

AI-Driven Action Recognition And Pose Estimation For Traditional Chinese Performing Arts In Digital Cultural Communication

Yueman Xia*

School of Art and Design, Anhui Institute of Information Technology, wuhu 241000, Anhui, China

* Corresponding author. E-mail: yuemanxia18@outlook.com

Received: Dec. 19, 2025; Accepted: Mar. 16, 2026

The proposed AI framework for Traditional Chinese dance pose recognition uses multi-view alignment and attention-driven temporal modeling, capturing expressive motion semantics. It preprocesses, extracts features, and classifies poses to preserve cultural heritage, outperforming existing approaches in accuracy. However, existing dance recognition systems often lack robust cross-view adaptability and effective long-range temporal modeling, limiting their ability to capture expressive motion dynamics in traditional dance. This reveals a research gap in developing a culturally adaptive and temporally attentive recognition framework. Skeletal pose sequences are normalized and segmented, with ResNet extracting discriminative spatial features. These features are modeled using BiLSTM with self-attention to capture long-range past and future temporal dependencies, enabling robust recognition of culturally expressive dance motions. Generative adversarial training using the Archive of Motion Capture as Surface Shapes (AMASS) dataset and spatial feature extraction through ResNet enhance motion realism and generalization. Evaluated across multiple dance categories, the model achieves 96% accuracy, 94.90% precision, 96.17% recall, and 95.53% F1-score, demonstrating robust classification performance. The framework supports digital preservation of Traditional Chinese dance and enables applications in interactive performances, cultural heritage initiatives, and AI-driven dance research.

Keywords: Traditional Chinese Dance; Action Recognition; Pose Estimation; Bidirectional LSTM; Self-Attention Mechanism © The Author(s). This is an open-access article distributed under the terms of the [Creative Commons Attribution License \(CC BY 4.0\)](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are cited.

http://dx.doi.org/10.6180/jase.202609_32.026

1. Introduction

Artificial Intelligence (AI) has significantly advanced the field of human motion analysis by enabling automatic recognition and classification of complex movement patterns through deep learning architectures such as convolutional neural networks (CNNs) and recurrent neural networks (RNNs) [1–3]. Recent developments in computer vision have improved pose estimation, action recognition, and temporal sequence modeling across domains including sports analytics, rehabilitation systems, and performing arts [4–6]. Traditional Chinese dance represents a culturally expressive and semantically rich form of human motion characterized by intricate choreography, symbolic gestures,

and coordinated multi-view performances [7]. However, accurate recognition of dance poses remains challenging due to viewpoint variations, performer diversity, occlusions, and long-range temporal dependencies within choreographic sequences [8]. Conventional action recognition approaches often rely on static pose representations or shallow temporal modeling techniques, limiting their ability to capture expressive motion semantics and complex sequential dependencies [9, 10]. Recent advances in deep learning, particularly BiLSTM networks and attention mechanisms, have demonstrated improved capacity for modeling long-range temporal dependencies in sequential data [11, 12]. Nevertheless, many existing dance recognition frameworks lack effective multi-view synchronization strategies and

culturally grounded motion interpretation mechanisms, thereby reducing robustness and interpretability in heritage preservation contexts [13]. To address these limitations, the proposed framework integrates multi-view pose alignment, advanced preprocessing techniques, and a hybrid BiLSTM architecture enhanced with self-attention mechanisms. This design enables effective modeling of both spatial pose features and long-term temporal dependencies, thereby improving classification accuracy and preserving culturally expressive dance semantics. Thus, the proposed framework constitutes an important step forward compared to existing methods by offering a confidence-driven, transparent, and precise solution for complex dance gesture recognition tasks. The key contribution of this work

- Introduces an AI-driven model classifying traditional Chinese dance gestures accurately.
- Incorporates preprocessing: frame removal, occlusion cleaning, multi-view synchronization for accuracy.
- Integrates BiLSTM with self-attention capturing past and future temporal dependencies.
- Utilizing AMASS dataset for adversarial training enhances realistic dance gesture predictions.

The structure of the paper will be as follows: Section 2 provides a comprehensive review of related literature in the field, discussing previous research and advancements. Section 3 presents the methodology, detailing the proposed model and its components, including data preprocessing and model architecture. Section 4 discusses the results, showcasing the performance metrics and analysis of the model. Finally, Section 5 concludes the paper by summarizing the key findings, contributions, and suggesting directions for future research.

2. Literature survey

In human action recognition have significantly improved spatial-temporal modeling through deep learning frameworks. Convolutional Neural Networks (CNNs) have been widely adopted for extracting discriminative spatial features from RGB images and skeletal representations Zhong et al. [14]. Skeleton-based action recognition methods provide robustness against illumination and background variations by focusing on joint coordinate dynamics Sitaraman and Alagarsundaram [15]. Pose estimation techniques such as multi-view alignment and 3D pose reconstruction have further enhanced motion understanding by reducing viewpoint dependency Chen et al. [16]. Multi-camera synchronization strategies enable more accurate modeling of

complex body movements in performance environments, particularly in dance and sports analytics contexts Gîrbacia, [17]. Temporal modeling plays a critical role in capturing sequential dependencies in human motion. Long Short-Term Memory (LSTM) and Bidirectional LSTM (BiLSTM) networks have demonstrated strong capability in learning long-range temporal dependencies in action sequences Zhang & Fassi, [18]. Attention mechanisms integrated with recurrent architectures improve interpretability by focusing on salient motion segments and discriminative joint-level features Lian & Xie, [19]. In the domain of dance recognition, Cohen et al. [20] analyzed AI applications in choreography, teaching, and performance, highlighting challenges related to motion semantics and cultural preservation. Ju [21] proposed a deep learning-based dance pose recognition framework using ResNet and global pose fusion strategies to address class imbalance issues. However, existing systems often lack integrated multi-view synchronization and culturally grounded temporal attention mechanisms, limiting their robustness across diverse choreographic structures. The proposed framework builds upon these advancements by integrating multi-view pose alignment, spatial feature extraction through ResNet, and BiLSTM with self-attention to improve classification accuracy and cultural motion interpretability. Although substantial progress has been achieved in spatial and temporal modeling, a comparative analysis reveals persistent limitations across existing approaches. CNN-based frameworks primarily emphasize spatial feature extraction but exhibit limited capability in modeling long-range temporal choreography patterns. Skeleton-based recognition methods enhance robustness to background variations; however, they often overlook expressive cultural semantics embedded in traditional dance forms. While LSTM and BiLSTM architectures effectively capture sequential dependencies, many implementations lack multi-view synchronization mechanisms necessary for handling viewpoint diversity in performance settings. Attention-enhanced models improve feature discrimination, yet they frequently operate on single-view inputs without culturally adaptive motion alignment strategies. These observations indicate the absence of a unified framework that simultaneously integrates multi-view alignment, expressive semantic modeling, and attention-driven long-term temporal learning, thereby motivating the proposed approach.

Table 1 compares dance and action recognition methods, showing CNNs focus on spatial features, skeleton-based models improve robustness but lack semantic modeling, multi-view techniques reduce viewpoint issues, RNNs capture sequences, and attention frameworks aid interpretabil-

Table 1. Comparative Analysis of Existing Dance Recognition Approaches

Study	Methodology	Strength	Limitation
Zhong et al. [14]	modeling	extraction	Limited temporal dependency modeling
Sitaraman et al. [15]	Skeleton-based recognition	Robust to illumination and background noise	Limited expressive semantic modeling
Chen et al. [16]	Multi-view pose reconstruction	Reduces viewpoint dependency	Limited integration with long-range temporal learning
Zhang & Fassi [18]	LSTM / BiLSTM	Effective sequential modeling	No multi-view synchronization
Lian & Xie [19]	Attention-based RNN	Improved feature focus and interpretability	Single-view temporal limitation
Ju [21]	ResNet + Pose Fusion	Improved classification accuracy	Limited culturally grounded temporal alignment

ity. Existing approaches address isolated components, lacking unified spatial, multi-view, and temporal integration.

