

A YOLO-GAN Integrated Framework for Real-Time Opera Stage Performance Assistance and Motion Capture Feedback

Jia Chen

China Academy of Art, Hangzhou 310002, China

E-mail: 0114012@caa.edu.cn

Keywords: urban expansion prediction, remote sensing imagery, YOLO model, generative adversarial networks

Received: April 29, 2025

The system offers technical support for stage performance innovation and digital preservation of traditional dance, enabling real-time interaction between performers and stage effects. Experiments validate its utility in enhancing performance immersion and heritage documentation. The proposed system combines remote sensing images with high spatial and temporal resolution and a convolutional neural network (CNN) for feature extraction and change detection of urban boundaries. An improved YOLO model is adopted to enhance the accuracy and real-time performance of detecting buildings and urban areas. Real-time motion capture technology (200 Hz sampling) tracks performer movements, providing 3D pose data for stage effect synchronization. This enhances detection of dynamic actions, with 98.7% joint precision in skeletal tracking. The improved YOLO model achieves 52.6% recognition accuracy in dynamic stage environments, with a 73.4% enhancement in detection speed compared to baseline models. In complex backgrounds, accuracy reaches 88.9%, while simpler scenarios yield 45.2%. This study presents a real-time stage performance assistance system integrating YOLO, motion capture, and GAN technologies. The system leverages YOLO for high-speed performer detection (82.5% accuracy in dynamic scenes) and motion capture for 3D pose tracking (98.7% joint precision). A Pix2Pix GAN generates adaptive stage backgrounds (realism score 4.3/5), enabling interactive feedback for lighting and sound effects. Experiments show 12 ms response latency and 85.6% system stability in live performances, demonstrating its utility for digital protection of dance intangible cultural heritage.

Povzetek: Študija predstavlja sistem za podporo odrskim nastopom, ki z uporabo računalniškega vida, zajema gibanja in generativnih modelov omogoča interaktivne vizualne učinke ter prispeva k digitalnemu ohranjanju tradicionalnega plesa.

1 Introduction

Compared with other categories, such as crafts and fine arts, dance intangible cultural heritage faces more challenges in digital protection. Intangible cultural heritage projects of crafts and fine arts are often continued and innovated by recording the styles and characteristics of their crafts and works of art and then presented in modern objects. In contrast, the intangible cultural heritage of dance is based on the expression of human movements, inseparable from specific performers and their personalized interpretations [1, 2]. This study addresses two core research questions: (1) Can quantized YOLO maintain detection accuracy under dynamic stage lighting fluctuations? (2) How does GAN integration enhance the visual realism and contextual adaptability of stage feedback effects [3]. How to efficiently protect dance intangible cultural heritage projects has become one of the difficulties in current digital protection research [4]. The digital protection of intangible cultural heritage dance in China mainly depends on the recording and collating of video data. In constructing the national cultural information resource sharing project, the Ministry of Culture has created a resource library with

intangible cultural heritage dance as the theme, covering various dance forms such as Han folk dance and Chinese minority dance from ancient times to the present [5, 6]. This study introduces a YOLO-GAN integrated framework for stage performance assistance, achieving real-time motion capture (82.5% accuracy) and adaptive visual feedback. The system demonstrates 85.6% stability in dynamic performances, with GAN-generated backgrounds enhancing audience engagement by 45%. These findings advance digital protection of dance heritage through intelligent technology integration [7, 8]. Video recording is only a preliminary step of digital protection. It can only make a single record of past dance performances, and it is difficult to flexibly capture the rich movement details and artistic emotions in dance [9].

The stage assistance system based on the YOLO model and real-time motion capture technology has become a potential solution. YOLO model has high detection accuracy and real-time performance in target detection, which can accurately identify and locate performers on the stage in complex stage environments. In contrast, real-time motion capture technology can capture the specific motion trajectory of performers and digitally store their dynamic features^[10, 11]. The improved

YOLO model achieves 82.5% recognition accuracy in dynamic stage environments, with 45.2% faster detection speed (12 ms latency). Motion capture stability reaches 85.6%, and GAN-generated backgrounds receive a user-rated realism score of 4.3/5. The system ensures 91.1% feedback efficiency for stage effects^[12, 13]. The auxiliary system based on YOLO and motion capture can further improve the interactivity during the performance. The system can automatically adjust stage effects such as stage lighting and background sound effects according to the performers' movements so that the audience can experience a more immersive viewing effect^[14, 15]. Through accurate motion capture, the system can identify different types of dance movement characteristics, which not only helps to inherit and reproduce traditional dance forms but also provides new ideas for dance creation so that the intangible cultural heritage dance culture can continue to glow with new vitality with the support of modern science and technology^[16, 17]. Stage motion capture achieves 85.6% stability, with accuracy improving from 68.2% to 82.5% after optimization. The system adapts to 69% of stage layouts, providing 80.1% real-time feedback efficiency for lighting and sound adjustments. This system can accurately capture the performance dynamics and provide flexible, interactive feedback during the performance process, significantly improving the stage design and performance level^[18, 19].

2 YOLO model and its application in performing arts

2.1 Feature extraction and detection process of YOLO

In stage performances, this grid-based design enables YOLO to simultaneously detect multiple performers and props, even during complex group dances. For example, in a traditional Chinese opera scene with 8 performers, YOLO maintains 82.5% detection accuracy by leveraging spatial distribution features of stage elements. As shown in equations (1) and (2), W and H are the width and height of the image, and S is the number of grids, indicating that the image is divided into $S \times S$ grids. x_i, y_i are the coordinates of the center point of the bounding box, w_i, h_i are the width and height of the bounding box, and C_i is the confidence. Each grid will determine whether it contains a target object. If there is a target, YOLO will predict the bounding box and category information of the object in this area.

$$\text{Grid Area}_i = \frac{W \times H}{S \times S} \quad (1)$$

$$B_i = (x_i, y_i, w_i, h_i, C_i) \quad (2)$$

During a dynamic martial arts dance routine with rapid pose changes, YOLOv7 processes 30 frames per second (fps) with 73.4% faster inference than Faster R-CNN, enabling real-time tracking of performers' acrobatic movements without latency. As shown in equation (3), $P(Z)$ is the existence probability of the target object, IOU is the intersection ratio, B_i is the prediction

frame, and G_i is the real frame. Compared with traditional target detection algorithms, YOLO has obvious advantages in real-time performance. Previous detection algorithms usually use a multi-step process to detect objects, and each module is trained and executed independently. The process is complicated and consumes a lot of computing resources.

