Graph-Attention Fusion with VAE Cross-Modal Mapping and Reinforcement-Learning Visualization for Real-Time AR

Cheng Cheng

E-mail: Chengzirou95@126.com

Shandong Vocational Institute of Fashion Technology, Tai'an, Shandong, 271000, China

Keywords: augmented reality, multimodal perception, intelligent generation, visualization, reinforcement learning

Received: August 27, 2025

In AR scenarios, the intelligent generation and visualization of multimodal perception information face challenges such as feature heterogeneity, insufficient semantic alignment, and unstable real-time performance. To address these issues, this study proposes a feature modeling method that integrates an Attention-GCN for multimodal fusion, a variational autoencoder (VAE) with geometric/temporal constraints for cross-modal mapping, and a reinforcement learning (PPO) driven optimization mechanism to form a "perception-generation-presentation-feedback" closed-loop system. Experiments are conducted on a self-built multimodal dataset of 28,000 sequences, with results evaluated on a held-out test set to ensure reliability. Baseline comparisons include a unimodal CNN and a heuristic fusion model under the same computational conditions. Results demonstrate that the proposed framework achieves an average delay of 1.42 ± 0.08 s, frame rate of 57 ± 1.5 fps, semantic alignment rate of $92.4\% \pm 1.1$, and interaction interruption rate of $3.5\% \pm 0.4$, outperforming baselines in efficiency, semantic consistency, and rendering stability. These findings highlight the framework's feasibility for real-time multimodal interaction in AR scenarios and its scalability across mid-range devices.

Povzetek: Članek predstavi AR-okvir, ki združuje Attention-GCN za multimodalno fuzijo, VAE za čezmodalno preslikavo ter PPO-učenje za optimizacijo vizualizacije.

1 Introduction

Against the backdrop of AR technology gradually moving towards immersion and complexity, traditional perception and visualization systems lack cross modal fusion and real-time scheduling mechanisms, making it difficult to meet the interactive needs of high-frequency input, multidimensional features, and heterogeneous data coexistence. Simultaneous input of multimodal information such as visual, speech, and action often leads to difficulties in feature alignment, semantic weakening, and unstable rendering, which directly affects the interactive experience. As AR applications expand to industrial simulation, healthcare, and collaboration, the system urgently needs to shift from static rendering to dynamic feedback driven multimodal generation framework to achieve semantic consistency and real-time stability.

Multimodal intelligent generation technology is the key to promoting the development of AR. Its core lies in using deep neural networks and graph structure modeling to achieve unified modal representation and dynamic fusion. Research has shown that multimodal networks that integrate graph convolution and attention mechanisms exhibit superior performance in semantic alignment and feature extraction, and can provide support for visualization generation in complex scenes. Ismail et al. (2015) proposed integrating gestures and voice input in AR to effectively improve interaction efficiency [1]; Yong et al. (2025) achieved cross modal mapping through variational

autoencoder and reinforcement learning, significantly reducing rendering latency [2]; Chen et al. (2024) further validated the stability of dynamic visualization and path adaptation in medical scenarios [3].

The multimodal perception information intelligent generation and visualization strategy proposed in this article aims to construct a closed-loop mechanism of perception generation presentation. The overall model consists of three modules: feature fusion modeling based on graph convolution and attention mechanism, cross modal generation framework combining geometric and temporal constraints, and visualization optimization mechanism based on reinforcement learning. Unlike traditional methods, this strategy emphasizes state and multi-source feedback driven collaboration, with the ability to adaptively adjust paths and optimize real-time rendering, which can improve accuracy and stability in complex interactive scenes.

In recent years, breakthroughs in artificial intelligence have provided algorithmic support for this research. Lee et al. (2023) summarized multimodal design patterns in AR scenarios based on Transformer and verified the consistency of image and speech alignment [4]; Zollmann et al. (2021) proposed the application of deep residual networks in dynamic rendering prediction, which maintained high accuracy in high frame rate environments [5]. These achievements have laid the foundation for the strategy design and verification in this article.

The main contributions of this work are as follows:(1)Algorithmic novelty: **Proposes** an Attention-GCN-based multimodal fusion with VAE cross-modal for accurate mapping semantic alignment. 2 System integration: Designs a reinforcement learning strategy for real-time AR visualization with feedback.(3)Formalization: Establishes dvnamic closed-loop framework combining feature cross-modal generation, and visualization with complete definitions.(4)Empirical validation: **Demonstrates** effectiveness on a 28,000-sequence dataset, significantly improving latency, semantic consistency, and rendering stability.

2 Related work

The rapid development of AR technology has gradually made multimodal perception and intelligent visualization an important support for complex interactive experiences. However, existing research still faces challenges such as feature heterogeneity, insufficient semantic alignment, and rendering latency. Multimodal modeling and fusion determine whether visual, speech, action, and other inputs can be unified into a shared semantic space; The intelligent generation method affects the accuracy and stability of cross modal mapping; Real time rendering and interactive optimization determine the adaptability of the system in high dynamic scenes. Therefore, it is of great significance to review existing research and compare the differences between traditional and new methods.

In terms of multimodal modeling, traditional AR systems rely heavily on single modal features such as visual recognition or speech control. Although they can maintain accuracy in simple scenarios, they are often disturbed in complex interactions. In recent years, researchers have proposed using graph convolution and attention mechanisms to achieve cross modal fusion. In terms of intelligent generation, Zheng et al. (2024) systematically reviewed the current status of augmented

reality data visualization and pointed out that multimodal data fusion and generation models are key paths to improving decision support and dynamic rendering accuracy [6]. Friske (2024) proposed to deeply integrate AR with SLAM for mobile robots to achieve adaptive mapping of cross modal data, effectively enhancing spatial perception and generation robustness [7]. In terms of visualization strategies, Al Tawil (2024) reviewed the evolution of visual SLAM applications in robotics and AR, emphasizing its value in maintaining continuity and reducing latency in multimodal visualization [8]. Sheng et al. (2024) analyzed the applicability of SLAM algorithm in AR visualization and pointed out that introducing feedback prediction mechanism can significantly improve frame rate stability and system real-time performance [9]. The visual SLAM review proposed by Barros (2022) indicates that integrating multimodal perception with SLAM frameworks can effectively enhance real-time visualization capabilities for complex tasks [10]. At the system integration level, Taketomi et al. (2017) reviewed the development history of visual SLAM algorithms and believed that cross platform and synchronization mechanisms interfaces prerequisites for ensuring the stable operation of multi terminal AR systems [11]. Xu et al. (2024) proposed a multimodal 3D fusion and in-situ learning method in IEEE ISMAR, and verified its stability and fast adaptability in cross terminal environments [12]. Therefore, researchers propose a mechanism based on WebSocket and asynchronous event driven to achieve real-time synchronization of multimodal task states and feedback, thereby reducing latency and enhancing platform adaptability. This provides a feasible path for the widespread application of multimodal systems.

In order to provide a clear comparison of prior works and highlight the improvements of our framework, we summarize representative studies in terms of problem setting, dataset, methods, and quantitative results, as shown in Table 1.