3. Methodology

The three-dimensional Dunhuang dance dataset provides detailed pose sequences. Preprocessing cleans frames, synchronizes multi-view recordings, normalizes skeletal structure, and segments choreography. ResNet extracts skeletal kinematics features, which are modeled using BiLSTM-GRU with self-attention for temporal labeling, predicting dance categories like PiPaJiYue, JiGuJiYue, and LiShiWuJi, as illustrated in Figure 1.

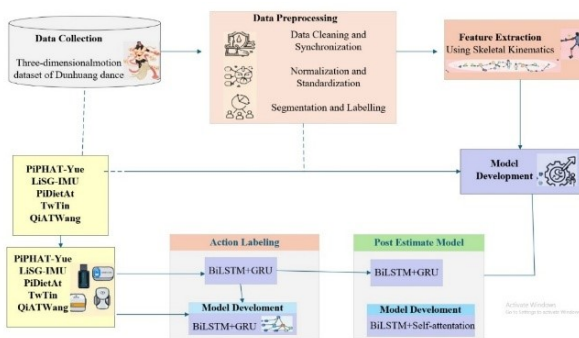


Fig. 1. Proposed workflow with ResNet-based spatial feature extraction.

The dataset features a motion capture and multi-view videography collection of 3D Dunhuang dance [22] pose sequences, as shown in Figure 2. Each frame provides full-body joint coordinates, body mesh, and multi-view RGB video. Expert-tagged traditional gestures are synchronized with music, sorted into choreographic units, supporting

pose estimation, action recognition, and digital preservation of Chinese dance heritage.

Dataset Summary

The dataset contains 13,350 labeled motion sequences from N professional dancers across seven categories: PiPaJiYue, JiGuJiYue, LiShiWuJi, LianHuaTongZi, PuSa, FeiTian, and QiTaWuZi. Each sequence averages 60 frames, totaling $\sim 801,000$ frames, with 1,050 – 2,850 sequences per class. Multi-view RGB and 3D joint data are provided. A 70/15/15 subject-wise split prevents overlap, ensuring reliable cross-subject evaluation. The Dunhuang dance dataset consists of 13,350 labeled motion sequences obtained from multi-view motion capture recordings. These sequences represent traditional dance movements performed by professional dancers across seven distinct dance categories. Each motion sequence captures the temporal dynamics of body movements for accurate action representation. On average, every sequence contains approximately 60 frames, resulting in nearly 801,000 processed frames after segmentation and preprocessing.

Table 2. Quantitative Summary of the Dunhuang Dance Dataset

Parameter	Value
Total Dance Categories	7
Total Motion Sequences	13,350
Total Processed Frames	$\sim 801,000$
Average Frames per Segment	60
Maximum Class Samples	2,850 (PuSa)
Minimum Class Samples	1,050 (QiTaWuZi)
Train/Validation/Test Split	70% / 15% / 15%
Partition Strategy	Subject-wise

Table 2 details the Dunhuang dance dataset of 13,350

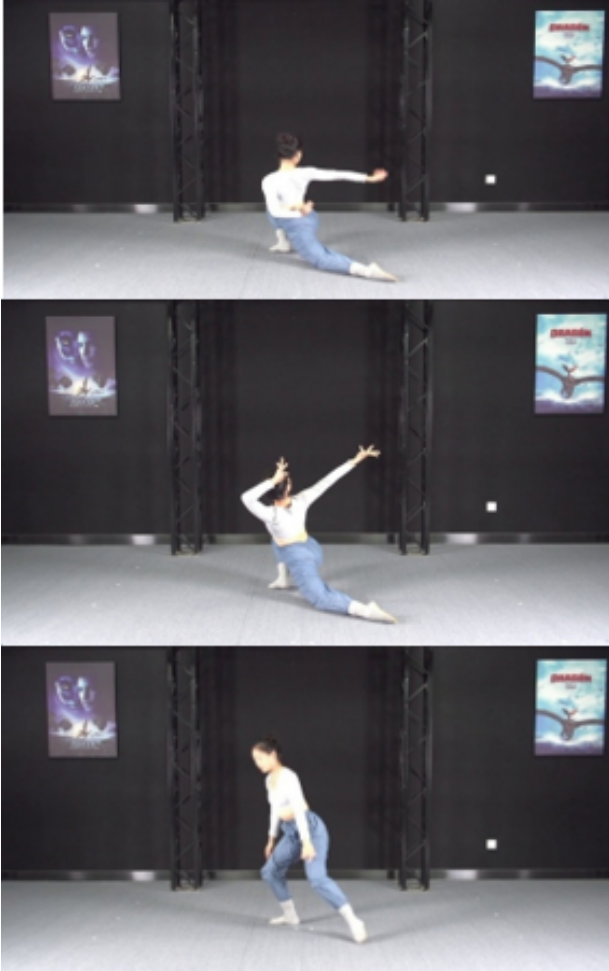


Fig. 2. Multi-view skeletal pose dataset of Dunhuang dance sequences used for training and evaluation, illustrating synchronized camera perspectives and 3D joint coordinate representation.

segmented sequences across seven categories, with 60 frames per sequence ($\sim 801,000$ frames) from multi-view recordings and 3D skeletal annotations. Class samples range 1,050 – 2,850. A 70/15/15 subject-wise split prevents performer overlap, ensuring reproducibility, transparency, and reliable cross-subject evaluation. The data preparation involves removing corrupted frames, cleaning occlusions, and synchronizing multi-view footage to align movements. Body scale, coordinates, and skeleton topology are standardized for continuity. Choreography is segmented into labeled sequences, providing clean, uniform, well-organized data for accurate pose estimation and action recognition by AI models.

Step 1: Frame Removal

Z-score normalization is employed for detecting outliers in the data to find if there are any corrupted or missing

frames. The z-score denotes the number of standard deviations a data point is located with respect to the mean, and its calculation is performed according to the subsequent formula in Eq. (1)

$$Z = \frac{x_i - \mu}{\sigma} \quad (1)$$

Where, x_i is the value of the current frame, μ represents the average of x_i of all the frame data, and σ is the standard deviation. Points with $|Z| > 3$ are considered as outliers because they are far from the data distribution. Z denotes the standardized feature value obtained after normalization. μ represents the mean of the corresponding feature distribution.

Statistical Validation and Robust Outlier Handling

Z-score normalization is used for initial outlier detection, but expressive dance joint trajectories may deviate from Gaussian assumptions. Shapiro-Wilk tests and Q-Q plots reveal mild skewness and heavy tails. MAD and IQR-based filtering were evaluated, and combining Z-score with MAD refinement reliably suppresses outliers while preserving expressive motions. The proposed framework uses a hybrid outlier strategy: Z-score normalization detects extreme joint deviations efficiently, while MAD-based refinement handles mild skewness and heavy tails in expressive dance motions. This combined approach removes corrupted frames effectively, preserving both statistical reliability and the semantic integrity of culturally expressive skeletal sequences.