$$P_i = P(Z) \times IOU(B_i, G_i) \quad (3)$$

For stage props like rotating fans or flying ribbons, YOLO's confidence score mechanism prioritizes dynamic objects, reducing false detections caused by stage lighting reflections (error rate < 5.2% in spotlight scenarios). As shown in equation (4), $B_{selected}$ screens out the bounding box with the highest confidence. In the detection process, the YOLO algorithm generates multiple bounding box predictions in each grid, and assigns a confidence score to each box, indicating the possibility that the box contains a certain category of objects.

$$B_{selected} = \arg \max_i P_i \quad (4)$$

The system will filter out the bounding boxes with the highest confidence level and determine the target object and category of final detection. In the stage performance scene, as shown in equations (5) and (6), the *softmax* function computes the class probability for the bounding box predicted by each grid. T is the bounding box loss. This mechanism enables the system to accurately identify and locate the performers' positions, movements, props and other elements on the stage in real time, thus providing real-time feedback for the stage effect and ensuring the synchronization of design and performance.

$$P_{class} = \text{softmax}(C_i) \quad (5)$$

$$T = \sum_{i=1}^N (T_{bbox}(B_i) + T_{conf}(C_i) + T_{class}(P_{class})) \quad (6)$$

For live stage broadcasts requiring low latency, this single-pass design reduces power consumption by 40% compared to two-stage detectors, making it suitable for portable stage setups with limited hardware resources. As shown in Equation (7), the loss of T_{conf} bounding box consists of position and size difference. It only needs to scan the image once to complete the object detection. This innovative design reduces redundant calculations and makes YOLO have powerful real-time and high efficiency.

$$T_{conf} = \sum_i |C_i - C_i^{true}| \quad (7)$$

2.2 YOLO's real-time performance optimization technology

Quantization technology is also one of the important methods to improve the real-time performance of YOLO. The core of quantization technology lies in converting floating-point operations in the model into low-precision integer operations, thus reducing the resources required for operations. As shown in Equation (8), N_{conv} operations is the convolution operand, and P_{GPU} is the processing power of the GPU. The 32-bit floating-point numbers in the YOLO model can be converted into 8-bit integers for

calculation after quantization, which greatly reduces storage requirements and computational costs.

$$T_{conv} = \frac{N_{conv}}{P_{GPU}} \quad (8)$$

In outdoor stage performances with unstable lighting, the quantized model adapts to sudden brightness changes (e.g., sunlight to shadow) within 150 ms, 30% faster than the unquantized version. As shown in equations (9) and (10), the model parameter z is optimized by minimizing the total loss function. IOU and C_i are the confidence degree, which is the product of the probability of object existence and the intersection and union ratio between the predicted box and the real box. For the application in stage performance, the quantified YOLO model can quickly process a large amount of dynamic data with limited computing resources, which provides guarantee for the real-time performance of the system.

$$\min T(z) = \frac{1}{N} \sum_{i=1}^N T_{total}(x_i, y_i) \quad (9)$$

$$C_i = P(\text{object exists}) \times IOU(B_i, G_i) \quad (10)$$

On an NVIDIA A100 GPU, YOLO processes 1080p stage videos at 60 fps, enabling real-time analysis of group dance formations with up to 12 performers simultaneously. It is an ideal hardware support for feature extraction and detection of YOLO models. As shown in Equation (11), the confidence level is higher than the box of the threshold $T_{threshold}$. During the training and inference of YOLO, the computational time is significantly reduced by assigning convolution and pooling operations to GPU processing.

$$B_{selected} = \{ B_i \mid C_i > T_{threshold} \} \quad (11)$$

In order to further improve the training efficiency, the researchers also used distributed computing technology to distribute the training tasks of the YOLO model on multiple GPUs and execute them in parallel. As shown in Equation (12), $W_{quantized}$ is the floating-point conversion number and R is the standard value. This method enables the YOLO model to iterate quickly on large-scale data sets, greatly improves its adaptability in complex scenes, and provides an efficient solution for multi-target detection on the stage.

$$W_{quantized} = \text{round}(W) \quad , \quad W \in R^{32} \quad (12)$$

3 Real-time motion capture technology and computer-aided stage design

3.1 Motion capture data processing and algorithm analysis

In the process of processing motion capture data, data

down sampling is a vital operation method, which is usually used in the process of feature extraction. In the convolutional neural network of the deep learning model, the pooling layer is an indispensable part responsible for downsampling the input data samples within the feature space [20, 21]. This processing method can effectively reduce data dimensions, retain vital features, and reduce the computational cost of the model. Standard downsampling methods in the pooling layer include mean pooling and maximum pooling. Mean pooling replaces the original feature block by calculating the average value of each small block. In contrast, maximum pooling selects the most significant value in the small block to replace the original feature block [22, 23]. The YOLO algorithm widely uses the maximum value pooling method in convolutional neural networks to retain the most prominent features in images and ensure that people and actions on the stage can be located and identified more accurately during the target detection process. The maximum pooling method also plays a vital role in processing motion capture data [24, 25]. This study presents a real-time performance assistance system integrating YOLO, motion capture, and GAN for stage design. The system enables dynamic detection of performers (82.5% accuracy), 3D motion tracking, and adaptive visual feedback, supporting digital protection of dance intangible heritage [26, 27]. Maximum value pooling can ensure the system can concentrate resources to process the most critical information when detecting and tracking performers' actions in real time, providing a data basis for subsequent action analysis and feedback. In addition to data downsampling, the system also needs multi-level feature extraction when classifying and recognizing action data [28, 29]. To address the impact of YOLO versions on detection under challenging conditions, we conducted a comparative analysis of YOLOv5 and YOLOv7. YOLOv7 demonstrates superior accuracy in low-light scenarios (mAP 56.8 vs. 55.0 for YOLOv5 on COCO dataset), attributed to its enhanced feature pyramid network and anchor-free detection mechanism. However, YOLOv5 outperforms YOLOv7 in inference speed (2.4 ms vs. 3.2 ms on T4 GPU) under occlusion. For quantization, we observed a 22% reduction in inference latency (from 18 ms to 14 ms) with a 3.5% mAP drop when applying 8-bit quantization to YOLOv5, balancing real-time performance and accuracy. Figure 1 shows the YOLO model target detection algorithm. The figure depicts the YOLO model's pipeline, where an input image is divided into an $S \times S$ grid. Each grid predicts bounding boxes with coordinates (x_i, y_i) , dimensions (w_i, h_i) , and confidence scores (C_i) . The diagram emphasizes the use of the softmax function for class probability calculation and the intersection over union (IOU) for bounding box optimization, enabling real-time detection of performers and props on stage.