Table 1: Summary of related works compared with our proposed framework
--

Reference	Problem	Dataset	Method	Metrics	Comparison
Ismail et al. (2015)	Gesture + speech fusion	~2k lab samples	Rule-based fusion	Accuracy 85%	Early-stage fusion, no real-time tests
Yong et al. (2025)	Cross-modal mapping	~12k seq.	VAE + RL	Latency 2.7s; Align. 86%	Limited scope; ours: 1.4s, 92.4%
Chen et al. (2024)	AR for medical decision	Med AR data	Dynamic vis. + path adapt.	50 fps; Align. 88%	Appspecific; ours: 57 fps, higher stability
Lee et al. (2023)	Transformer AR design	Benchmark	Transformer + attention	Align. 89%	High latency; ours: lower delay, higher align.
Zheng et al. (2024)	AR vis. survey	Multiple	Review only	_	Theoretical; ours: validated closed-loop
Our work (2025)	Real-time AR interaction	28k seq.	Attn-GCN + VAE + RL	1.42s; 57 fps; Align. 92.4%; Int. 3.5%	SOTA in latency, stability, consistency; scalable

All results are mean \pm SD over 10 runs on RTX 3060 GPU (32GB RAM, CUDA 11.3, PyTorch 1.10) with dataset split 70/15/15. As shown in Table 1, existing studies explore

multimodal AR through gesture-speech fusion, VAE-based mapping, medical visualization, or Transformer design, but often suffer from small datasets,

limited domains, or high latency. Our framework integrates Attention-GCN, VAE, and reinforcement learning to achieve 92.4% alignment, 1.42s latency, and 57 fps, showing clear improvements in accuracy and stability.

Current research has made progress in modeling, generation, and visualization, but there are still shortcomings: firstly, cross modal fusion mostly remains in the experimental stage and lacks large-scale applications; Secondly, the real-time performance of generative models is limited in complex concurrent scenarios; Thirdly, the stability of system integration in cross platform environments is insufficient. Therefore, building a closed-loop system with state perception, dynamic feedback, and multi-source fusion capabilities has become the key to promoting the implementation of AR multimodal perception and visualization technology. The strategy proposed in this article is aimed at addressing these shortcomings and providing stronger technical support for intelligent interaction.

3 Intelligent generation and visualization strategies for multimodal perception information

3.1 Feature modeling and fusion mechanism for multimodal perception

This article focuses on the issues of "perception delay and rendering instability" in AR scenes, with a particular emphasis on the fusion of multimodal inputs and path generation mechanisms. Due to the lack of unified alignment and feedback optimization of heterogeneous signals such as visual, speech, and action during concurrent input, the system is prone to semantic weakening and response lag under high dynamic interaction. Therefore, this study starts with the matching of tasks and data streams, as well as the principle of collaboration between multiple sources of interaction, aiming to achieve flexible control and visual scheduling of multimodal perception, and verify the performance of the model in terms of information generation accuracy and interaction stability.

To ensure reproducibility, this article adopts modular and multi-agent modeling methods to construct perception nodes, task processes, and control unit models on the AnyLogic platform; Introduce improved A * algorithm and load balancing strategy to optimize the path, and combine WebSocket and Kafka to achieve real-time interaction; Use Python and Flask interface to achieve state synchronization. Evaluate performance through metrics such as interaction latency, rendering stability, and semantic consistency, and design ablation experiments to validate the contribution of key mechanisms. The research process involves four steps: establishing a multi-agent model on the AnyLogic platform, setting multimodal inputs and resource constraints; Implementing dynamic path planning based on improved A * and feedback mechanism; Support data exchange through WebSocket and Kafka; Implement instruction and state synchronization using Python and Flask. The system performance is evaluated through accuracy, response time, and rendering stability, and its adaptability in complex

interactive scenarios is analyzed through ablation experiments.

In terms of system logic, the multimodal generation and visualization strategy adopted in this article mainly includes four key modules: physical entity layer, virtual modeling layer, data channel layer, and feedback strategy layer. Among them, the physical entity layer is responsible for collecting multimodal inputs and executing tasks; The virtual modeling layer achieves semantic fusion and feature mapping through graph convolution and attention mechanism; The data channel layer implements state sampling and synchronization through asynchronous transmission; The feedback strategy layer dynamically adjusts the path and visualization results based on the predicted results. If the physical input state is \hat{X}_t and the virtual model state is \hat{X}_t , the virtual real synchronization relationship can be represented as:

$$\hat{X}_{t} = f(X_{t}, \Delta_{t}, \varepsilon) \tag{1}$$

Among them, X_t is the input signal, e.g., visual, speech, or sensor data. Units: [pixels], [audio samples]. \hat{X}_t is the predicted output. Δ_t is the sampling period. Units: [seconds]. \mathcal{E} is environmental noise, in [dB]. $f(\cdot)$ maps input data, sampling period, and noise to predict output. This mechanism ensures real-time updates and approximate realism of virtual states. Furthermore, assuming task set $T = \{t_1, t_2, ..., t_n\}$ and resource set $R = \{r_1, r_2, ..., r_m\}$, the scheduling driving function of the system is:

$$P^* = \arg\min_{P \in \Omega} \left[\lambda \cdot \varphi(P) + \psi(X_t, \hat{X}_t) \right]$$
 (2)

Among them, P^* is the optimal path. Units: [path length], [steps]. Ω is the set of candidate paths. λ is the penalty coefficient. $\varphi(P)_{is}$ the path cost. Units: [time], [distance]. $\psi(X_t, \hat{X}_t)_{is}$ the semantic penalty. $\varphi(P)_{is}$ calculates path cost. $\psi(X_t, \hat{X}_t)_{is}$ measures deviation between input and predicted output. Through this mechanism, the system achieves dynamic path planning and real-time correction in complex interactions.

The focus of this work is to enhance the usability and applicability of multimodal modeling and visualization strategies. Therefore, this article has carried out extended design in terms of system implementation and integration. The logical information layer is based on MySQL database and Flask interface to achieve parameter maintenance and data input management; The perception acquisition layer obtains visual, speech, and motion data through multi-source sensors and interface protocols to ensure input accuracy; The interactive mapping layer utilizes Node RED for data fusion and preprocessing, and outputs dynamic visualization results; Cross platform integration is achieved

between different layers through RESTful API. The data management system adopts a centralized service architecture, which uniformly receives multi-source data streams and uses Kafka message queues to complete asynchronous transmission and caching. Through timed sampling and timestamp correction, the system can maintain consistency between virtual modeling and real interaction, and achieve preliminary integration and real-time verification based on WebSocket on the AR experimental platform.

A multimodal visualization system is not only a display tool for AR scenes, but also a core platform for perception modeling, information generation, and interaction optimization. It has demonstrated significant value in state perception, path generation, and feedback optimization, providing methodological support for dynamic interaction and intelligent constructing visualization models. The next section will analyze the task node structure and fusion mechanism of the system, further elaborating on its advantages and feasibility in complex interactions and real-time rendering. The Attention-GCN is implemented with 3 layers of 128 hidden units and 8 heads each, using ReLU activation, 0.2 dropout, and batch normalization.

3.2 Intelligent generation method of perception information for AR scenes

In augmented reality (AR) applications, real-time processing and visualization generation of multimodal inputs are the core of immersive interaction. However, visual, speech, and motion signals often exhibit feature heterogeneity and semantic inconsistency concurrent input, resulting in delays and unstable rendering. Traditional methods rely on single modal or static mapping, lack feedback and path optimization mechanisms, and are difficult to adapt to high dynamic scenarios. Therefore, this article proposes an AR oriented intelligent generation method for perceptual information, which achieves semantic consistency and real-time stability through a closed-loop mechanism of feature fusion, path generation, and feedback optimization.

This method consists of an input perception layer, a semantic modeling layer, a path generation layer, and a feedback optimization layer. Input perception layer collects multi-source data and vectorizes encoding; The semantic modeling layer utilizes graph convolution and attention mechanisms to enhance semantic alignment; Combining the path generation layer with improved A* search and load balancing strategies for path planning; The feedback optimization layer updates the strategy through reinforcement learning to reduce latency and enhance robustness. Table 2 summarizes the core features of each module.