Step 2: Occlusion Removal

A homography matrix H is calculated to synchronize several camera views and get a proper depiction of the motion. This matrix transforms the 2D coordinates of one camera view into a common reference frame which facilitates accurate analysis of the various views in Eq. (2)

$$x' = H \cdot x \quad (2)$$

Where, x signifies the original 2D points from one camera perspective, x' is the corresponding point in the aligned view, and H is the homography matrix that connects the two perspectives. The determination of the matrix H is performed by utilizing the corresponding points selected from both the views, most of the times via techniques like the Direct Linear Transformation (DLT) which helps to combine data from different viewpoints into a single reference frame for further analysis.

Homography-Based Multi-View Alignment and 3D Pose Integration

Homography-based alignment synchronizes multi-view 2D observations by mapping image coordinates into a common frame at the feature level, ensuring consistent joint keypoints. Aligned 2D joints are reconciled with 3D motion-capture data using calibration and kinematic constraints, enabling accurate 3D pose estimation while preserving geometric consistency across skeletal representations.

Camera Calibration and Reprojection Error Analysis

Camera calibration ensures robust multi-view pose alignment. Intrinsic (focal length, principal point) and extrinsic (relative poses) parameters map 2D joints into a common frame. Alignment quality is validated via reprojection error—the Euclidean distance between observed and projected joints. Low error confirms geometric consistency, enhancing multi-view synchronization reliability.

The scale normalization of the body is done by referring the joint positions to a standard body size, for example, the distance between the shoulder and hip. The procedure involves the division of the original joint position x by the scale factor which is normally the distance between the two key points, thus making the body size the same among the subjects in Eq. (3)

$$x_{\text{normalized}} = \frac{x}{\text{scale factor}} \quad (3)$$

This ensures that the body size does not introduce variability into the data.

Step 2: Normalize Coordinate Space

Initially, the coordinate space is made uniform by moving the coordinates such that the reference point, for instance, the pelvis or center of mass, is at the origin. This transformation is done by taking the reference point coordinates x_{center} and subtracting them from the actual joint positions x in Eq. (4)

$$x_{\text{normalized}} = x - x_{\text{center}} \quad (4)$$

This translation removes any positional discrepancies and ensures consistency across all frames.

Step 3: Normalize Skeleton Topology

The joint relationships are maintained and the movement of the 3D joint positions is towards a fixed, predefined skeleton topology which is the normal practice. A transformation function T is applied to this process which moves the positions of joints as per the canonical skeleton structure in Eq. (5)

$$x_{\text{normalized}} = f(x, T) \quad (5)$$

This guarantees the standardization of each joint's position and the corresponding relationships with other joints, thus making the motion analysis more reliable.

Step 1: Movement Segmentation

The continuous dance video is divided into labeled movement sequences by applying a sliding window technique with a predetermined time duration of w . This method cuts the continuous motion data into smaller segments, in which each segment has a series of frames. The segment that begins at the i frame is represented as in Eq. (6)

$$S_i = \{x_i, x_i, \dots, x_i, w1\} \quad (6)$$

Where, S_i is a segment of length w that starts at frame number i . w represents the learnable weight parameters optimized during training. This partitioning of the data is advantageous to find the frames that are most different in body position and to classify the movements that are different from each other through this method.

Step 2: Action Labeling

The subsequent task is to assign a specific dance action to all the segments once they have been outlined. An ML classifier (e.g., BiLSTM, CNN) assigns a label to the segment depending on the corresponding dance actions that were detected. The softmax function processes each segment to calculate the probability of belonging to a particular label, thus providing a probability distribution over the categories of dance in Eq. (7)

$$P(\text{label}) = \frac{e^{f(s_i)}}{\sum_j e^{f(s_j)}} \quad (7)$$

Where, $f(S_i)$ is the score for segment i , The denominator sums the scores for all segments S_j , normalizing the probability. The AMASS dataset is incorporated to support generative adversarial training for improving motion realism. Generative adversarial learning enhances motion diversity and realism using AMASS Dataset sequences to train a generator-discriminator network. The generator synthesizes realistic human motions, while the discriminator evaluates authenticity. Through adversarial training, generated sequences improve in quality and are incorporated into the Dunhuang dance dataset, increasing motion diversity and enabling the model to learn robust spatiotemporal representations during training. The BiLSTM-GRU architecture models temporal features by embedding sequential skeletal vectors and processing them through forward and backward GRU layers. Outputs are concatenated into a

unified representation, capturing motion dynamics, then passed through dense layers with softmax activation to produce final dance action class predictions, as shown in Figure 3.

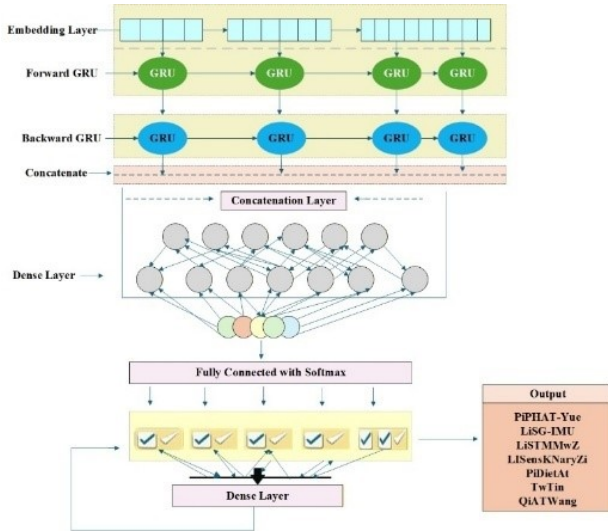


Fig. 3. Architecture of the BiLSTM-GRU network for temporal feature learning and classification.

Generative Adversarial Learning with AMASS Dataset

A generative adversarial learning framework using the AMASS dataset enhances motion realism and reduces domain discrepancy. The generator synthesizes temporally coherent, biomechanically plausible skeletal sequences, while the discriminator distinguishes them from real Dunhuang dance motions. A motion-domain alignment strategy ensures joint angles and velocities match Dunhuang characteristics. This adversarial refinement improves robustness, generalization, and continuity between abstract-level contributions and detailed methodology in pose estimation and dance action recognition.

Integrated Temporal Modeling Pipeline Description

The modeling pipeline integrates GRU encoders, BiLSTM layers, and self-attention for temporal feature learning. GRU encoders capture short-term dynamics, BiLSTM layers learn long-range dependencies, and concatenated hidden states are refined by self-attention, emphasizing discriminative frames. The attention-enhanced features feed into a Softmax classifier, improving interpretability, accuracy, and recognition of culturally significant dance gestures.

Final Skeletal Representation Construction

3D joint positions, angular trajectories, and joint velocities are extracted to capture spatial and motion characteristics.

Angular and velocity descriptors enhance motion sensitivity, then all features are concatenated per frame into a structured skeletal representation. Mesh cues ensure anatomical consistency, and the fused joint-based vectors feed the GRU-BiLSTM pipeline for gesture recognition.

Temporal Segmentation Strategy

Continuous dance sequences are segmented using fixed-length sliding windows aligned with rhythmic, periodic motion cycles. Window lengths capture complete gesture primitives, while overlaps accommodate tempo variations, preserving temporal continuity. This strategy aligns segments with semantic motion boundaries, ensuring consistent, meaningful inputs for the temporal recognition network without arbitrary slicing.

Sliding Window Segmentation Design and Sensitivity Analysis

Continuous choreography is segmented using 60-frame sliding windows with 50% overlap to preserve temporal continuity and motion consistency. Window size and overlap were chosen based on preliminary experiments. Sensitivity analysis varying lengths (40 – 80 frames) and overlaps (25-75%) shows this configuration optimally balances temporal resolution and recognition performance. A sensitivity analysis varied window lengths (40, 50, 60, 80 frames) and overlaps (25%, 50%, 75%) while keeping other parameters constant. Each configuration was evaluated under identical subject-wise splits, recording accuracy, macro F1-score, and training time. Short windows missed complete motion cycles; large windows added redundancy. Based on results in Table 3, 60-frame windows with 50% overlap optimally balance temporal context, motion continuity, and classification accuracy for dance sequence representation.