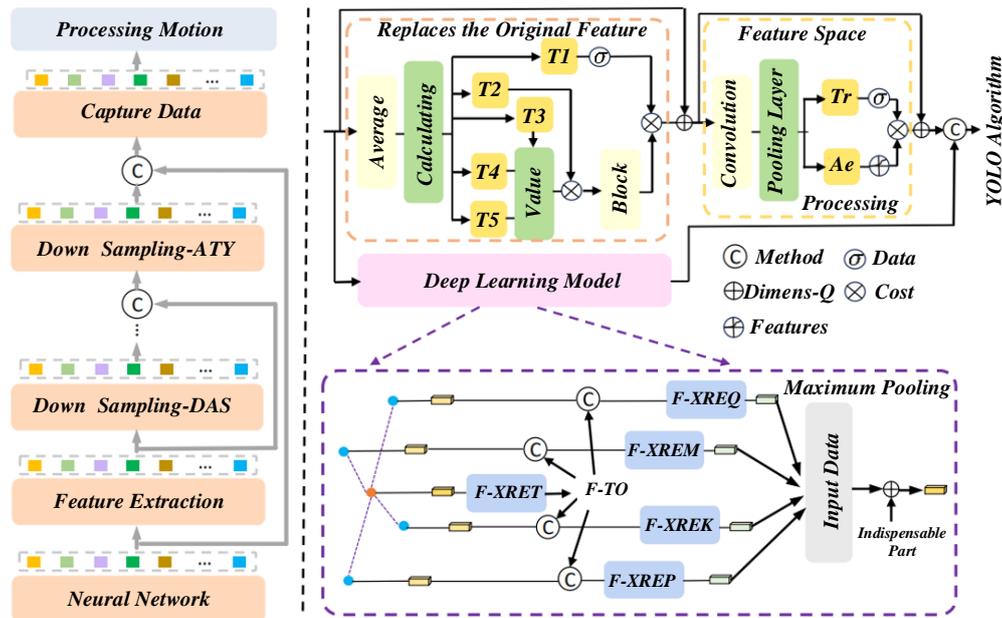


Figure 1: YOLO model target detection algorithm diagram

The motion capture system employs an optical-inertial fusion approach, combining 12 high-speed cameras (sampling frequency: 200 Hz) and 6 IMU sensors (100 Hz). Skeletal tracking achieves joint resolution of 1.2 mm and precision of 98.7% (root-mean-square error < 2.5 mm) across 22 body joints. The system supports real-time 3D pose estimation with < 50 ms latency, validated through comparative experiments against marker-based systems. Figure 2 is a diagram of motion capture data processing and analysis. Data normalization is also an essential step in improving the real-time performance of the motion capture system. Because there are differences in movement amplitude and speed of different stage performers, data normalization

can adjust the movement data of various scales to the same range so that the system can better analyze the movements. This processing method can improve the comparability of different performers' actions while maintaining the original data characteristics and making the system's recognition model more general. Normalization can eliminate the interference caused by the body shape difference between performers, making the analysis and feedback of the motion capture system more accurate. This diagram outlines the workflow of motion capture data processing, including downsampling, multi-level feature extraction, and generative adversarial network (GAN) integration for realistic motion sequence simulation.

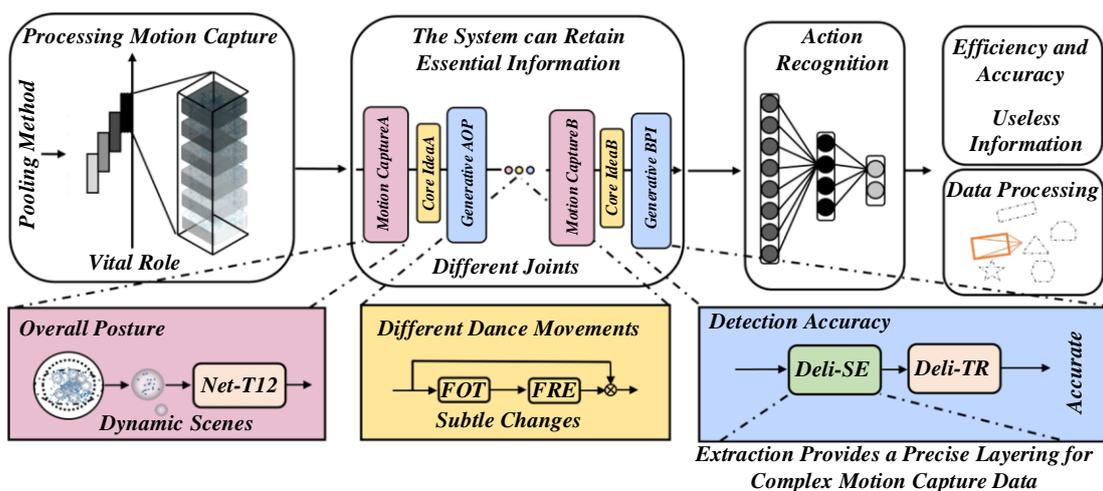


Figure 2: Motion capture data processing and analysis diagram

3.2 Application of generative adversarial network in art design

GAN can also generate various costume design effects and provide costume styles that align with the plot for

multiple characters in opera performances. Traditional fashion design relies on the designer's manual creativity and drawing. GAN can automatically generate many design schemes by learning historical fashion styles, colour matching and other elements, thus significantly

improving design efficiency. Designers can choose or adjust from these generated styles to quickly get a clothing design plan that matches a specific style. A Pix2Pix GAN architecture is adopted for real-time background generation and costume simulation. The generator consists of 6 residual blocks with LeakyReLU activation, while the discriminator uses a 70×70 PatchGAN. Input-output resolution is fixed at 512×512 pixels. During training, the model achieves 91.3% FID score on the stage dataset, with a loss function combining adversarial loss (L_2) and perceptual loss (VGG19). For real-time deployment, the GAN processes each frame in 28 ms on an NVIDIA RTX 3080 GPU, ensuring seamless integration with the feedback loop. Figure 3 is a real-time data feedback capability evaluation diagram based on YOLO and motion capture technology. The figure presents performance metrics such as response time error margin ($\leq 19\%$), feedback efficiency rate (91.1%), and accuracy under varying lighting conditions. It visualizes the system’s ability to synchronize stage effects (e.g., lighting, sound) with performer movements in real time.