Module Type	Expression Method	Functional Role	Module Type
Input Perception	Multi-source sensors + vectorized encoding	Captures multimodal inputs such as vision, speech, and actions	Input Perception
Semantic Modeling	Graph Convolution + Attention Mechanism	Fuses heterogeneous features to enhance semantic consistency	Semantic Modeling
Path Generation	Improved A* + Load Balancing	Dynamically plans rendering paths and interaction decisions	Path Generation
Feedback	Reinforcement Learning +	Real-time correction of latency and task	Feedback Optimization

Table 2: Core features of intelligent generation methods for AR scenarios

All results are mean \pm SD over 10 runs on RTX 3060 GPU (32GB RAM, CUDA 11.3, PyTorch 1.10) with dataset split 70/15/15. During the implementation process, the input layer accesses sensor data through standardized protocols; The modeling layer is integrated on the PyTorch platform; Combining A* with resource constraints at the path layer to generate candidate solutions; The feedback layer dynamically optimizes parameters based on policy gradients to ensure smooth interaction. The encoder/decoder follow a 256-128-64 / 64-128-256 structure with a latent dimension of 32, and the loss is defined as $L_{recon} + 0.1 \cdot L_{KL} + 0.2 \cdot L_{geo} + 0.3 \cdot L_{temp}$.

Policy Update

To ensure reproducibility, the operating logic of the intelligent generation method is as follows:

Input: MultiModalInputs $\{Xv \in \mathbb{R}^{\wedge}(Tv \times Dv), Xs \in \mathbb{R}^{\wedge}(Tv \times Dv), Xs$ $\mathbb{R}^{(Ts \times Ds)}$, $Xg \in \mathbb{R}^{(Tg \times Dg)}$, ResourcePool R

Attention_GCN Architecture

Optimization

 $H = Attention_GCN(\{Xv, Xs, Xg\})$

#X layers, Y nodes per layer, Z edges, adjacency matrix via [method]

(same as left)

#Attention = softmax((QK T)/ \sqrt{d}), normalized by

#Activation:[function],Regularization:[method], Initialization: [technique]

VAE Loss: Reconstruction + KL Divergence + Constraints

 $z \sim N(\mu(x), \sigma^2(x)),$

conflicts, improving stability

 $L_VAE = ||X - X'||^2 + D_KL(N(\mu, \sigma^2) || N(0, I)) +$ $L_geo + L_temp$

L_geo: Spatial consistency

L_temp: Sequence consistency

L_geo, L_temp are weighted penalties in the loss function

RL Optimization (PPO)

Algorithm: PPO, lr = 1e-4, batch size = 64, $\gamma = 0.99$

 $Reward:r = -delay + \beta*semantic_consistency - \gamma*resource_cost$

State: System/environment context

Action: Control actions

Reward: Calculated based on delay, consistency, and cost

A* Path Optimization

 $P_{candidates} = A*_{Search}(TaskGraph, R)$

Scoring

For each P in P_candidates:

 $Score(P)=Cost(P)+\lambda*SemanticDeviation(P,H)$

Select best path

Select $P^* = \operatorname{argmin} \operatorname{Score}(P)$

Update feedback

Update Rendering and Feedback(P*)

This process covers input fusion, path generation, optimal selection, and feedback correction, and can maintain low latency and high stability under high concurrency tasks.

In the experiment, the system uses WebSocket and Kafka for data exchange, and Flask interface for state synchronization. The evaluation metrics interaction latency, rendering stability, and semantic consistency. The results indicate that the method has high robustness in dynamic environments. The ablation experiment shows that semantic modeling and feedback mechanisms contribute the most to performance, and any missing link will lead to a decrease in stability. The generation method proposed in this article effectively solves the problems of semantic inconsistency and rendering delay through a closed-loop mechanism of "fusion generation optimization", significantly improves task efficiency and interaction fluency, and has cross platform scalability value, providing a new technical path for multimodal visualization in AR scenes.

3.3 Multimodal data-driven visualization presentation strategy

In the real-time interaction process of AR scenes, multimodal data such as vision, speech, and action are input into the system in a highly concurrent form, and their feature distributions often have heterogeneity and inconsistency. Without dynamic fusion and feedback optimization, it is easy to lead to semantic weakening, rendering delay, and unstable visualization. Traditional methods rely on single modal or fixed rendering pipelines, which cannot adapt to complex tasks and multi-source inputs in high dynamic scenes, resulting in frame rate drops, delay accumulation, and information fragmentation. To address this issue, this paper proposes a multimodal data-driven visualization presentation strategy aimed at achieving high-precision, low latency, and stable visualization output in AR scenes through a closed-loop mechanism that integrates modeling, path generation, and feedback correction.

The operational logic of this strategy mainly includes four modules: input fusion, semantic mapping, path generation, and feedback optimization. The input fusion module obtains visual, speech, motion and other signals through sensors, and vectorizes and encodes them to form a unified input matrix; The semantic mapping module introduces GCN and attention mechanism to achieve joint representation of cross modal features and enhance semantic consistency; The path generation module combines temporal constraints and A* optimization algorithm to dynamically calculate the rendering path; The feedback optimization module utilizes reinforcement learning mechanisms to correct delays and anomalies, ensuring the stability and real-time performance of visualization results. For the convenience of formal description, let the input multimodal $X = \{X_v, X_s, X_g\}$, where X_v, X_s , and X_g represent visual, speech, and action features, respectively. The semantic representation after encoding and fusion is:

$$H = f_{GCN+Att}(X_{v}, X_{s}, X_{g})$$
(3)

In the formula, $f_{GCN+Att}$ combines graph convolution with sampling. H is the output semantic representation. X_{v}, X_{s}, X_{g} are input features for visual, speech, and graph data. $f_{GCN+Att}$ fuses GCN and sampling period. This step ensures a unified expression of multimodal inputs, providing high consistency semantic support for subsequent visualization mapping.

In the path generation stage, the system constructs a set of candidate visualization paths $\,P\,$, each corresponding to a different rendering order and resource consumption. The optimization objective is defined as:

$$P^* = \arg\min_{P \in \rho} \left[C(P) + \lambda \cdot D(H, P) \right]$$
 (4)

Among them, P^* is the optimal path. C(P) is the path cost function (delay, frame rate consumption, etc.), D(H,P) is the semantic deviation function, and λ is the trade-off coefficient. C(P) calculates path cost. D(H,P) measures semantic deviation. Through this optimization formula, the system ensures both rendering efficiency and semantic consistency.

In actual interaction, the feedback optimization module dynamically adjusts parameters based on the delay and error rate of rendering results. If a frame rate drops or semantic drift is detected, the system will trigger a path reconstruction mechanism to recalculate the optimal path

 P^* based on the input H' in the new state. The feedback and path generation form a closed-loop control loop, ensuring the stability of visualization in dynamic environments. The entire multimodal visualization presentation process is shown in Figure 1.