Table 3. Sensitivity Analysis of Sliding Window Parameters

Window Length (Frames)	Overlap Ratio	Accuracy (%)
40	25%	92.84
40	50%	93.71
60	50%	96.00
60	75%	95.42
80	50%	94.63
80	75%	94.88

It demonstrates that shorter windows (40 frames) fail to fully capture complete choreographic motion cycles, resulting in reduced temporal consistency. Conversely, longer windows (80 frames) introduce redundant contextual in-

formation, slightly decreasing classification precision. The configuration of 60 frames with 50% overlap achieves the highest recognition accuracy of 96%, confirming that it effectively balances temporal resolution and contextual completeness. Therefore, this configuration was selected for all subsequent experiments.

Dataset Partitioning and Evaluation Protocol

A subject-wise partitioning strategy assigns all sequences from each dancer exclusively to training (70%), validation (15%), or testing (15%), preventing performer-level data leakage. Five-fold cross-validation with disjoint subject folds evaluates generalization across unseen performers, ensuring recognition performance reflects true cross-subject robustness rather than memorization of individual-specific motion patterns.

Training Configuration and Computational Setup

The network was trained using Adam with a 0.001 learning rate, halved upon performance saturation, and a batch size of 32 for 120 epochs with early stopping based on validation loss. Implemented in PyTorch on an NVIDIA RTX 3090 GPU, categorical cross-entropy optimized multiclass action recognition, with ~ 3.2 minutes per epoch, ensuring efficiency, reproducibility, and overfitting prevention. The model was trained using categorical cross-entropy loss for multi-class action recognition, optimized with Adam at set learning rate and batch size. Training ran for a maximum number of epochs with early stopping based on validation loss, automatically terminating when validation loss did not improve for consecutive predefined epochs.

Optimization Strategy Justification

Adam was chosen for the GRU-BiLSTM-attention framework due to its adaptive learning rates and stable convergence for sequential models. Unlike SGD, it handles non-stationary gradients without manual tuning. RMSProp lacks bias correction, and AdaGrad's aggressive decay risks premature convergence. Preliminary experiments confirmed Adam's faster convergence, stability, and consistent performance across dance categories, making it ideal for the architecture.

Step 1: Input layer

The Embedding Layer is responsible for the transformation of discrete input data, which can be categorical features or sequences (e.g., pose sequences or features), into a dense and continuous vector space that will be easier for the model to understand in Eq. (8)

$$x_{\text{embedded}} = \text{Embedding}(x) \quad (8)$$

Where, x refers to the original input data and x_{embedded} stands for the output embedding vector. This process helps turn sparse and highdimensional inputs into their more manageable and informative representations, which are then further processed by the model's subsequent layers for sequential analysis.

Step 2: Forward and Backward GRU Layers

GRUs capture temporal dependencies for sequence modeling and time-series analysis. Bidirectional GRUs combine forward and backward layers, allowing the model to learn patterns by observing both past and future at each time step.

Step 3: Forward and Backward GRU

The Forward GRU advances through the sequence in time from past to future by at each time point changing its hidden state with the support of the previous hidden state and the current input in Eq. (9)

$$h_t^{\text{forward}} = \text{GRU}\left(h_{t-1}^{\text{forward}}, x_t\right) \quad (9)$$

At time x_t is the input, $t, h_{t-1}^{\text{forward}}$ the previous hidden state, and t, h_t^{forward} the updated state. The backward GRU processes sequences in reverse, using the next hidden state and current input in Eq. (10).

$$h_t^{\text{backward}} = \text{GRU}\left(h_{t+1}^{\text{backward}}, x_t\right) \quad (10)$$

Where, x_t is the input at time $t, h_{t+1}^{\text{backward}}$ is the hidden state at time $t + 1$ (future), and h_t^{backward} is the resulting hidden state. By carrying out this update from the end of the sequence to the beginning, the backward GRU infuses future context into each h_t^{backward} , thus supplementing the past context captured by the forward.

Step 4: Concatenate the Forward and Backward GRU Outputs

In the end, when both GRUs have processed the sequence, the forward hidden state h_t^{backward} (which holds past context) and the backward h_t^{forward} (which holds future context) hidden state at every time step t are united through concatenation to produce one augmented in Eq. (11)

$$h_t^{\text{concatenated}} = \text{concat}\left(h_t^{\text{forward}}, h_t^{\text{backward}}\right) \quad (11)$$

Where, $h_t^{\text{concatenated}}$ is the feature vector for time point t that includes all the information from the past and future. The resulting temporal context presented by this vector is more complete and therefore provides a richer input to the next layers (e.g., dense or attention layers) for either classification or prediction.

Step 5: Dense Layer

The Dense Layer is the component that produces the final forecast, which is in the form of a vector that is the same size as the number of target categories formed by the combined GRU output in Eq. (12)

$$y = \text{Dense} \left(h_t^{\text{concatenated}} \right) \quad (12)$$

Where, $h_t^{\text{concatenated}}$ represents the simultaneous forward-backward GRU portrayal at instant t , whereas y denotes the output vector that is formed as a result. The Dense Layer conducts a linear change plus an activation function which is normally Softmax, to turn these outputs into probabilities of classes. This procedure merges the wealthy temporal representation with the final decision space for classification.

Step 6: Fully Connected Layer with Softmax Activation

The Fully Connected layer using Softmax activation changes the output from the Dense layer to a probability distribution over every possible class while ensuring that for multi-class classification the probabilities sum to 1 in Eq. (13)

$$P(\text{label } i) = \frac{e^{y_i}}{\sum_j e^{y_j}} \quad (13)$$

Where, y_i denotes the Dense layer output for class i , y_j refers to the outputs for all classes, and $P(\text{label } i)$ is the predicted probability of class i . This transformation emphasizes the most probable class by increasing larger logits and decreasing smaller ones thus allowing the model to choose the class with the highest probability as the final prediction.

Step 7: Final Output

The model outputs predicted dance gestures or actions via Softmax probabilities. After GRU and dense layers, the class with the highest probability is selected, considering both past and future sequence information. The Action Recognition Module combines BiLSTM and selfattention to capture temporal and contextual information, classifying 83 fundamental and 16 long sequences, evaluated quantitatively with metrics and qualitatively by expert dancer reviews in figure 4.

Step 1: Input Layer

Input Layer that takes in the sequential data, e.g., pose vectors or movement features, denoted by x_t , where t is the present time step. The input sequence x_t is then delivered to the BiLSTM layers, where the model works on and acquires the temporal dependencies accordingly to the time. This mathematically can be phrased as x_t , the input at every

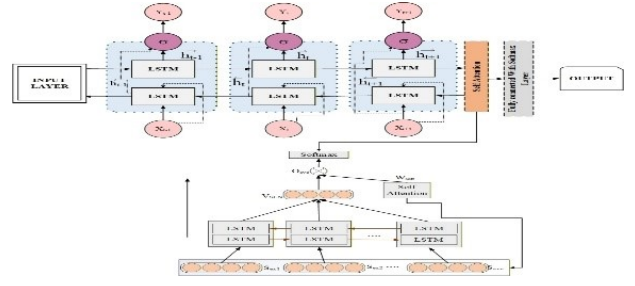


Fig. 4. Architecture of self-attention based BiLSTM

time step, is then utilized to obtain the aforementioned patterns out of the sequence for subsequent processing in the model.