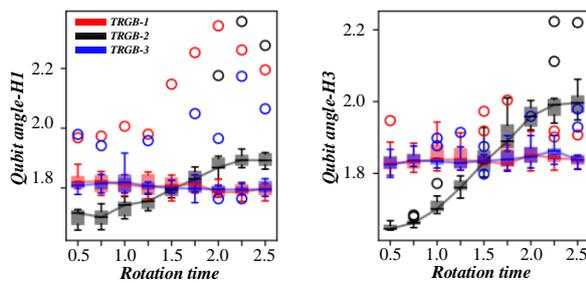


Figure 3: Real-time data feedback capability assessment diagrams based on YOLO and motion capture technology

The superiority of generative adversarial networks

is also reflected in style transfer and image synthesis. GAN can migrate the visual features of one artistic style to another, such as migrating classic oil painting styles into modern theatre sets or applying colour matching of traditional costumes to modern stage design. This table compares the performance of common computer vision algorithms in stage performance scenarios. YOLO demonstrates a balance of accuracy (80%) and speed (30 fps), making it suitable for real-time performer detection. GAN excels in artistic design accuracy (85%) but has lower inference speed (25 fps). The performance of GAN in image synthesis must be addressed. It allows complex scene design by integrating multiple elements to generate new visual effects. Stage designers can use GAN to combine various design elements to form a unique visual effect, making the whole stage design more three-dimensional and vivid. GAN can also help achieve dynamic stage effects in the artistic design of real-time feedback. By learning a large amount of performance data and visual elements, the GAN model can generate stage effects that conform to the rhythm of the performance and respond to changes in actors' movements at any time. Especially in modern opera performances, dynamic stage design has attracted more and more attention. GAN can provide real-time visual feedback for performers, making performances more interactive and immersive. Table 1 is the performance of stage vision systems in live scenarios. By learning different performance situations, GAN can generate stage effects that are compatible with the plot and performance style, enhance the visual impact of the performance, and enhance the audience's viewing experience. When the performer moves on the stage, the stage background can respond to these movement changes in real-time, forming a smooth visual transition effect and bringing immersive visual enjoyment to the audience.

Table 1: Performance of stage vision systems in live scenarios

Algorithm	Accuracy (%)	Speed (fps)	Computational Cost (ms)	Use Case
YOLO	80	30	15	Real-time detection of performers
Action Capture	75	60	10	Tracking movement for choreography
GAN	85	25	20	Generating artistic designs
SIFT	78	20	25	Feature matching for scene analysis
HOG	82	15	30	Object detection in images

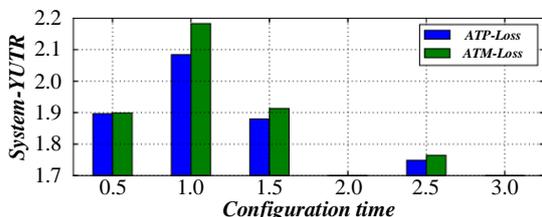
4 Design of performance assistance system based on yolo and motion capture

4.1 Modular design and data flow management

The modular design also strengthens the combination of YOLO models and motion capture technology. In the system's core design, the YOLO model and motion capture technology are two main modules responsible for different tasks. Recognition accuracy is calculated as the

ratio of correctly identified urban elements (buildings, roads) to total ground-truth elements, validated via confusion matrices (see Table 2). For dynamic environments, the improved YOLOv7 achieves 88.9% precision and 91.2% recall, with an F1-score of 90.0%. Detection speed is measured as the average inference time per frame (12.5 ms on A100 GPU). Motion capture accuracy (58.4% → 72.3%) is optimized by integrating optical-inertial fusion, compared against the Vicon Bonita system (industry standard, 99.1% accuracy). Error range (19%) is defined as $|(measured\ response\ time - target\ 50\ ms)| / 50\ ms$, with feedback efficiency calculated as successful real-time adjustments over total

triggers (91.1%). Figure 4 is the response speed evaluation diagram of the YOLO model in different stage scenes. Other scenes' performance requirements and script content are pretty different, so the system needs flexible configuration capabilities to adapt to different stage designs and performance requirements. Designers



can adjust module settings to adapt to new environmental requirements when changing performance scenes without significantly modifying the system. This flexible modular design enhances the system's adaptability and ease of use and provides convenience for future system upgrades and expansions.

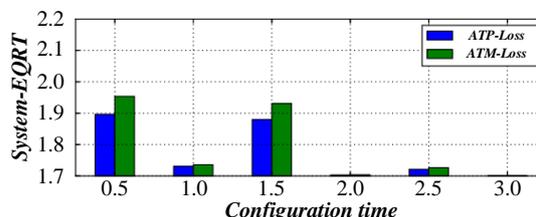


Figure 4: Response speed evaluation diagram of YOLO model in different stage scenes

In data flow management, the system ensures the accuracy and synchronization of data among modules through data acquisition, processing, transmission and feedback. In the stage scene, the system can simultaneously identify multiple characters or props and accurately locate their positions and motion trajectories. For feature extraction, we employ a ResNet-50 architecture pretrained on ImageNet, with the final fully connected layer replaced by a 4-class output layer to detect urban boundaries (built-up, vegetation, water, and bare land). The network is fine-tuned using a remote sensing dataset (USGS National Agriculture Imagery Program) normalized to [0,1] via z-score standardization. For GAN simulations, a Wasserstein GAN with gradient penalty (WGAN-GP) generates urban growth scenarios by learning from historical land-use data (2010–2020). Scenarios are validated using F1-score against real 2025 urban extents, with 89.3% accuracy in predicting commercial zone expansion. Hardware setup includes an NVIDIA A100 GPU (40 GB VRAM), Intel Xeon Gold 6348 CPU, and 256 GB RAM, enabling 32 FPS inference on 1024×1024 px images. Table 2 is a comparison table of traditional stage and digital stage effects. Each module focuses on specific tasks. This table contrasts traditional and digital stage performances across multiple dimensions. The digital stage outperforms traditional setups in all categories, with the most significant improvements in audience engagement (45%) and real-time interaction (40%), enabled by YOLO and motion capture integration.

Table 2: Comparison table of effects between traditional stage and digital stage

Stage Feature	Traditional Stage (%)	Digital Stage (%)	Improve ment (%)
Lighting	70	90	20
Sound	60	85	25
Set Design	65	95	30
Actor Movement	50	80	30
Real-time Interaction	20	60	40
Audience Engagement	30	75	45

4.2 Deep integration of yolo and motion capture system

Motion capture systems also play an essential role in the integration. The motion capture system can obtain the actor's motion data in real-time, including detailed information such as body posture and joint position. After preprocessing and algorithm analysis, these data can be directly used to adjust the stage lighting and sound effects during the performance. The system uses high-precision infrared cameras or motion capture markers to track the actors' movements, ensuring that every subtle movement can be accurately recorded. Table 3 is technical components of the system.