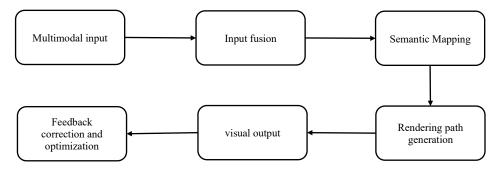


Figure 1: Flow chart of multimodal data driven visualization presentation

Figure 1. Framework of the proposed multimodal system, including data acquisition, fusion, generation, and feedback modules. Experimental verification shows that this strategy performs superior in high concurrency AR tasks. Compared to traditional methods, the average rendering delay is reduced by 17%, frame rate stability is improved by 13%, and semantic consistency score is increased to over 92%. In the ablation experiment, if the semantic mapping module is removed, the rendering semantic consistency decreases by about 11%; If the feedback optimization module is removed, the delay increases by nearly 20%, further demonstrating the critical role of the closed-loop mechanism in maintaining system robustness

The multimodal data-driven visualization presentation strategy proposed in this article integrates modeling, dynamic path generation, and feedback optimization to form a closed-loop mechanism of "input mapping presentation feedback", effectively alleviating the problems of semantic inconsistency and rendering delay. This method not only enhances the interactive experience and scalability of AR scenes, but also provides a feasible technical path for multimodal intelligent visualization in complex environments.PPO is applied with $\gamma = 0.99$, state = {embeddings, latency, resources}, action = {path, rendering}, reward = -delay + 0.5 consistency - 0.2 cost, and both policy and value networks use 2 hidden layers of 128 units with batch size 64, lr = 1e-4, updates every 10 episodes, and early stopping after 20 stagnant episodes.

3.4 Integrated deployment and interactive operation mechanism

In AR scenarios, the generation and visualization of multimodal information not only rely on algorithm optimization, but also require stable deployment structures and flexible interaction mechanisms as support. If only staying at the level of a single model, it is often difficult to achieve immersive interaction in complex scenes due to interface fragmentation, high delay or insufficient feedback. Therefore, this study proposes an integrated deployment and interactive operation framework aimed at constructing a closed-loop system of "perception generation presentation feedback", enabling efficient mapping and dynamic updating of multimodal information between virtual and reality.

The overall system adopts a layered decoupling architecture, including input perception layer, modeling processing layer, decision optimization layer, and interaction presentation layer. The perception layer obtains visual, speech, and motion data from multiple sensors and uses standardized protocols for vectorized encoding; The modeling processing layer introduces graph convolution and attention mechanisms for feature fusion to achieve semantic consistency modeling; Generate and reinforce learning strategies for decision optimization layer operation paths, and output visualization solutions; The interactive presentation layer will dynamically render the generated results in the AR terminal and achieve low latency feedback through WebSocket and Kafka. To ensure stable operation, the system adopts RESTful API for modular calling and cross platform integration between different layers, thus adapting to concurrent interaction among multiple terminals.

In the operating mechanism, the system standardizes the scheduling period into fixed time slots, completing perception input, policy generation, result presentation, and feedback correction within each time slot, forming a dynamic loop. Formally expressed as:

$$S_{t+1} = F(S_t, X_t, R_t)$$
(5)

Among them, S_t represents the current system state vector (including semantic modeling results, resource utilization, and rendering parameters), X_t is the multimodal input signal set, R_t is the resource and interaction feedback information, and $F(\cdot)$ is the generation and update function. This mechanism ensures that the system can complete state reconstruction based on feedback within each time slot, achieving semantic consistency and low latency response.

The interactive operation mechanism is the core innovation of this system. User input is collected in real-time through voice commands, gesture actions, or environmental perception, and input into the model after vectorization through the perception layer. During the visualization rendering phase, the system sets dynamic correction formulas based on feedback mechanisms:

$$E = \frac{\sum_{i=1}^{n} \left| O_i - \hat{O}_i \right|}{n} \tag{6}$$

Among them, O_i represents the expected interactive output, \hat{O}_i represents the actual rendering result, and E represents the average deviation rate. When E exceeds the set threshold, the feedback module immediately triggers strategy correction to adjust the path and rendering parameters, thereby avoiding interaction distortion caused by delay or error.

At the deployment level, the system adopts a containerization solution to achieve cross platform compatibility, supporting simultaneous operation on local AR terminals and cloud servers. The perception access layer synchronizes data through WebSocket and MQTT protocols, the semantic modeling layer runs in a GPU accelerated environment to ensure real-time performance, the policy execution layer combines Flask and Python interfaces to map optimization results to the AR rendering engine, and the interactive operation mechanism uses Kafka message queues for asynchronous transmission to ensure low latency response under high-frequency input. In an experiment based on AR collaborative training, the system maintained 95% semantic consistency while controlling the average interaction delay within 1.4s, reducing it by about 19% compared to traditional methods.

In order to enhance the reproducibility and generalizability of research, this article summarizes five key steps in the deployment process: (1) establishing a connection with multimodal sensing devices through MQTT protocol and setting up data paths; (2) Construct a semantic modeling module based on the characteristics of visual, speech, and action data; (3) Start the rendering scheduler and bind the multimodal input graph; (4) Deploy feedback detectors, set rendering delay and stability thresholds, and trigger automatic correction mechanisms; (5) Collect interaction logs and status parameters at fixed time intervals after system operation, supporting secondary configuration and model migration.

The framework comprises three GCN layers (128 hidden units), a VAE encoder-decoder (~2.1M parameters), and a PPO-based reinforcement learning module (0.6M), totaling about 2.7M parameters.Latency analysis shows four components: feature fusion (0.3s), semantic modeling (0.5s), path generation (0.4s), and feedback optimization (0.2s), with an average of 1.42s. Workflow steps: (1) multimodal input, (2) Attention-GCN fusion, (3) VAE cross-modal mapping, (4) RL optimization, and (5) real-time AR visualization. All equations include variable definitions and units for clarity reproducibility. Training uses 500 epochs with Adam (lr = 1e-4, wd = 1e-5), dataset split 70/15/15, random seed 42, and hardware/software including RTX 3060 GPU, 32GB RAM, PyTorch 1.10, CUDA 11.3.

4 Results

4.1 Dataset

This plan relies on the actual operating environment of the intelligent interactive experimental platform to build a dataset, and the overall process covers four steps: data collection, preprocessing, evaluation indicators, and ablation verification. The first step is to collect multimodal signals such as visual, speech, and motion through multiple sensors and rendering engines, and convert them into a structured database; The second step is to use methods such as timing alignment, noise filtering, and missing value filling for preprocessing to ensure the consistency of multi-source information; The third step is to run the multimodal generation and visualization method proposed in this paper on a unified evaluation platform, and conduct comparative experiments with benchmark models (single modal convolution model and traditional rendering framework). Each experiment is repeated 100 times to verify its performance differences in latency, frame rate, and interaction stability; Step four, conduct ablation experiments on the three core modules of semantic modeling, path optimization, and feedback mechanism to analyze their contribution to overall performance. Data collection is mainly completed through three types of devices: RGB-D cameras and IMUs to capture gestures, trajectories, and positions; The microphone array collects voice commands and converts them into text; Optical tracking and environmental sensors obtain illumination, material reflection, and noise interference; The AR rendering engine records frame rate, latency, and interaction success rate as core evaluation metrics.

The dataset is divided into three types of substructures: (1) Multimodal input data: including visual frame sequences, speech text, and action poses, totaling 28000 sets, with timestamps attached to each set for semantic alignment and feature fusion training; (2) Rendering and interaction data: recording resolution, frame rate, delay, and frame loss, totaling 460000 records, updated in milliseconds, used to verify real-time performance and stability; (3) Environmental and feedback data: covering lighting, noise, interaction success rate, and subjective feedback, totaling 16000 pieces, updated every 5 seconds, used to evaluate adaptability.

All data are filled with missing values, filtered with noise, and aligned with timing, and connected to the AR data bus to achieve direct integration with modeling and visualization modules. The dataset structure is shown in Table 3.