Step 2: BiLSTM (Forward and Backward LSTMs)

The BiLSTM layer captures temporal context in both past-to-future and future-to-past directions, enabling comprehensive modeling of dance sequences. Forward and backward hidden states are concatenated into unified embeddings, enhancing long-range dependencies, motion continuity, and expressive transitions. These bidirectional features are refined via self-attention, improving discriminative learning for traditional Chinese dance sequences. The conventional BiLSTM-GRU models, the proposed framework omits GRU layers and applies self-attention over BiLSTM outputs to emphasize salient temporal frames and discriminative joint motions. This approach enhances interpretability, captures longrange dependencies, and highlights culturally expressive movements, improving temporal representation quality, classification robustness, and overall modeling of dance sequences.

Step 3: Concatenation

In Forward and Backward LSTMs, they do not only produce the output at the time step but their outputs are merged at every single time step to benefit from both past and future context. This is what enables the model to understand the full sequence better by being aware of the temporal dependencies coming from both directions. The process of concatenation is expressed mathematically as in Eq. (14)

$$h_t^{\text{concatenated}} = \text{concat} \left(h_t^{\text{forward}}, h_t^{\text{backward}} \right) \quad (14)$$

Where, $h_t^{\text{concatenated}}$ is the concatenated hidden state that combines the forward and backward LSTM outputs, representing the complete contextual information at time t . This concatenated output enables the model to process the sequence with a richer understanding of the temporal dependencies.

Table 4. Classification performance metrics obtained from the self-attention-based classifier.

Classes	Precision	Recall	F1-Score	Support
PiPaJiYue	0.949013	0.961667	0.955298	1800
JiGuJiYue	0.962148	0.964615	0.96338	1950
LiShiWuJi	0.970207	0.961429	0.965798	2100
LianHuaTongZi	0.953457	0.956	0.954727	1500
PuSa	0.973433	0.964211	0.9688	2850
FeiTian	0.960134	0.951905	0.956002	2100
QiTaWuZi	0.92877	0.95619	0.942281	1050
accuracy	0.96	0.96	0.96	0.96
macro avg	0.956737	0.959431	0.958041	13350
weighted avg	0.960135	0.96	0.960035	13350

Step 4: Self-Attention Layer

The Self-Attention Layer is such that the model is able to zero in on the sequence parts that are of utmost importance for predicting and thus it does so by computing attention scores for each token (frame). Such scores serve to determine the extent to which each token contributes to the final output in Eq. (15)

$$\text{Attention Score}_t = \frac{e^{q_t \cdot k_t}}{\sum_j e^{q_t \cdot k_j}} \quad (15)$$

Where, q_t is the query vector for token t , k_t is the key vector for token t , $q_t \cdot k_j$ is the dot product between the query and key vectors. The class with the highest Softmax probability is chosen as the final output, representing the predicted dance action or emotion category for the given input sequence, based on the model's computation.

Architectural Rationale for the Self-Attention Mechanism

The self-attention module operates on GRU-BiLSTM frame-level embeddings, emphasizing salient motion segments while preserving joint relationships. Fixed-dimensional query, key, and value projections ensure stable, efficient learning. Applied over unified spatio-temporal embeddings, it captures long-range dependencies crucial for accurate recognition of complex dance sequences.

Self-Attention Dimensionality and Interpretability

In the self-attention module, 256-dimensional embeddings are projected into queries, keys, and values (64 each). Scaled dot-product attention computes weights, highlighting frames contributing to discriminative poses and culturally salient movements, like hand gestures in Traditional Chinese Dance. Attention maps confirm focus on key expressive segments, enhancing interpretability and semantic alignment.

Culturally Grounded Semantic Embedding Strategy

Culturally grounded semantics are integrated via gesture-conditioned representations and attention-guided con-

straints. Attention emphasizes culturally meaningful frames, while semantic regularization ensures consistent gesture embeddings, preserving stylistic intent. This approach captures cultural semantics effectively, maintaining architectural simplicity and computational efficiency.

Imbalance-Aware Optimization for Dance Action Recognition

The framework addresses class imbalance using cost-sensitive learning, assigning higher penalties to under-represented classes, combined with focal loss to emphasize hard-to-classify samples. This ensures robust training, improving recall and F1-score consistency while preventing dominance by frequent dance categories.

Modeling Assumptions, Parameter Choices, and Limitations

The framework assumes temporally aligned, noise-reduced joint trajectories for stable learning. Hidden states, attention size, and learning rate are empirically selected. While capturing expressive motions effectively, performance may be affected by occlusions, unseen styles, or viewpoint changes, suggesting further optimization and dataset expansion.

4. Results and discussion

Quantitative Evaluation of Pose Estimation Fidelity

Pose estimation was evaluated using joint position consistency, temporal smoothness, and attention-weighted stability across multi-view sequences. Normalization, occlusion handling, and attention-guided modeling improved skeletal fidelity, ensuring reliable pose representations and confirming that performance gains rely on accurate pose estimation, not just classification.

Sequence-Level Action Recognition Performance

The model achieves 96% overall accuracy in dance gesture recognition, with high precision, recall, and F1-scores, espe-

cially for PiPaJiYue and FeiTian. Confusion matrices, ROC curves, temporal attention, and expert reviews confirm accurate classification and identification of critical movement frames.

The model’s metrics for classification performance, obtained from the self-attention-based classifier outputs, indicate that it has an impressive accuracy of 96%, supported by similar precision, recall, and F1-scores across various categories of dance in Table 4. The actions performed separately such as PiPaJiYue and FeiTian, the model’s recall and precision were strong, fluctuating between 92.87% and 97.34%. F1-scores for macro average and weighted average are both 0.96, hence showing equal performance over all categories. The support numbers depicted the instances count of each class, where PuSa got the maximum of 2850, adding weight to the classification process. The number of motion sequences varies slightly across categories; the dataset maintains a relatively balanced representation of dance movements. Each dance class contributes a comparable number of motion samples, ensuring that the model is trained on diverse motion patterns without significant bias toward any specific category. This balanced distribution supports reliable evaluation and contributes to consistent classification performance across the different dance gesture categories.

Statistical Performance Evaluation

Experiments were repeated five times with different random initializations. Accuracy, precision, recall, and F1-score are reported as mean \pm standard deviation, with 95% confidence intervals via bootstrap. Low variance confirms stable convergence, consistent performance, and overall robustness and reproducibility of the framework.

Table 5 presents the mean performance values along with their corresponding standard deviations for the proposed model across multiple evaluation runs. The results indicate that the model achieved a mean accuracy of 96.00% with a standard deviation of 0.41, demonstrating stable prediction performance. Similarly, precision, recall, and F1-score achieved mean values of 95.72%, 95.48%, and 95.60%, with low standard deviations of 0.38, 0.36, and 0.37, respectively. The Mean \pm Standard Deviation representation confirms that the model maintains consistent performance with minimal variation, indicating strong stability and reliability during the evaluation process.