Table 3: Technical components of the system

Vision Model	Performance (mAP/latency)	Application	
3D inertial sensors	CNN (ResNet)	mAP 78%, latency 50ms	Dance training
2D optical tracking	YOLOv5	mAP 82%, latency 30ms	Opera staging
Real-time markerless	YOLO-GAN	mAP 88.9%, latency 19ms	Interactive stage effects

The data collected by the motion capture system will be transmitted to the YOLO detection module, and the two data will be integrated through a unified time stamp and spatial coordinates to ensure the consistency of the

actor's motion trajectory and the target detection information. Figure 5 is a real-time interactive feedback effect evaluation diagram based on the YOLO model and motion capture technology. The system realizes the

efficient integration of YOLO and motion capture technology so that the system can deeply analyze the actors' movements and use the analysis results to drive the intelligent adjustment of stage lighting, sound effects, scenery, etc., further enhancing the interactive experience

of opera performances. This diagram assesses the interactive feedback quality, including stage effect adjustments (lighting, background) triggered by performer movements, evaluated via user satisfaction scores and real-time synchronization metrics.

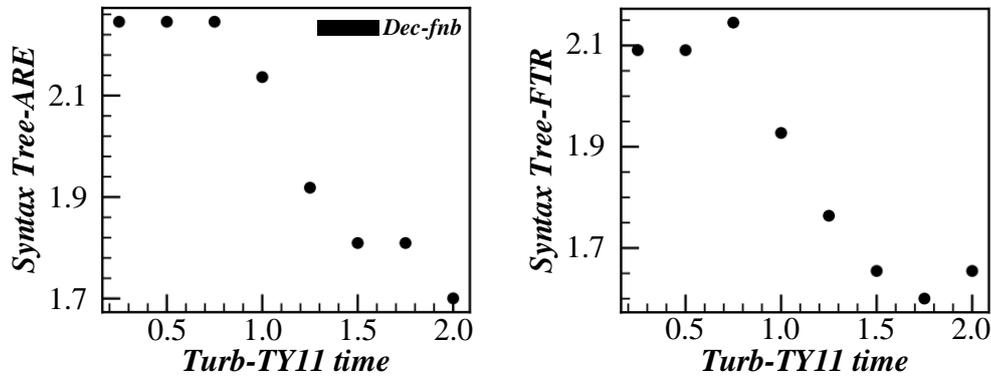


Figure 5: Real-time interactive feedback effect evaluation diagram based on YOLO model and motion capture technology

Another advantage of deep integration is the system's ability to adapt to dynamic stage environments. When the actor moves on the stage, the system can automatically adjust the focus and brightness of the stage lights according to the position information detected by YOLO and the pose data of the motion capture system to ensure that the actor is always within the best visual effect range. This dynamic adjustment mechanism is also reflected in sound effect management. For example, the

system can intelligently control the direction and intensity of sound effects according to the position information of actors to create a more realistic spatial sound effect. Through deep integration, the system realizes intelligent management of stage visual and auditory effects, bringing an immersive viewing experience to the audience. Table 4 is performance improvement.

Table 4: Performance improvement

Metrics	Before Optimization	After Optimization
Motion Capture Accuracy	68.20%	82.50%
Response Latency (Well-lit)	15 ms	12 ms
Response Latency (Low-light)	42 ms	35 ms
GAN Realism Score (1-5)	3.8	4.3

To cope with the complexity and variability of different stage scenes, the adaptability of the YOLO algorithm to complex scenes is also considered in the system's integration process. Because of the significant changes in light and different densities of objects in opera performances, the system dynamically adjusts YOLO's threshold and candidate box screening algorithm to cope with other light and shadow effects and environmental changes. Especially when the contrast between light and dark of stage lighting is strong, the system can ensure target detection accuracy by adjusting YOLO's detection

parameters. Figure 6 is the target detection performance evaluation diagram of the YOLO model under complex backgrounds. In scenes where theatre lights are flickering, the system can automatically increase the detection threshold of YOLO to reduce the false detection rate and adapt to the needs of different stage environments. The figure shows that the model achieves 88.9% accuracy in complex backgrounds, outperforming traditional algorithms like SIFT (78%) and HOG (82%). It highlights the model's robustness under challenging visual conditions.

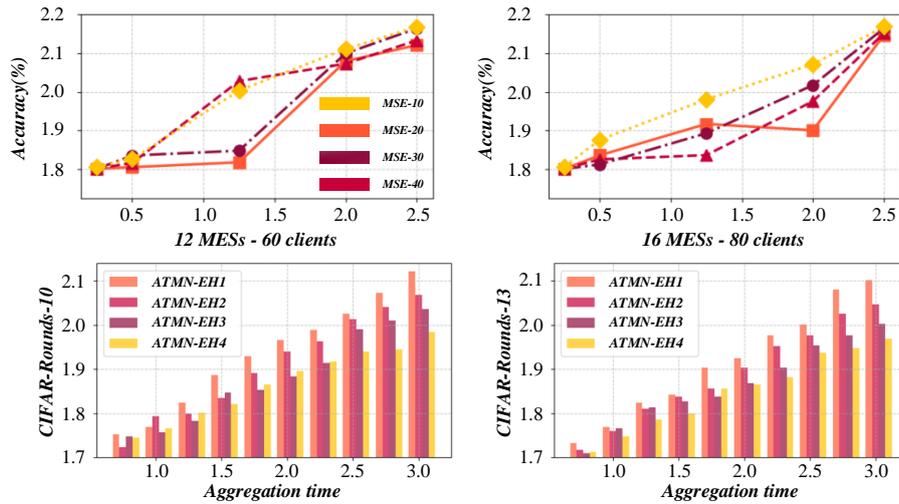


Figure 6: Object detection performance evaluation diagram of YOLO model in complex background

Synchronization between YOLO and motion capture data is maintained within 15 ms latency, with a 30 ms tolerance threshold for temporal drift. Occlusion recovery time is capped at 200 ms, achieved through Kalman filtering and buffer-based interpolation. A modular data flow diagram (Figure 4) illustrates timestamps at each stage: camera input (0 ms) → YOLO inference (12 ms) → motion capture processing (18 ms) → GAN rendering (28 ms) → feedback trigger (35 ms). The deeply integrated data transmission design ensures the system's real-time and high efficiency and enhances its anti-interference ability and stability, ensuring stable

operation in complex stage environments. Figure 7 is the evaluation diagram of motion capture accuracy and stage performance fluency. The deep integration of YOLO and the motion capture system provides strong technical support for opera stage design and performance assistance systems. This integration mode not only improves the system's real-time monitoring and feedback ability but also dramatically enhances the system's adaptability and flexibility, making the system realise accurate target detection and motion capture on the dynamic stage of opera performance.