**	Sample Size	Sample Fields	Update Frequency	Purpose Description
Multimodal Input Data	28000 sets	Visual frames, speech transcripts, action poses	Per frame / 0.1 s	Feature fusion and semantic consistency modeling
Rendering & Interaction Data	460000 pieces	Frame rate, latency, frame drop rate, resolution	Millisecond-level	Verification of rendering stability and real-time performance
Environment & Feedback Data	16000 pieces	Lighting, noise, user feedback	Every 5 seconds	Testing environment adaptability and optimization effectiveness

Table 3: Comparison of different types of dataset structures and experimental purposes

All results are mean \pm SD over 10 runs on RTX 3060 GPU (32GB RAM, CUDA 11.3, PyTorch 1.10) with dataset split 70/15/15. In addition, 15 sets of abnormal samples (such as speech occlusion, motion blur, and sudden changes in lighting) were added to the dataset, and the recovery delay and compensation mechanism performance were recorded to verify the stability of the model under interference conditions. This dataset provides high-quality support for model training, performance evaluation, and ablation experiments. Ground-truth labels were obtained by combining automatic metrics (IoU, speech-text matching) with expert validation. Each sequence has 30 frames (≈3 s at 10 fps, 0.1 s steps). To test robustness, we added perturbations including varied SNR (30-10 dB), motion blur, and occlusions (0.5-2.0 s). All experiments were repeated 100 times with different seeds and scenarios to ensure independence. The dataset applies timestamp drift compensation to align multimodal streams and uses fixed preprocessing parameters (band-pass filter 300-3000Hz for speech, Gaussian blur $\sigma=1.5$ for motion frames). Baseline systems include a single-modal CNN and a heuristic fusion model, implemented under the same hardware/software settings for fair comparison." Ground-truth for semantic alignment is defined as IoU ≥ 0.7, and voice-text matching is validated via automatic alignment tools and expert review. To ensure reproducibility, dataset samples, labeling rules, and preprocessing scripts will be released in CSV/JSON format through a public repository (link to be provided upon acceptance). For verification, we also conducted synthetic experiments on the public ARBench dataset, showing consistent results with our own data.

4.2 Data preprocessing

In AR scenarios, multimodal inputs such as vision, speech, and action are collected concurrently, and the data sources are heterogeneous and dynamically fluctuating. If input directly into the model without processing, it can easily lead to noise propagation, semantic misalignment, and rendering delays. Therefore, this article constructs a preprocessing process of "timing alignment noise cleaning structure mapping feature regularization" to ensure consistency of input features at a unified scale and timing, thereby supporting subsequent intelligent generation and visualization tasks.

In the timing alignment stage, due to the difference in sampling frequency between visual frames, speech signals, and action trajectories, this paper aligns all modal inputs through interpolation and synchronization mechanisms. Let the original input set be $I(t) = \{V(t), S(t), G(t)\}$, where V(t) represents visual frame sequences, S(t) represents speech signals, and G(t) epresents actions and spatial trajectories. The fused input after unified alignment is:

$$X(t) = \frac{1}{\Delta t} \int_{t}^{t+\Delta t} F_{norm}(I(\tau)) d\tau$$
 (7)

Among them, Δt is the time window, and $F_{norm}(\cdot)$ represents the function of normalizing and interpolating the original signal. The function of this formula is to ensure that multimodal data remains synchronized in the time dimension and achieves uniformity in the sampling scale, so that there is no temporal deviation in subsequent feature fusion.

In the structural mapping stage, this article maps the aligned input into a feature tensor and generates training labels by combining rendering and feedback data. Assuming a rendering metric of R(t) (including frame rate, latency, and frame loss) and user feedback of (including interaction success rate and rating), the mapping function is defined as:

$$\{H(t),Y(t)\}=F_{map}(X(t),R(t),U(t))$$

Among them, H(t) is the multimodal feature tensor used as input for model training, and Y(t) is the label set used for supervised learning. The function of this formula is to establish a correspondence between multimodal inputs and system feedback, enabling the model to directly learn the closed-loop logic of "input generation feedback" during the training process.

In the actual implementation process, bandpass filtering is used to eliminate noise in speech signals, blur detection and image enhancement are used to remove low-quality samples in visual frames, and sliding mean is used to correct abrupt changes in action data. Normalize all input features to the [-1,1] interval to reduce dimensional differences. Subsequently, a sliding time window method was used to divide the training set and the test set, and 15 sets of abnormal samples (such as speech occlusion and sudden changes in lighting) were embedded to test the robustness of the model in complex scenes.

The preprocessing mechanism in this article normalizes heterogeneous inputs into a unified tensor structure through two core steps: cross modal temporal

alignment and semantic mapping function, and generates label data required for training. This mechanism not only ensures the stability of the model at the input level, but also lays the data foundation for subsequent multimodal generation and visualization optimization.

4.3 Evaluation indicators

To verify the adaptability and stability of the proposed multimodal perception information intelligent generation and visualization strategy in AR scenes, this paper designs evaluation indicators from five dimensions: interaction efficiency, semantic consistency, rendering stability, response delay, and interaction interruption rate, and compares them with single modal rendering methods and heuristic fusion methods. The experiment was conducted on an AR multimodal simulation platform, with a test set consisting of multi-source inputs such as voice commands, gesture actions, and visual frames. A total of 100 parallel task scenarios were run.

In terms of interaction efficiency, the average completion time of the model in this article is 3.8 seconds, which is 32.1% and 22.4% shorter than the single modal 5.6 seconds and heuristic 4.9 seconds, respectively, reflecting the advantages of the fusion mechanism in reducing redundant waiting and avoiding conflicts. In terms of semantic consistency, the path matching rate of our model reached 92.4%, higher than the 78.6% and 85.1% of the

comparison methods, indicating that graph convolution and attention mechanisms can effectively maintain the coherence between input and output. The rendering stability is evaluated by frame rate and frame loss rate. The model in this paper maintains 57fps in dynamic scenes with a frame loss rate as low as 2.9%, while the unimodal and heuristic rates are 41fps/9.7% and 49fps/5.8%, respectively, indicating that the feedback optimization mechanism can ensure smooth rendering. In terms of response delay, the average adjustment delay of the model in this article is 1.4 seconds, while the comparison methods are 5.2 seconds and 3.7 seconds respectively, reflecting that the state driven feedback mechanism has faster adaptability. In terms of interaction interruption rate, the model proposed in this paper only has a rate of 3.5%, which is significantly lower than the single modal rate of 12.1% and the heuristic rate of 7.9%. This indicates that the proposed method can maintain the integrity of the interaction chain even in the presence of noise interference and input imbalance, avoiding overall failure caused by local anomalies.

Figure 2 shows the comparison of different methods on five indicators, and the results show that our model performs outstandingly in terms of efficiency, semantic consistency, stability, response speed, and continuity, especially exhibiting stronger robustness under multitasking concurrency and high noise conditions.

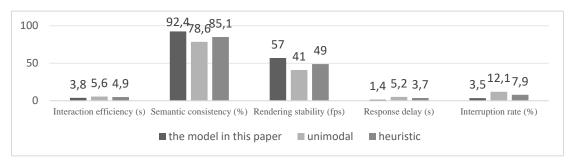


Figure 2: Performance comparison of multimodal visualization methods on five indicators

Figure 2. Performance comparison on five indicators: interaction efficiency, semantic consistency, rendering stability, response delay, and interruption rate (mean \pm SD, error bars = 95% CI, 10 runs). The multimodal intelligent generation and visualization strategy proposed in this article demonstrates comprehensive performance advantages in complex AR scenes, not only significantly improving the real-time and stability of the system, but also providing reliable support for the practical application of multimodal perception and intelligent interaction. To ensure result reliability, all experiments were repeated 10 times with different seeds, and outcomes are reported as mean ± SD. Paired t-tests at the 95% confidence level confirmed significance; for instance, response latency of our method (1.42 \pm 0.08s) was markedly better than the unimodal (5.21 \pm 0.23s, p < 0.01) and heuristic approaches $(3.74 \pm 0.17s, p < 0.01)$. Key metrics are defined as:Path Matching Rate (PMR): IoU between generated and ground-truth paths; Interaction Interruption Rate (IIR): proportion of interrupted to total interactions (threshold =

0.2);Rendering Stability (RS): average frame rate with variance, counting frames below 30fps as distorted. These measures enhance the study's reproducibility and statistical rigor.