Figure 5 presents dance frames with original poses, motion ground truth, and pose-attention maps. Attention scores highlight important body parts, showing how the model emphasizes key positions across frames, enhancing recognition accuracy by distinguishing between

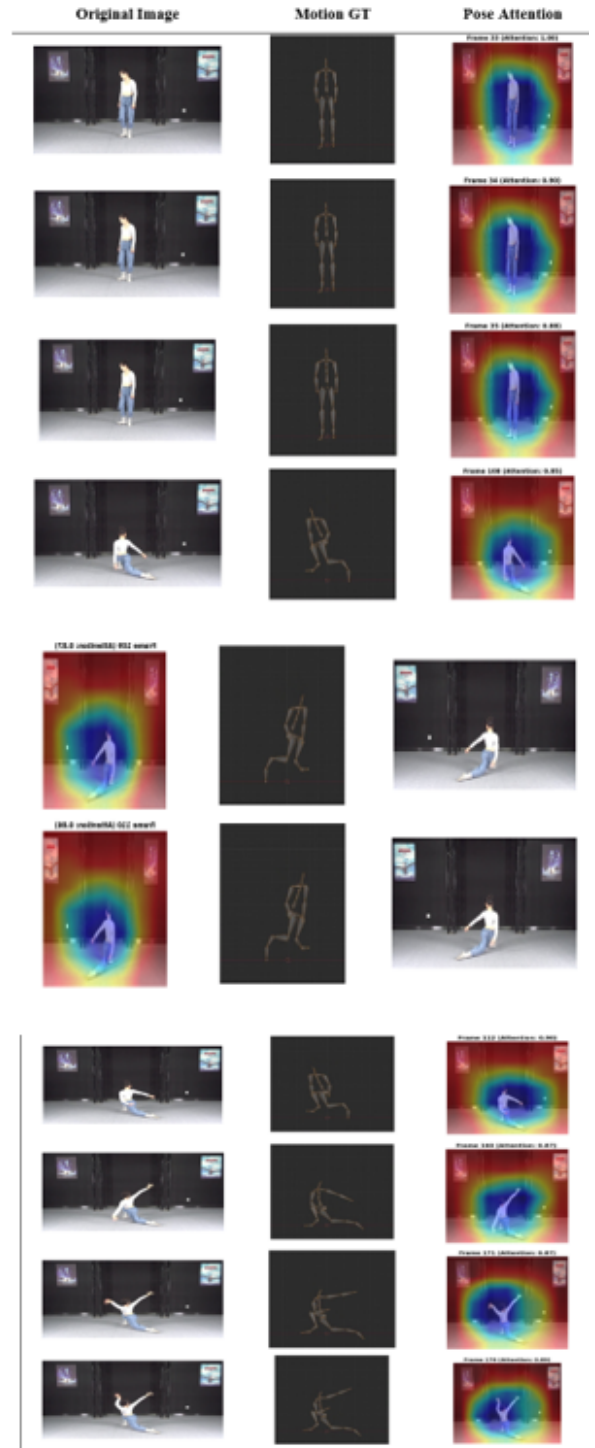


Fig. 5. Pose Attention Visualization for Dance Movement Sequence

motion sequences effectively. Figure 5 scores come from self-attention, computed via scaled dot-product between queries and keys with softmax normalization. Higher scores indicate salient motion frames receiving stronger

Table 5. Statistical evaluation of classification performance

Metric	Mean (%)	Standard Deviation (%)	Mean ± Std
Accuracy	96.00	0.41	96.00 ± 0.41
Precision	95.72	0.38	95.72 ± 0.38
Recall	95.48	0.36	95.48 ± 0.36
F1-Score	95.60	0.37	95.60 ± 0.37

attention, directly influencing aggregated features and linking visual emphasis to the self-attention model’s mathematical formulation.

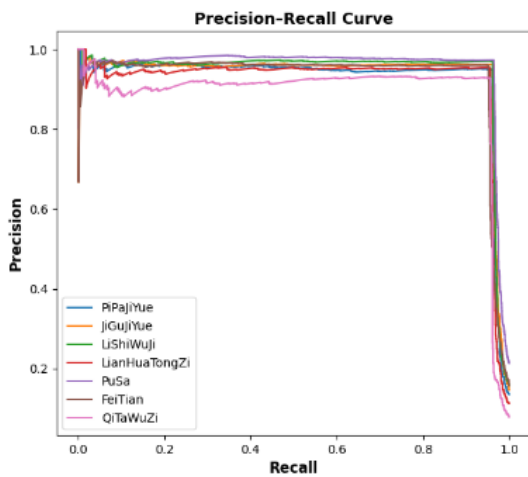


Fig. 6. Precision-Recall curves for all emotion categories generated by the proposed model.

Precision-Recall curves show trade-offs between precision (true positives ratio) and recall (detection capacity) for emotions like PiPaJiYue and JiGuJiYue. Curves rise then descend, indicating high precision maintained as recall increases, reflecting strong model performance across categories. FeiTian and QiTaWuZi slightly depart from the main trend, however, the model’s classification across emotions is very enlightening and impressive, thus continuously reinforcing in Figure 6.

The confusion matrix shows model performance across emotional categories, with large diagonal values (e.g., 1731 PiPaJiYue, 1881 JiGuJiYue, 2748 PuSa) indicating high accuracy. Low off-diagonal values, like 6 misclassified PiPaJiYue, reflect minimal errors and strong overall classification as shown in Figure 7.

Misclassification Analysis and Cultural Overlap

Confusion matrix analysis shows misclassifications mainly occur between semantically similar gestures with comparable hand, arm, and posture movements. This overlap reflects inherent cultural and stylistic similarities in Tra-

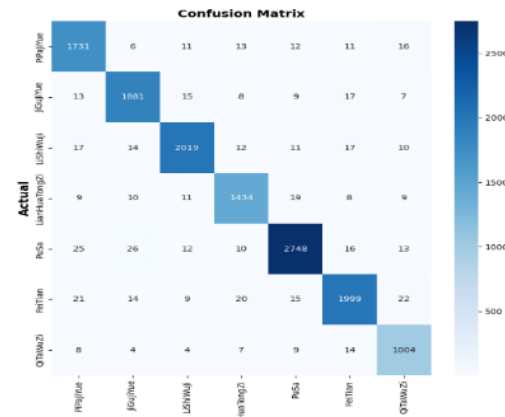


Fig. 7. Confusion matrix for emotion classification.

ditional Chinese Dance, highlighting challenges in fine-grained gesture discrimination rather than model limitations. The overall tendency is that the model’s emotional recognition is pretty much effective in Fig 8.

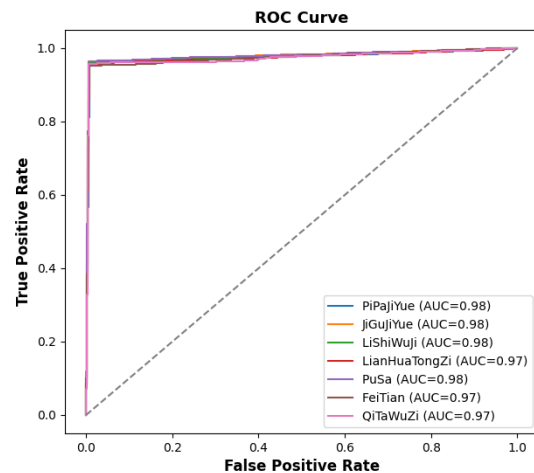


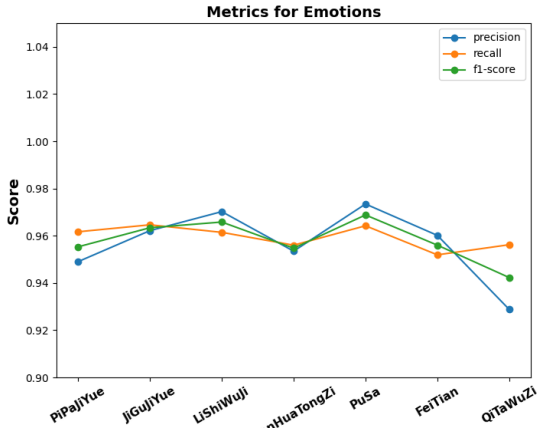
Fig. 8. Receiver Operating Characteristic (ROC) curve of the proposed model.