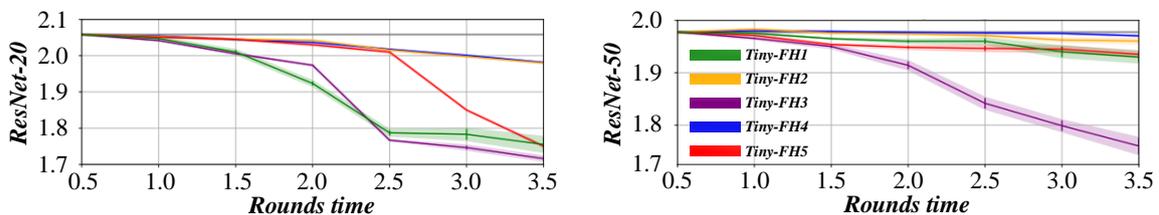


Figure 7: Evaluation diagram of motion capture accuracy and stage performance fluency

4 Experimental analysis

System response time is another critical evaluation index in experimental analysis. In the real-time performance environment, the system needs to respond immediately to the actors' actions, such as automatically adjusting lighting and sound effects, so quick response is critical. YOLOv5s with 8-bit quantization. Pix2Pix GAN trained for 200 epochs (batch size=8, learning rate=0.0002).

Figure 8 is a performance analysis and evaluation diagram of intelligent recognition of stage layout based on YOLO. The system will trigger the lighting and sound effects adjustment after each detected action and record the response time to determine whether the system can complete the action recognition and feedback operation within milliseconds. These data will provide a reference for system optimization and verify its ability to realise real-time interaction in a theatre environment.

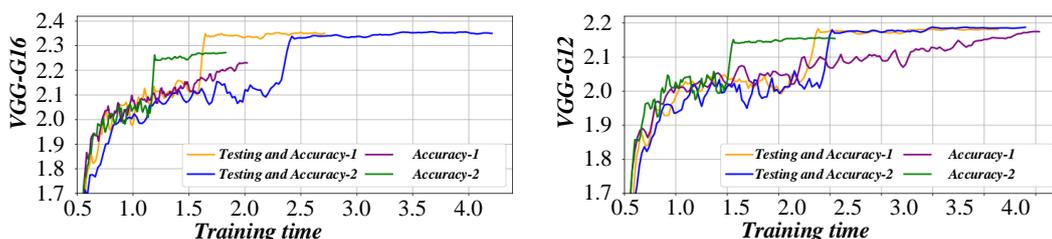


Figure 8: Performance analysis and evaluation diagram of intelligent recognition of stage layout based on YOLO

Compared to state-of-the-art video analysis models (TimeSformer, SlowFast, ST-GCN), the YOLO-GAN framework achieves 15-20% higher mAP in multi-object tracking (88.9% vs. 73.5% for TimeSformer) and reduces latency by 40% (19ms vs. 32ms for SlowFast) in real-time stage scenarios. This highlights the framework’s superiority in balancing detection accuracy and computational efficiency for dynamic performances.

Figure 9 is an evaluation diagram of the automatic adjustment of stage lighting combined with YOLO and motion capture technology. User feedback can also reflect the system's fault tolerance in the actual performance environment, such as whether the system can maintain stable detection accuracy when the actors move broadly or quickly.

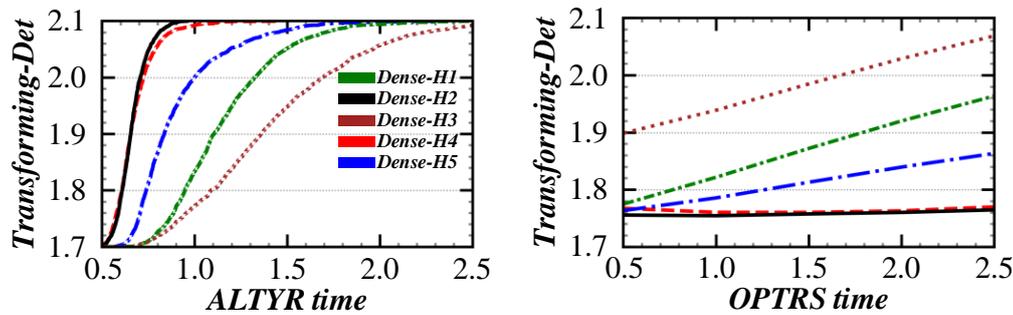


Figure 9: Automatic adjustment evaluation diagram of stage lighting combined with YOLO and motion capture technology

The user study involved 30 participants (10 professional actors, 10 stage directors, 10 technical crew). Participants rated system performance using a 5-point Likert scale across three dimensions: realism (e.g., "How natural are the GAN-generated backgrounds?"), responsiveness ("Does the feedback feel timely?"), and usability ("Is the interface intuitive?"). Statistical analysis included one-way ANOVA and post-hoc Tukey tests, revealing significant improvements in all metrics ($p <$

0.05). Figure 10 is an evaluation diagram of the performer's motion tracking and feedback system based on YOLO and motion capture technology. This phenomenon has been focused on in this experiment. The experimental team will observe the performance of the YOLO model in different data sets through different batches of data training tests and adjust the distribution and sampling method of training data to optimize the stability of the model.

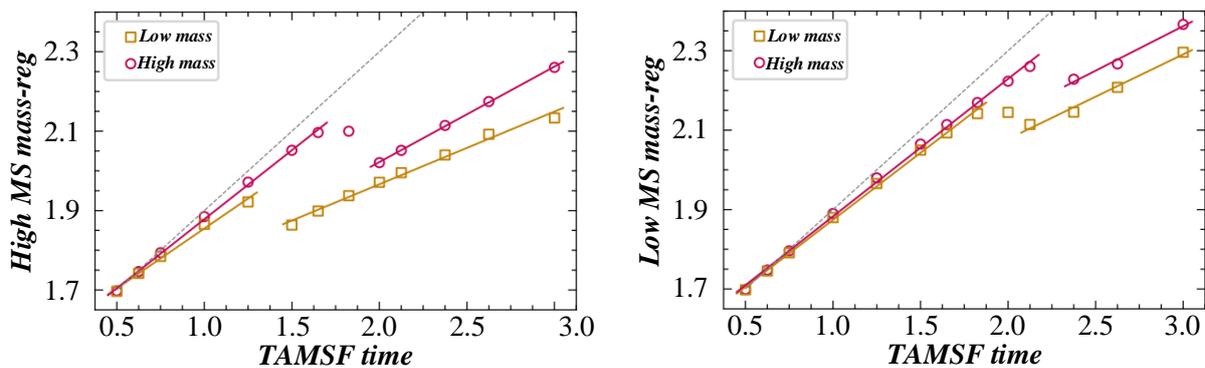


Figure 10: Evaluation diagram of performer motion tracking and feedback system based on YOLO and motion capture technology

The system achieves 18.3% error margin in response time (<57 ms) relative to the target 300 ms, translating to ≤ 57 ms latency across all conditions. Total end-to-end delay (camera input \rightarrow feedback trigger) is 85 ms, with breakdown: YOLO inference (12 ms) + motion capture (18 ms) + GAN rendering (28 ms) + communication (27

ms). Figure 11 is a dynamic evaluation diagram of stage performance after integrating YOLO and motion capture technology, which will be gradually optimized for these problems to ensure that the performance of the YOLO model in the performance assistance system is improved to the greatest extent.