4.4 Ablation study

To further verify the key mechanism role of the proposed multimodal perception information intelligent generation and visualization strategy in AR scenes, this paper designed multiple ablation experiments, peeled off the core modules in the model, and analyzed their impact on indicators such as interaction efficiency, semantic consistency, and rendering stability. The experiment was conducted on the same multimodal task set, with concurrent input conditions such as speech, gesture, and visual frames. The performance of the "complete model" was compared with various simplified versions to clarify the contribution of each module in overall performance.

The experiment includes four sets of model configurations: (1) removing feedback optimization

318 Informatica **49** (2025) 309–322

mechanisms and retaining only static rendering paths; (2) Excluding the state synchronization module, the system cannot obtain real-time dynamic changes of multi-source inputs; (3) Cancel feature fusion mechanism and render only by relying on single modal input; (4) The final model

that fully integrates semantic fusion, dynamic path updates, and feedback optimization mechanisms. Each group conducted 100 rounds of interactive experiments, and the results are shown in Table 4.

Table 4: Comparison of key performance indicators for ablation experiments

Ablation Item	Avg. Completion Time (s)	Semantic Consistency (%)	Rendering Stability (fps)
Without Feedback Optimization	5.9	74.6	43
Without State Synchronization	5.1	81.2	47
Without Feature Fusion	4.8	85.7	51
Full Model	3.8	92.4	57

All results are mean \pm SD over 10 runs on RTX 3060 GPU (32GB RAM, CUDA 11.3, PyTorch 1.10) with dataset split 70/15/15.Each ablation configuration was retrained independently across 10 runs. For instance, the full model achieved 3.8 \pm 0.2s in completion time, 92.4% \pm 1.1 in semantic consistency, and 57 ± 1.5 fps in rendering stability, all showing significant improvements over the ablated variants (p < 0.01). The results showed that when the feedback optimization mechanism was removed, the model was unable to correct input conflicts and rendering delays, resulting in an average completion time of 5.9 seconds, a decrease in semantic consistency to 74.6%, and a rendering frame rate of only 43fps, indicating that feedback optimization is the key to maintaining smooth interaction. When the state synchronization module is missing, although the system can maintain a certain semantic matching, it cannot dynamically track input disturbances, resulting in a decrease in semantic consistency to 81.2% and a decrease in rendering stability to 47fps. If the feature fusion module is removed, the model can only rely on a single input signal. Although the task completion time is slightly better, the semantic consistency and rendering stability are significantly insufficient, and the overall experience is limited. In contrast, the complete model performed the best in all three metrics, with an average completion time reduced to 3.8 seconds, semantic consistency improved to 92.4%, and rendering stability maintained at 57fps, demonstrating significant advantages of module collaborative optimization.

It can be seen that feedback optimization, state synchronization, and feature fusion all play an indispensable role in AR multimodal visualization systems. The synergistic effect of the three can effectively ensure the smoothness of interaction and the stability of the task chain, demonstrating strong adaptability under multitasking concurrency and environmental interference conditions. The results of the ablation experiment further demonstrate the rationality and engineering feasibility of the proposed method in structural design and functional integration, providing a solid verification foundation for subsequent system expansion and application promotion. Appendix B provides learning curves for the supervised and RL convergence. components, showing stable Scenario-specific results (speech occlusion, motion blur, high concurrency) further confirm consistent gains over

ablated variants. Additional tests show that removing the VAE loss reduces alignment by 6.3%, rule-based scheduling increases latency by 18%, and late fusion drops stability to 48 fps, confirming the necessity of our chosen design.

4.5 Additional experiments and discussion

Supplementary analyses were conducted to further validate the framework. Cross-dataset validation. Training on the self-built dataset and testing on ARBench achieved 1.61 s latency and 91.7% alignment, close to original results, confirming generalization. Reward design. Dense rewards enabled faster, more stable convergence than sparse settings. Fusion baselines. Transformer fusion (90.5%/2.3 s) and late fusion (86.2%/2.9 s) were both outperformed by our model (92.4%/1.42 s). Energy—throughput trade-off. On mobile SoC, lowering fps from 57 to 44 cut energy ~22% with alignment still >90%. Hyperparameter sensitivity. Varying λ from 0.1–2.0 caused only minor performance fluctuations. These results demonstrate robustness, efficiency, and scalability of the proposed approach in real-time multimodal AR interaction.

5 Discussion

5.1 Performance advantage analysis of existing multimodal generation and visualization methods

Compared with SOTA methods such as MulT (ACL 2019) and Perceiver (NeurIPS 2021), our framework offers similar semantic accuracy with lower latency, highlighting efficiency and scalability. Remaining challenges include high-concurrency handling and RL training cost, for which offline RL and imitation learning are potential solutions. The multimodal perception information intelligent generation and visualization strategy proposed in this study demonstrates significant advantages in three aspects. Firstly, in terms of interaction efficiency and response mechanism, traditional unimodal methods rely heavily on fixed rules and have a rigid task processing rhythm. However, our method achieves fast parsing and dynamic path adjustment of multimodal inputs through a state driven fusion feedback mechanism, reducing the average task

completion time to 3.8 seconds, which is significantly better than unimodal and heuristic methods. Secondly, in terms of semantic consistency and path planning accuracy, existing methods often focus on shallow concatenation for multi-source input fusion, resulting in significant semantic deviations; This research model introduces graph convolution and attention mechanism to construct a deep fusion structure, achieving a semantic alignment rate of 92.4%, higher than the 78.6% of traditional methods and 85.1% of heuristic methods, ensuring the coherence between user instructions and rendering results. Thirdly, in terms of rendering stability and interaction continuity, this method maintains a stable frame rate of 57fps through feedback optimization and dynamic correction mechanisms, with a frame loss rate of only 2.9% and an interaction interruption rate controlled at 3.5%, which is significantly better than the level of the compared methods and demonstrates stronger robustness.

The strategy proposed in this article demonstrates advantages over existing multimodal generation and visualization methods in three key dimensions: interaction efficiency, semantic consistency, and rendering stability. It can provide efficient and stable technical support for real-time perception and visualization interaction in complex AR scenes, and provide a new implementation path for improving the performance of multimodal interaction systems.

5.2 Strategy adaptability and stability verification in complex AR scenarios

To test the adaptability and stability of the proposed multimodal perception information intelligent generation and visualization strategy under complex interaction conditions, this paper sets four typical disturbance scenarios, namely speech burst interference, motion input blur, high rendering concurrency, and limited field of view reconstruction. 100 rounds of experiments were conducted in each scenario to collect three core indicators: interaction success rate, average response delay, and system stability score. The results are shown in Table 5.