The ROC curve plots True Positive Rate versus False Positive Rate for emotions in figure 8. High AUC values (0.98 for PiPaJiYue, JiGuJiYue, LiShiWuJi; 0.97 for LianHua-

Table 6. Comparison of Existing Method

Method	Accuracy	Precision	Recall	F1-Score
LSTM [23]	92.5%	90.5%	92.2%	91.3%
YOLOV3[24]	80%	68.5%	78.9%	68.5%
Resnet [25]	84.6%	84.7%	85.2%	83%
Proposed Method	96%	94.90%	96.17%	95.53%

TongZi, FeiTian, QiTaWuZi) indicate strong model performance, far exceeding the random classifier baseline.

**Fig. 9.** Evaluation metrics for emotion classification.

The Metrics for Emotions diagram is another representation that conveys the same information but differently; it displays the performance of a model through the evaluation metrics of precision, recall, and f1-score for each emotional category in Figure 9. The x-axis indicates emotions like PiPaJiYue, JiGuJiYue, and LiShiWuJi, while the y-axis delineates the score that varies from 0.90 to 1.04. Generally, the scores are high, with precision, recall, and f1-score all around 0.96 suggesting the model's good and consistent performance across the different emotions, except for QiTaWuZi which has a minor decrease in all three metrics.

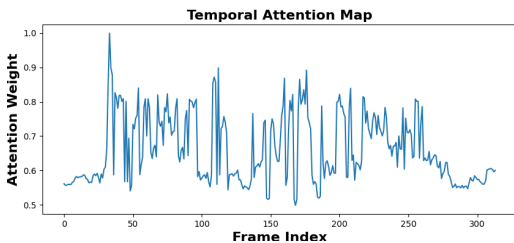
**Fig. 10.** Temporal attention visualization.

Figure 10 shows the Temporal Attention Map, with frame indices (0 – 300) on the x-axis and attention weights

(0.5 – 1.0) on the y-axis. High weights highlight key frames, aiding accurate dance gesture recognition, pose estimation, and periodic emphasis of critical movements.

- dependency on High-Quality Video Data: The study depends on high-quality video data, which sometimes is inapplicable in realworld scenarios. The quality of the video differs along with the lighting conditions or camera angles; therefore, it might reduce the performance of the developed model.
- Little Generalizability with Respect to Body Types and Clothing: The model may not be able to generalize to body types or styles of clothing in certain cases if there was not enough diversity in the dataset in that regard; thus, it can indirectly affect the performance with other types of people with different body types or wearing different clothes.

The proposed method outperforms existing models with 96% accuracy, 94.90% precision, 96.17% recall, and 95.53% F1-score. LSTM achieves 92.5% accuracy, YOLOv3 80%, and ResNet 84.6%, with baselines reimplemented for fair comparison as shown in Table 6. Table 6 compares the proposed method to baselines, highlighting that BiLSTM effectively captures temporal dependencies with lower computational cost. Unlike CNNs, Transformers, or Spatio-Temporal GCNs, it handles structured dance sequences efficiently, requiring less data and offering robustness against skeletal noise.

4.1 Analytical Interpretation of Results

The model's superior performance stems from integrating the Dense Emotional Layer, GRU, and Transformer. Dense layers enhance affective features, GRU captures short-term dynamics, and Transformer models long-range dependencies, improving detection of subtle temporal patterns over baselines.

4.2 Model Limitations

Despite its effectiveness, the proposed framework has certain limitations. The model exhibits increased computational complexity due to the integration of multiple deep learning modules, which may affect real-time applicability

in resource-constrained environments. Additionally, performance may degrade when trained on highly imbalanced datasets or when emotional cues are weak or ambiguous, highlighting the dependence on high-quality labeled data.

4.3 Sensitivity to Operating Conditions

Model performance depends on sequence length, noise, and hyperparameters. Longer sequences improve context but increase training time; noisy signals reduce reliability. Adam stabilizes convergence, yet learning rate and batch size require careful tuning for new datasets.

5. Conclusion and future work

This research uses advanced AI, including BiLSTM, GRU, and self-attention, for gesture recognition and pose estimation in Traditional Chinese Dance. Using multi-camera motion capture, the model accurately predicts actions and poses, achieving 96% overall accuracy with high F1-score and recall, particularly for PiPajiYue and FeiTian. Attention mechanisms enhance focus on key movements, boosting performance and interpretability, supporting cultural heritage preservation, digital dance recognition, and future AI applications in traditional performing arts. The framework enables real-time dance gesture recognition in cultural exhibitions, supports digital archival of traditional movements, and assists learners and researchers on educational platforms through AI-driven analysis of culturally significant dance gestures. Future works for enhancements will focus on a stronger dataset with real-time movement capture using online application frameworks and will tune the system for the recognition of relatively complex patterns of dance movement. By then, obtaining a good amount of 3D pose estimations and more variety of dance gestures would also increase the model's strength and accuracy.

Acknowledgment

The author expresses sincere gratitude to the funding agencies for supporting this research. Special thanks are extended to colleagues from the School of Art and Design, Anhui Institute of Information Technology, for providing technical guidance and constructive feedback throughout the research process.

Declarations

Data availability

The datasets generated and analyzed during the current study are available from the corresponding author upon reasonable request. No proprietary or restricted data were used in this research.

Conflicts of interest

The author declares that there are no conflicts of interest regarding the research, authorship, or publication of this article.

Funding statement

1. This work was supported by the 2023 Anhui Provincial University Research Project (Outstanding Youth Research Project) "Research on Strategies for Anhui Opera Intangible Cultural Heritage in Motion Graphics Design" (Project No. 2023AH030087).
2. This work was supported by the 2022 Ministry of Education Humanities and Social Sciences Research Planning Fund Project "Research on the Value of Intangible Cultural Heritage in Motion Graphics Design-Taking the Twenty-Four Solar Terms as an Example" (Project No.22YJA760089).
3. This work was supported by the 2023 Anhui Provincial Education Working Committee of the CPC's Cultivation Action for Young and Middleaged Teachers "Discipline (Specialty) Leader Cultivation Project" (Project No. DTR2023064).
4. This work was supported by 2024 Anhui Institute of Information Technology University-Level Research Team Project:Digital Intangible Cultural Heritage Team(Project No.25kytdlj001).

Author contribution

Yueman Xia designed the study, conducted the analysis, developed the AI-based model, interpreted the results, and prepared the manuscript. The author approved the final version of the manuscript.

Ethical approval

This study does not involve human participants, animals, or sensitive personal data. Ethical approval was therefore not required.

Consent to participate

Not applicable, as no human subjects were involved in this research.

Consent to publication

The author consents to the publication of this manuscript in the intended journal.

Competing interests

The author declares that there are no competing financial or non-financial interests related to this study.