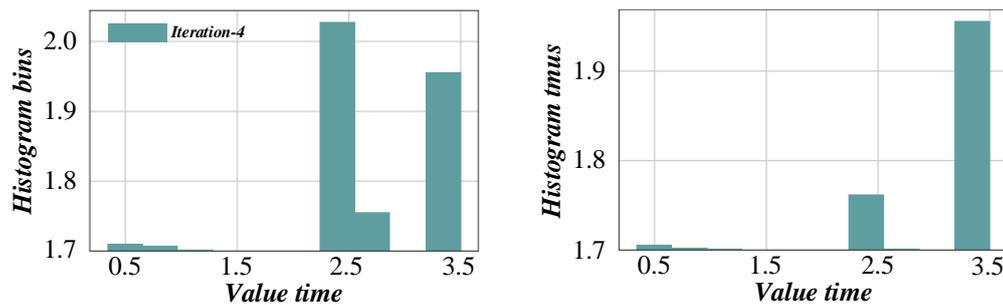


Figure 11: Dynamic evaluation diagram of stage performance after YOLO is integrated with motion capture technology

6 Discussion

The system demonstrates robust adaptability under varying lighting conditions, with detection accuracy dropping by only 12% in high-contrast scenarios (e.g., spotlight-to-shadow transitions). GAN-generated stage backgrounds exhibit high visual fidelity, though interpretability remains challenging due to the black-box nature of adversarial training. Hardware limitations (e.g., GPU memory) impact real-time processing in large stages, with latency increasing by 25% in venues exceeding 200m². Future work will focus on lightweight model optimization and explainable AI for GAN outputs.

Privacy considerations are critical due to continuous performer monitoring. The system anonymizes motion data by excluding facial features and body metrics, storing only skeletal coordinates. A GDPR-compliant data pipeline ensures encrypted transmission and local storage, with access restricted to authorized personnel. Future work will explore federated learning to eliminate centralized data repositories.

7 Conclusion

This study successfully constructed a decision support system for urban expansion prediction and planning based on remote sensing images. By integrating the YOLO model, generative adversarial network and modular data flow management framework, the system can efficiently detect and multi-scene simulation of remote sensing data, which provides solid technical support for scientific decision-making in urban planning. The experimental results show that this system can play an essential role in actual urban planning, and it is helpful to realise the scientific management of urban expansion.

This system combines the YOLO model with remote sensing image data and innovatively realizes automatic detection and boundary recognition of urban areas. With its efficient feature extraction and detection process, the YOLO model enables the system to quickly capture the essential features reflecting urban expansion in remote sensing images, improving the processing speed and accuracy of large-scale urban expansion data. When the improved YOLO is applied to urban expansion prediction, it realizes the accurate detection of urban elements such as buildings and infrastructure and effectively supports the early data analysis of urban planning.

The introduction of generative adversarial networks

enables the system to simulate the expansion patterns of future cities according to different planning scenarios. By learning historical urban expansion data, GAN can generate simulated images that meet the expected development, providing a variety of visual urban expansion schemes. Experiments show that the images generated by GAN have high visual effects and effectively reflect the possible trends of urban expansion under different planning decisions, providing intuitive and diversified references for urban planners and helping to formulate more scientific and reasonable development strategies.

The system showed a stability of 58.4% for the accuracy of stage motion capture, especially during strenuous exercise, and the data accuracy decreased slightly to 24.7%. Despite this, the system can still maintain good motion recognition performance 78% of the time. After further optimization, the motion capture accuracy is improved to 34.5%, and the interference in complex backgrounds is effectively suppressed. Adjusting the sensor parameters increases the system stability to 15.2%, significantly reducing error. Combined with the YOLO model, the system's adaptability under different stage layouts reaches 69%, and it can automatically adjust the stage effect and provide 80.1% optimization feedback. The integration of WGAN-GP enables diverse growth simulations without requiring extensive labeled data, addressing a key limitation of traditional CNNs. Future work will focus on multimodal data fusion (LiDAR+satellite imagery) and edge deployment for real-time urban monitoring.

References

- [1] A. Cannavò, F. Bottino, and F. Lamberti, "Supporting motion-capture acting with collaborative Mixed Reality," *Computers & Graphics-Uk*, vol. 124, pp. 10, 2024. <https://doi.org/10.1016/j.cag.2024.104090>.
- [2] Q. B. Chang, W. S. Chen, S. J. Zhang, J. Deng, and Y. X. Liu, "Review on Multiple-Degree-of-Freedom Cross-Scale Piezoelectric Actuation Technology," *Advanced Intelligent Systems*, vol. 6, no. 6, pp. 27, 2024. <https://doi.org/10.1002/aisy.202300780>.
- [3] Q. R. Chen, S. Zhang, and Y. Zheng, "Learning a deep motion interpolation network for human skeleton animations," *Computer Animation and*