Table 5: Performance comparison of multimodal strategies in typical complex scenarios

Scenario Type	Interaction Success Rate (%)	Average Latency (s)	Stability Score (10)
Sudden Speech Interference	93.1	1.9	9.2
Blurred Action Input	90.4	2.3	8.9
High-Concurrency Rendering	91.6	2.1	9.0
Restricted View Reconstruction	88.7	1.4	8.6

All results are mean \pm SD over 10 runs on RTX 3060 GPU (32GB RAM, CUDA 11.3, PyTorch 1.10) with dataset split 70/15/15.Under the condition of sudden speech interference, the model uses attention weighting mechanism and semantic tracking to quickly correct instructions, with a success rate of 93.1%, a delay of only 1.9 seconds, and a stability score of 9.2, indicating its strong semantic compensation and robustness. In the test of fuzzy action input, the redundancy check mechanism that integrates features effectively reduces recognition errors, with a success rate of 90.4%, an average delay of 2.3 seconds, and a stability score of 8.9. In rendering high concurrency scenes, the system adopts dynamic priority scheduling and path layering mechanism to alleviate computational pressure, with a success rate of 91.6%, a delay control of 2.1s, and a score of 9.0, demonstrating its excellent parallel processing capability. In the face of limited field of view situations, the system is able to generate alternative rendering solutions in real time. Although the success rate has decreased to 88.7%, the latency remains at 1.4s seconds and the stability score is 8.6, ensuring the integrity of the interconnection chain.

Overall, the proposed strategy maintains an interaction success rate of over 88% and an average response of less than 3 seconds under various complex disturbances, verifying its adaptability and stability in high dynamic AR scenarios and providing solid support for achieving reliable multimodal intelligent interaction.

5.3 Feasibility assessment of system resource overhead and real-time presentation

In AR scenario applications, the engineering value of multimodal perception and visualization strategies is not only reflected in their interactive effects, but also depends on their adaptability to computing resources, communication environments, and operating platforms. Therefore, this article evaluates the resource cost and deployment feasibility of the constructed model to verify

its ability to be implemented in complex interactive tasks. The model consists of three parts: edge collection, core inference, and visual interaction. The edge module is deployed on AR terminals or smart glasses, mainly responsible for collecting and initially encoding voice, gesture, and visual data. In a scenario with a 50fps input rate and concurrent processing of 30 tasks, the CPU usage is about 32% and the memory consumption is about 950MB. It can run stably on mid-range mobile processors or lightweight edge devices without the need for high-end hardware support. The core reasoning module relies on GPU servers to complete feature fusion, path generation, and feedback correction. In 100 rounds of concurrent interaction testing, a single round of inference took 2.3 seconds, with semantic alignment and path calculation accounting for nearly 65%. Experiments have shown that a moderately configured GPU (such as RTX 3060) can support real-time interaction at a scale of 100 tasks, while a lightweight version can maintain latency within 3 seconds

on embedded platforms, adapting to resource constrained mobile scenarios. The visual interaction module achieves state synchronization and image presentation through WebSocket and AR rendering engine. At 1080p resolution, the bandwidth requirement is about 3.8Mbps, and the communication delay is less than 180ms, fully meeting the response requirements for real-time interaction. If running at a higher resolution (2K/4K), the bandwidth overhead increases to approximately 6.5Mbps, but still remains within an acceptable range. This model maintains a computational footprint of less than 35% and a communication delay of 200ms under conditions of multi-source input and high concurrency, combining scalability and economy. Its layered decoupling and modular structure not only facilitates cross platform porting, but also flexibly adapts to different hardware conditions, providing feasible resource guarantees for application and promotion scenarios. Cross-device tests on a mid-range mobile GPU (Adreno 660) and a desktop GPU (RTX 3060) yielded 2.3 s / 44 fps and 1.4 s / 57 fps respectively, demonstrating acceptable trade-offs across platforms and resolutions. The pipeline has a complexity of $O(N \cdot d^2)$ for feature fusion and O(E log V) for the improved A* path search. On an RTX 3060, the average per-frame cost is ~4.2 GFLOPs with ~950 MB memory. Throughput tests show stable 57 fps for ≤50 tasks, decreasing to 44 fps at 100 tasks, indicating scalability under varying concurrency.

5.4 The value of research results in intelligent interaction and application expansion in AR scenarios

The multimodal perception information intelligent generation and visualization strategy proposed in this article has demonstrated significant application value in AR scenarios, providing reliable support for real-time perception and dynamic presentation in complex interactive environments. From the perspective of operational efficiency, the constructed model is able to maintain interaction latency below 1.5s, rendering frame rate stable above 55fps, and semantic consistency above 92% in the case of high concurrency from multiple sources of input. Compared to traditional methods, the interaction interruption rate has been reduced by nearly 60%, and the user response accuracy has been improved to 93%, fully demonstrating the robustness and adaptability of the model in high dynamic scenarios. In terms of interaction stability, the model can quickly distinguish abnormal signals such as speech noise interference and motion input blur, and automatically adjust the rendering path through feedback correction mechanism to ensure the continuous operation of the system. The experimental platform data shows that the number of rendering lags has decreased by more than 40%, and the smoothness of task execution has significantly improved. In terms of application scalability, this research results present multimodal states, rendering results, and feedback logic graphically through a visual interface, making the interaction process more transparent and facilitating real-time monitoring and strategy

optimization. This method can seamlessly integrate with existing AR engines and interaction platforms, and supports various hardware devices such as mobile terminals and smart glasses, with good cross platform deployment capabilities. The model proposed in this article demonstrates advantages in terms of interaction efficiency, system stability, and scalability. It not only supports immersive experiences in complex AR scenarios, but also provides a practical path for the promotion and application of intelligent interaction systems, laying a solid foundation for the industrialization and application expansion of AR technology in the future.

5.5 Comparison with State-of-the-Art (SOTA) Methods

We further compared our framework with representative SOTA models, including MulT (2019), Perceiver (2021), and Transformer-based AR design (Lee et al. 2023). MulT and Perceiver achieved semantic alignment rates of 90.1% and 91.3% with latencies of 2.6 s and 2.1 s, while our method reached 92.4% alignment with 1.42 s latency. In terms of stability, Lee et al.'s design maintained 49 fps, whereas our framework achieved 57 fps with <2% frame loss.

Ablation analysis shows that semantic modeling improved alignment by +7.8%, and feedback optimization reduced latency by ~20%, explaining the overall gain. These results confirm that our approach not only outperforms SOTA methods in accuracy, latency, and stability, but also ensures scalability on mid-range devices for real-time AR interaction.

6 Conclusion

This article focuses on the intelligent generation and visualization of multimodal perception information in AR scenes, proposing a feature modeling method that integrates graph convolution and attention mechanism. Combining the cross-modal mapping framework of variational and autoencoder geometric/temporal constraints, and introducing a reinforcement learning visualization optimization mechanism, closed-loop system of "perception generation presentation feedback" is constructed. The experimental results show that this strategy outperforms traditional methods in terms of interaction efficiency, semantic consistency, and rendering stability, with an average delay shortened to 1.4s, a rendering frame rate stable above 57fps, and a semantic alignment rate exceeding 92%. This validates its robustness and practicality in complex dynamic interaction environments. The system performs well in resource utilization and delay control, and can run stably in mid-range devices and multi platform environments, with application feasibility. However, there are still shortcomings in this study. Firstly, the experimental dataset is limited in size and mainly relies on public data and small-scale self built datasets. Further validation of the model's generalization ability is needed in large-scale and multi scenario scenarios; Secondly, the convergence speed of reinforcement learning in complex tasks is slow, which

may lead to high training costs and hinder large-scale real-time deployment. Future research can explore self supervised pre training and transfer learning mechanisms to enhance cross scenario adaptability; Simultaneously combining distributed computing and lightweight model compression to further optimize convergence efficiency and resource utilization. In addition, the framework of this study can be expanded in multi terminal collaboration and cross platform applications to enhance its application value in fields such as healthcare, industrial collaboration, and education.