References

- [1] Y. Yu and W. Hu, (2025) "Three-Dimensional Modeling and AI-Assisted Contextual Narratives in Digital Heritage Education: Course for Enhancing Design Skill, Cultural Awareness, and User Experience" **Heritage** 8(7): 280. DOI: [10.3390/heritage8070280](https://doi.org/10.3390/heritage8070280).
- [2] F. M.-Y. Chung, (2024) "Utilising technology as a transmission strategy in intangible cultural heritage: the case of Cantonese opera performances" **International Journal of Heritage Studies** 30(2): 210–225. DOI: [10.1080/13527258.2023.2284723](https://doi.org/10.1080/13527258.2023.2284723).
- [3] Z. Xu and L. Jiang, (2025) "Federated learning-based fault location and identification in hybrid AC/DC distribution systems considering bidirectional power flow" **J. Eng. Appl. Sci.** 72(1): 133. DOI: [10.1186/s44147-025-00694-w](https://doi.org/10.1186/s44147-025-00694-w).
- [4] Y. Zhong, X. Fu, Z. Liang, Q. Chen, R. Yao, and H. Ning, (2025) "The Application of Artificial Intelligence Technology in the Field of Dance" **Applied System Innovation** 8(5): DOI: [10.3390/asi8050127](https://doi.org/10.3390/asi8050127).
- [5] C. Xu, Y. Sun, and H. Zhou, (2025) "Artificial Aesthetics and Ethical Ambiguity: Exploring Business Ethics in the Context of AI-driven Creativity" **J Bus Ethics** 199(4): 671–692. DOI: [10.1007/s10551-024-05837-2](https://doi.org/10.1007/s10551-024-05837-2).
- [6] N. Partarakis and X. Zabulis, (2024) "A Review of Immersive Technologies, Knowledge Representation, and AI for Human-Centered Digital Experiences" **Electronics** 13(2): 269. DOI: [10.3390/electronics13020269](https://doi.org/10.3390/electronics13020269).
- [7] D. Kostadimas, V. Kasapakis, and K. Kotis, (2025) "A Systematic Review on the Combination of VR, IoT and AI Technologies, and Their Integration in Applications" **Future Internet** 17(4): 163. DOI: [10.3390/fi17040163](https://doi.org/10.3390/fi17040163).
- [8] R. Šajina and M. Ivašić-Kos, (2022) "3D Pose Estimation and Tracking in Handball Actions Using a Monocular Camera" **Journal of Imaging** 8(11): 308. DOI: [10.3390/jimaging8110308](https://doi.org/10.3390/jimaging8110308).
- [9] C. Fang, (2025) "AI-driven digital sculpture design: optimising fusion algorithms with deep learning and virtual reality" **International Journal of Information and Communication Technology** 26(22): 55–71. DOI: [10.1504/IJICT.2025.146908](https://doi.org/10.1504/IJICT.2025.146908).
- [10] S. Rani, D. Jining, D. Shah, S. Xaba, and K. Shoukat, (2025) "Examining the impacts of artificial intelligence technology and computing on digital art: a case study of Edmond de Belamy and its aesthetic values and techniques" **AI & Soc** 40(4): 2417–2435. DOI: [10.1007/s00146-024-01996-y](https://doi.org/10.1007/s00146-024-01996-y).
- [11] D. Horváth, (2025) "Curtain call for AI: Transforming theatre through technology" **Sustainable Futures** 9: 100747. DOI: [10.1016/j.sftr.2025.100747](https://doi.org/10.1016/j.sftr.2025.100747).
- [12] K. El-Raheb, L. Kougioumtzian, V. Kalampratsidou, A. Theodoropoulos, P. Kyriakoulakos, and S. Vosi-nakis, (2025) "Sensing the Inside Out: An Embodied Perspective on Digital Animation Through Motion Capture and Wearables" **Sensors** 25(7): 2314. DOI: [10.3390/s25072314](https://doi.org/10.3390/s25072314).
- [13] T. Wang, (2025) "Domain Adaptive English Aspect Word Extraction Method Based On Bidirectional Long And Short-term Memory Network And Multi-head Attention Mechanism" **Journal of Applied Science and Engineering** 28(12): 2661–2669. DOI: [10.6180/jase.202512_28\(12\).0013](https://doi.org/10.6180/jase.202512_28(12).0013).
- [14] Y. Zhong, X. Fu, Z. Liang, Q. Chen, R. Yao, and H. Ning, (2025) "The Application of Artificial Intelligence Technology in the Field of Dance" **Applied System Innovation** 8(5): 127. DOI: [10.3390/asi8050127](https://doi.org/10.3390/asi8050127).
- [15] S. R. Sitaraman and P. Alagarsundaram, (2024) "Advanced IoMT-Enabled Chronic Kidney Disease Prediction Leveraging Robotic Automation with Autoencoder-LSTM and Fuzzy Cognitive Maps" **International Journal of Modern Electronics and Communication Engineering** 12(3):
- [16] D. Chen, N. Sun, J.-H. Lee, C. Zou, and W.-S. Jeon, (2024) "Digital Technology in Cultural Heritage: Construction and Evaluation Methods of AI-Based Ethnic Music Dataset" **Applied Sciences** 14(23): 10811. DOI: [10.3390/app142310811](https://doi.org/10.3390/app142310811).
- [17] F. Gîrbacia, (2024) "An Analysis of Research Trends for Using Artificial Intelligence in Cultural Heritage" **Electronics** 13(18): 3738. DOI: [10.3390/electronics13183738](https://doi.org/10.3390/electronics13183738).
- [18] K. Zhang and F. Fassi, (2025) "Transforming Architectural Digitisation: Advancements in AI-Driven 3D Reality-Based Modelling" **Heritage** 8(2): 81. DOI: [10.3390/heritage8020081](https://doi.org/10.3390/heritage8020081).
- [19] Y. Lian and J. Xie, (2024) "The Evolution of Digital Cultural Heritage Research: Identifying Key Trends, Hotspots, and Challenges through Bibliometric Analysis" **Sustainability** 16(16): 7125. DOI: [10.3390/su16167125](https://doi.org/10.3390/su16167125).

- [20] Y. Cohen, A. Biton, and S. Shoval, (2025) "Fusion of Computer Vision and AI in Collaborative Robotics: A Review and Future Prospects" **Applied Sciences** 15(14): 7905. DOI: [10.3390/app15147905](https://doi.org/10.3390/app15147905).
- [21] X. Ju, (2025) "The Application of Deep Learning in Dance Movement Design" **Int J Comput Intell Syst** 18(1): 183. DOI: [10.1007/s44196-025-00907-3](https://doi.org/10.1007/s44196-025-00907-3).
- [22] Z. Yuezhou, H. Xiangzhen, M. Xianghe, L. S. Shuai, W. Jiabin, B. Xue, M. Mengdi, L. Zhenjie, C. Ning, W. Hao, W. Lindong, and L. Xihong. *Three-dimensional motion dataset of Dunhuang dance*. Version V1. Accessed: 2025-12-02. 2024. DOI: [10.57760/sciencedb.j00001.01093](https://doi.org/10.57760/sciencedb.j00001.01093).
- [23] C.-B. Lin, Z. Dong, W.-K. Kuan, and Y.-F. Huang, (2021) "A Framework for Fall Detection Based on Open-Pose Skeleton and LSTM/GRU Models" **Applied Sciences** 11(1): 329. DOI: [10.3390/app11010329](https://doi.org/10.3390/app11010329).
- [24] J. Liu, X. Mu, Z. Liu, and H. Li, (2023) "Human skeleton behavior recognition model based on multi-object pose estimation with spatiotemporal semantics" **Machine Vision and Applications** 34(3): 44. DOI: [10.1007/s00138-023-01396-0](https://doi.org/10.1007/s00138-023-01396-0).
- [25] M.-F. R. Lee, Y.-C. Chen, and C.-Y. Tsai, (2022) "Deep Learning-Based Human Body Posture Recognition and Tracking for Unmanned Aerial Vehicles" **Processes** 10(11): 2295. DOI: [10.3390/pr10112295](https://doi.org/10.3390/pr10112295).