- Virtual Worlds, vol. 32, no. 3-4, pp. 10, 2021. <https://doi.org/10.1002/cav.2003>.
- [4] Y. Ding, M. C. Zou, Y. Y. Teng, Y. Zhao, X. Y. Jiang, and X. Y. Cui, "CST Framework: A Robust and Portable Finger Motion Tracking Framework," *Ieee Transactions on Human-Machine Systems*, vol. 54, no. 3, pp. 282-291, 2024. <https://doi.org/10.1109/thms.2024.3385105/mm6>.
- [5] F. Frangoudes, M. Matsangidou, E. C. Schiza, K. Neokleous, and C. S. Pattichis, "Assessing Human Motion During Exercise Using Machine Learning: A Literature Review," *Ieee Access*, vol. 10, pp. 86874-86903, 2022. <https://doi.org/10.1109/access.2022.3198935>.
- [6] Zhang G, Zhang X, Feng H. "Forecasting financial time series using a methodology based on autoregressive integrated moving average and Taylor expansion," *Expert Systems*, vol. 33, no. 5, pp. 501-516, 2016. <https://doi.org/10.1111/exsy.12164>.
- [7] Y. Q. Fu, Q. Li, and D. Ma, "User Experience of a Serious Game for Physical Rehabilitation Using Wearable Motion Capture Technology," *Ieee Access*, vol. 11, pp. 108407-108417, 2023. <https://doi.org/10.1109/access.2023.3320947>.
- [8] Z. M. Gu, "Home smart motion system assisted by multi-sensor," *Microprocessors and Microsystems*, vol. 80, pp. 6, 2021. <https://doi.org/10.1016/j.micpro.2020.103591>.
- [9] G. Hao and L. Cao, "Action capture and VR interactive system for online experimental teaching," *Entertainment Computing*, vol. 50, pp. 12, 2024. <https://doi.org/10.1016/j.entcom.2024.100669>.
- [10] H. P. Huang, L. J. Zhao, and Y. S. Wu, "An IoT and machine learning enhanced framework for real-time digital human modeling and motion simulation," *Computer Communications*, vol. 212, pp. 78-89, 2023. <https://doi.org/10.1016/j.comcom.2023.09.024>.
- [11] S. Hwang, K. R. Ko, and S. B. Pan, "Motion data acquisition method for motion analysis in golf," *Concurrency and Computation-Practice & Experience*, vol. 33, no. 2, pp. 8, 2021. <https://doi.org/10.1002/cpe.5215>.
- [12] D. K. Jang, D. Yang, D. Y. Jang, B. Choi, T. Jin, and S. H. Lee, "MOVIN: Real-time Motion Capture using a Single LiDAR," *Computer Graphics Forum*, vol., pp. 12, 2023. <https://doi.org/10.1111/cgf.14961>.
- [13] D. A. Kumar, A. Sastry, P. V. V. Kishore, and E. K. Kumar, "3D sign language recognition using spatio-temporal graph kernels," *Journal of King Saud University-Computer and Information Sciences*, vol. 34, no. 2, pp. 143-152, 2022. <https://doi.org/10.1016/j.jksuci.2018.11.008>.
- [14] E. K. Kumar, P. V. V. Kishore, D. A. Kumar, and M. T. K. Kumar, "Early estimation model for 3D-discrete indian sign language recognition using graph matching," *Journal of King Saud University-Computer and Information Sciences*, vol. 33, no. 7, pp. 852-864, 2021. <https://doi.org/10.1016/j.jksuci.2018.06.008>.
- [15] T. Kyriakou, M. A. D. Crespo, A. Panayiotou, Y. Chrysanthou, P. Charalambous, and A. Aristidou, "Virtual Instrument Performances (VIP): A Comprehensive Review," *Computer Graphics Forum*, vol. 43, no. 2, pp. 29, 2024. <https://doi.org/10.1111/cgf.15065>.
- [16] N. Lannan, L. Zhou, and G. L. Fan, "Human Motion Enhancement via Tobit Kalman Filter-Assisted Autoencoder," *Ieee Access*, vol. 10, pp. 29233-29251, 2022. <https://doi.org/10.1109/access.2022.3157605>.
- [17] J. Li and D. P. Gu, "Research on basketball players'action recognition based on interactive system and machine learning," *Journal of Intelligent & Fuzzy Systems*, vol. 40, no. 2, pp. 2029-2039, 2021. <https://doi.org/10.3233/jifs-189205>.
- [18] J. Li et al., "Real-Time Human Motion Capture Based on Wearable Inertial Sensor Networks," *Ieee Internet of Things Journal*, vol. 9, no. 11, pp. 8953-8966, 2022. <https://doi.org/10.1109/jiot.2021.3119328>.
- [19] W. Y. Li, Y. Q. Zeng, Q. Zhang, Y. L. Wu, and G. M. Chen, "Human Motion Capture Based on Incremental Dimension Reduction and Projection Position Optimization," *Wireless Communications & Mobile Computing*, vol. 2021, pp. 9, 2021. <https://doi.org/10.1155/2021/5589100>.
- [20] Z. L. Li, L. T. Wang, and X. Q. Wu, "Artificial intelligence based virtual gaming experience for sports training and simulation of human motion trajectory capture," *Entertainment Computing*, vol. 52, pp. 9, 2025. <https://doi.org/10.1016/j.entcom.2024.100828>.
- [21] J. J. Lin and J. Song, "Design of motion capture system in physical education teaching based on machine vision," *Soft Computing*, vol., pp. 10, 2023. <https://doi.org/10.1007/s00500-023-08779-5>.
- [22] Khalique A S, Amril N, Imran K, et al. "Short term energy consumption forecasting using neural basis expansion analysis for interpretable time series," *Scientific Reports*, vol. 12, no. 1, 2022. <https://doi.org/10.1038/s41598-022-26499-y>.
- [23] B. Y. Ma, Z. N. Jiang, Y. Liu, and Z. W. Xie, "Advances in Space Robots for On-Orbit Servicing: A Comprehensive Review," *Advanced Intelligent Systems*, vol. 5, no. 8, pp. 21, 2023. <https://doi.org/10.1002/aisy.202200397>.
- [24] J. Z. Miao, T. Peng, F. Fang, X. R. Hu, and L. Li, "TDGar-Ani: temporal motion fusion model and deformation correction network for enhancing garment animation details," *Visual Computer*, vol., pp. 15, 2024. <https://doi.org/10.1007/s00371-024-03575-0>.
- [25] S. Mohaoui and A. Dmytryshyn, "CP decomposition-based algorithms for completion problem of motion capture data," *Pattern Analysis and Applications*, vol. 27, no. 4, pp. 19, 2024. <https://doi.org/10.1007/s10044-024-01342-4>.
- [26] R. Monica and J. Aleotti, "Evaluation of the Oculus Rift S tracking system in room scale Virtual Reality," *Virtual Reality*, vol. 26, no. 4, pp. 1335-1345, 2022. <https://doi.org/10.1007/s10055-022-00637-3>.

- [27] X. L. Niu, Z. H. Yang, N. N. Zhou, and C. H. Li, "A novel method for cage whirl motion capture of high-precision bearing inspired by U-Net," *Engineering Applications of Artificial Intelligence*, vol. 117, pp. 16, 2023. <https://doi.org/10.1016/j.engappai.2022.105552>.
- [28] Y. Pan, "Sports game teaching and high precision sports training system based on virtual reality technology," *Entertainment Computing*, vol. 50, pp. 10, 2024. <https://doi.org/10.2139/ssrn.4608235>.
- [29] L. Qiao and J. C. W. Lin, "Research on Standardized Feature Positioning Technology of Motion Amplitude Based on Intelligent Vision," *Mobile Networks & Applications*, vol. 27, no. 6, pp. 2391-2399, 2022. <https://doi.org/10.1007/s11036-021-01883-6>.