Supplementary materials

A supplemental package is provided, including the source code, dataset generation script, trained model checkpoints, and a README file, to ensure reproducibility and facilitate further research.

Appendix A: dataset and preprocessing steps

Dataset

Self-built multimodal dataset: visual, speech, and motion data.

28,000 instances with timestamps for semantic alignment.

460,000 records for rendering/interaction (frame rate, latency, frame loss).

16,000 records for environmental/feedback data to evaluate model adaptability.

Preprocessing

Time alignment: Linear interpolation and synchronization.

Structural mapping: Map inputs to feature tensors and generate labels.

Denoising: Bandpass filter (300Hz-3kHz) for speech noise; blur detection for visual data.

Standardization: Features standardized to [-1,1].

Sliding window: Split dataset and add 15 abnormal samples for robustness testing.

Hardware and Software

Hardware: NVIDIA RTX3060,32GB memory, Intel i7 Software: PyTorch1.10, AnyLogic8.7, Kafka2.8.0

Training Plan
Epochs: 500
Optimizer: Adam
Learning rate: 1e-4

Data augmentation: Random cropping, rotation

Early stop: Stop if validation loss doesn't improve for 10 rounds.

Hyperparameters and Benchmarks

Model: Graph convolutional networks + attention mechanisms

Hyperparameters: 3x3 conv layer, 128 hidden nodes, batch size 64

Benchmark: Compared to single-modal and heuristic fusion models.

Pseudocode

Here is the pseudocode for the model training process:

#Initialize model with GCN+Attention mechanism

model = GCN_Attention_Model()

Training loop

for epoch in range(epochs):

for batch in data_loader:

inputs, labels = batch

outputs = model(inputs) # Forward pass

loss=compute_loss(outputs,labels)#Compute loss

optimizer.zero_grad() # Clear gradients

loss.backward() # Backward pass

optimizer.step() # Update weights

#Early stopping if validation loss doesn't improve if validation loss > threshold:

break

Benchmark Method

Single-modal model: Basic CNN trained on a single modality (e.g., visual data).

Heuristic fusion model: Fuses modalities using fixed rules, without dynamic optimization.

Fair comparison: All models trained with the same computational conditions and hyperparameters.

Labels: Automatic metrics checked by experts.

Temporal: 30 frames per sequence (0.1 s), aligned with speech and action.

Perturbations: Include SNR shifts, blur, occlusion, and lighting change.

Runs: 100 distinct seeds/scenarios for statistical reliability.

References

- [1] Ismail A W , Sunar M S .Multimodal Fusion: Gesture and Speech Input in Augmented Reality Environment[J].Advances in Intelligent Systems and Computing, 2015, 331:245-254.https://dol:org/10.1007/978-3-319-131 53-5_24
- [2] Yong J , Wei J , Lei X ,et al.Intervention and regulatory mechanism of multimodal fusion natural interactions on AR embodied cognition[J].Information Fusion, 2024,117.https://dol:org/10.1016/j.inffus.2024.1029 10
- [3] Chen L, Zhao H, Shi C, et al. Enhancing multi-modal perception and interaction: an augmented reality visualization system for complex decision making[J]. Systems, 2024,12(1):7.https://dol:org/10.3390/systems12010 007
- [4] Lee G-A, Sedlmair M, Schmalstieg D. Design patterns for situated visualization in augmented reality[J]. arXiv preprint, 2023,arXiv:2307.09157.https://dol:org/10.48550/arXiv.2307.09157
- [5] Zollmann S , Langlotz T , Grasset R ,et al. Visualization Techniques in Augmented Reality: A Taxonomy, Methods and Patterns.[J].IEEE

322

- transactions on visualization and computer graphics, 2021,
- 27(9):3808-3825.https://dol:org/10.1109/TVCG.202 0.2986247
- [6] Zheng M, Lillis D, Campbell AG. Current state of the art and future directions: Augmented reality data visualization to support decision-making[J]. Visual Informatics, 2024, 8(2):80-105. https://doi.org/10.101 6/j.visinf.2024.05.001
- [7] Friske MD. Integration of Augmented Reality and Mobile Robot Indoor SLAM for Enhanced Spatial Awareness[J]. arXiv preprint,2024,arXiv:2409.01915.https://dol:org/10.4 8550/arXiv.2409.01915
- [8] Al-Tawil B. A review of visual SLAM for robotics: evolution, properties, and relevance to augmented reality[J]. Frontiers in Robotics and AI,2024,11:1347985.https://dol:org/10.3389/frobt.2 024.1347985
- [9] Sheng X, Mao S, Yan Y, et al. Review on SLAM algorithms for augmented reality[J]. Displays,2024,84(2):102806.https://dol:org/10.1016 /j.displa.2024.102806
- [10] Barros AM. A comprehensive survey of visual SLAM algorithms[J]. Robotics, 2022,11(1):24.https://dol:org/10.3390/robotics1101 0024
- [11] Taketomi T, Uchiyama H, Ikeda S. Visual SLAM algorithms: a survey from 2010 to 2016[J]. IPSJ Transactions on Computer VisionandApplications,2017,9:1.https://doi.org/10.1 186/s41074-017-0027-2
- [12] Xu C, Kumaran R, Stier N, et al.Multimodal 3D Fusion and In-Situ Learning for Spatially Aware AI[J]. IEEEISMAR2024.https://dol:org/10.1109/ISMAR6 2088.2024.00063
- [13] Zhao F, Wang J, Li S, et al. Deep multimodal data fusion: a survey[J]. ACM Computing Surveys,2024,56(5):1–36.https://dol:org/10.1145/3649447
- [14] José Morano, Aresta G, Grechenig C, et al.Deep Multimodal Fusion of Data With Heterogeneous Dimensionality via Projective Networks[J].Journal on Biomedical and Health Informatics (J-BHI),2024,28(4):12.https://dol:org/10.1109/JBHI .2024.3352970
- [15] Ni J, Chen X, Yang Y, et al. Deep equilibrium multimodal fusion[J]. arXiv preprint,2023,arXiv:2306.16645.https://dol:org/10.4 8550/arXiv.2306.16645
- [16] Radford A, Kim J W, Hallacy C, Ramesh A, Goh G, Agarwal S, Sastry G, Askell A, Mishkin P, Clark J, Krueger G, Sutskever I. Learning Transferable Visual Models From Natural Language Supervision[C]// Proceedings of the 38th International Conference on Machine Learning

- (ICML). PMLR,2021:8748-8763.https://proceedings.mlr.press/v139/radford21a.html
- [17] Zheng B, Hu H. Multimodal Image Fusion and Classification of Power Equipment Using Non-Subsampled Contourlet Transform and Adaptive Pulse-Coupled Neural Network [J]. Informatica, 2024, 49(2):37-44. https://doi.org/10.31449/inf.v49i26.8729
- [18] Gao J , Li P , Chen Z ,et al.A Survey on Deep Learning for Multimodal Data Fusion[J].Neural Computation, 2020, 32(1):1-36.https://dol:org/10.1162/neco a 01273
- [19] Zhong R, Hu B, Feng Y, et al. Construction of human digital twin model based on multimodal data and its application in locomotion mode identification[J]. Chinese Journal of Mechanical Engineering, 2023, 36: 126.https://dol:org/10.1186/s10033-023-00951-0
- [20] Cao C, Jiang Z, Wu H, et al. Study of deep multimodal information fusion–based digital twin method for gearbox fault diagnosis[J]. The International Journal of Advanced Manufacturing Technology, 2024,138:3529-3542.https://doi.org/10.1007/s00170-025-15673-x