educational needs, as shown in Table 2. Environmental interference tests simulated typical classroom noises such as sudden lighting changes, equipment movement, and background conversations, assessing the anti-interference performance of visual and features. Cross-scenario transfer acoustic experiments verified the model's knowledge transfer efficiency from laboratory environments to real classrooms, and from online learning platforms to offline teaching scenarios-a capability that directly determines the model's potential for large-scale application. Analysis results demonstrate that the hierarchical fusion architecture better resists performance degradation caused by local mode failures compared to traditional methods, while the dynamic attention mechanism significantly enhances the model's generalization recognition capability for unknown learning patterns.

Table 2. Comparison of Model Capability Analysis

=							
Assessment dimensions	test method	Key performance indicators	Advantage features				
Data robustness	Cross-group subset	Range of fluctuation of	Dynamic feature weighting				
	testing	recognition accuracy	mechanism				
Environmental robustness	Noise injection pressure	Performance decline	Multimodal complementary				
	test	magnitude	properties				
Short-term generalization	Test incremental data in	The improvement rate after	Feature decoupling				
ability	the same scenario	parameter tuning	representation				
Long-term generalization	Cross semester data	Concept drift adaptation speed	Online learning mechanism				
ability	validation	Concept diffi adaptation speed	Online learning mechanism				

Cross scenario	Different teaching	Zero sample recognition	Domain invariant feature
generalization capability	scenario migration test	accuracy	extraction

In-depth analysis reveals that the model's robustness advantages stem from three design features: The complementarity of multimodal data prevents noise from specific modalities from causing system-wide failures, hierarchical feature extraction preserves state representations at different granularities, and online fine-tuning mechanisms enable continuous adaptation to environments. Enhanced teaching generalization capabilities are attributed to attention mechanisms focusing on core features and lightweight design avoiding overfitting issues. These characteristics collectively ensure reliable performance across various educational informatization scenarios. establishing technical foundation for deploying personalized learning support systems. Notably, the model still exhibits recognition biases cross-cultural learners, pointing to directions for future research improvements.

5. Conclusions

This study systematically integrates multimodal data to develop an efficient and accurate personalized learning state recognition model, addressing the limitations of traditional methods in representing complex learning scenarios. Theoretical analysis and experimental validation demonstrate that the model can dynamically capture learners' cognitive and emotional changes, providing technical support for intelligent decision-making in adaptive learning systems. Future research could further explore cross-scenario transfer capabilities and privacy protection mechanisms to promote sustainable development in intelligent education

applications.

References

- [1] Xie Dingfeng, Zhou Anzhong, Li Jieqin. Research on Personalized Education Evaluation Empowered by Multimodal Data [J]. Journal of Hubei Open Vocational College, 2025,38(10):149-151.
- [2] Xue Yaofeng, Qiu Yisheng, Chen Zhan. System Framework for Multimodal Data Fusion and Its Educational Applications [J]. Basic Education, 2024,21(05):62-70.
- [3] Xie Dingfeng, Zhou Anzhong, Li Jieqin, et al. Research on Precision Intervention for Personalized Learning Based on Multimodal Data [J]. Computer Knowledge and Technology, 2024,20(16):98-100+104.
- [4] Jiang Jie, Yu Wenting, Wang Haiyan. Research on Student Learning Behaviors in Smart Classrooms Based on Multimodal Data [J]. China Education Informatization, 2024,30(04):107-117.
- [5] Zhang Lele and Gu Xiaoqing. A Classroom Teaching Behavior Analysis Model and Practical Framework Supported by Multimodal Data [J]. Open Education Research, 2022,28(06):101-110.
- [6] Hu Wenting. Research on Learning and Analysis Applications for Multimodal Data [J]. China New Communications, 2022,24(17):107-109.
- [7] Zhang Lizhao. Analysis of Learning Investment in Multimodal Data-Driven Contexts [D]. Central China Normal University, 2022.