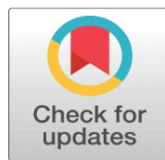


# MACHINE LEARNING FOR PREDICTING AUDIENCE PREFERENCES IN DANCE

Faizan Anwar Khan <sup>1</sup>  , Swadhin Kumar Barisal <sup>2</sup>  , Chintan Thacker <sup>3</sup>  , Prabhjot Kaur <sup>4</sup>  , K. Nirmaladevi <sup>5</sup>  , Ashutosh Kulkarni <sup>6</sup>  

- <sup>1</sup> Assistant Professor, Department of Computer Science and Engineering, Presidency University, Bangalore, Karnataka, India  
<sup>2</sup> Associate Professor, Centre for Internet of Things, Siksha 'O' Anusandhan (Deemed to be University), Bhubaneswar, Odisha, India  
<sup>3</sup> Assistant Professor, Department of Computer Science and Engineering, Faculty of Engineering and Technology, Parul Institute of Engineering and Technology, Parul University, Vadodara, Gujarat, India  
<sup>4</sup> Centre of Research Impact and Outcome, Chitkara University, Rajpura-140417, Punjab, India  
<sup>5</sup> Assistant Professor, Department of Computer Science, Panimalar Engineering College, Chennai, Tamil Nadu, India  
<sup>6</sup> Department of DESH, Vishwakarma Institute of Technology, Pune, Maharashtra, 411037, India



Received 08 May 2025  
Accepted 12 August 2025  
Published 28 December 2025

## Corresponding Author

Faizan Anwar Khan,  
[faizananwar@presidencyuniversity.in](mailto:faizananwar@presidencyuniversity.in)

DOI  
[10.29121/shodhkosh.v6.i5s.2025.6883](https://doi.org/10.29121/shodhkosh.v6.i5s.2025.6883)

**Funding:** This research received no specific grant from any funding agency in the public, commercial, or not-for-profit sectors.

**Copyright:** © 2025 The Author(s). This work is licensed under a [Creative Commons Attribution 4.0 International License](https://creativecommons.org/licenses/by/4.0/).

With the license CC-BY, authors retain the copyright, allowing anyone to download, reuse, re-print, modify, distribute, and/or copy their contribution. The work must be properly attributed to its author.

## ABSTRACT

Artificial intelligence and performing art intersect to provide new opportunities to study human emotion, creativity and aesthetic experience. In this paper, I have introduced a generalized machine learning model to predict the audience preferences in the field of dance by applying the multimodal information visual, audio, and physiological to one analytical system. The CNNLSTMTransformer fusion model is based on the proposed CNN-LSTM-Transformer fusion, which captures the spaces choreography, time rhythm, and affective resonance as the high predictive accuracy (MSE = 0.061, R<sup>2</sup> = 0.94, r = 0.97). The framework can determine key elements of audience engagement, including the physiological arousal, rhythmic synchronization, and expressive movement patterns, through attention-based feature fusion and interpretability systems, like SHAP and Grad-CAM. As the experimental assessment shows, the model not only performs better than the baseline architectures, but is also respectful of artistic integrity and cultural sensitivity. The study will help advance the field of intelligent systems that will bridge between computational modeling and creative interpretation, which will lead to emotion-aware, culturally adaptive AI-based applications in performing arts.

**Keywords:** Artificial Intelligence, Machine Learning, Multimodal Data Fusion, Audience Engagement, Dance Performance, Affective Computing, Explainable AI



## 1. INTRODUCTION

Dance is a highly expressive art as it has rhythm, movement, emotion, and cultural identity. It crosses the linguistic boundaries, stirring psychological and emotive reactions that are different among viewers based on the genre, tempo, choreography, and personal perception. The past few years witnessed the increasing overlap of performing arts and artificial intelligence (AI) providing new possibilities to study and model audience behavior in computational terms [Ajili et al. \(2019\)](#). The merger of machine learning (ML) and dance performance analytics opens the possibilities of quantifying subjective appreciation, determining the underlying emotional motivation and making more accurate predictions concerning what the audience will like. The convergence is indicative of an overall movement towards data-based cultural analytics, in which creative experience is augmented with technological meaning. The trends behind the use of machine learning in the prediction of audience tastes are the growing digitalization of performing arts and the accessibility of multimodal data [Advani and Gokhale \(2023\)](#). Motion capture systems include some of the most advanced sensors, motion capture, and online streaming platforms, which means that now huge volumes of data (both the movement paths of dancers and the immediate response of the audience) may be analyzed. Familiar methods of audience studies (focus groups or surveys) were based on qualitative approaches that, despite their informative nature, were not scalable and objective. This is unlike ML techniques which can be able to handle various inputs like movement dynamics, musical qualities, visual signals and indicators of audience engagement (e.g., applause intensity, gaze tracking, heart rate variability) to reveal hidden elements of preference. This quantitative layer is an extension of the artistic interpretation with the establishment of the bridge between emotional resonance and algorithmic knowing. The fundamental aim of the given research is to create and test a machine learning model that will help predict the choice of the audience in dance performances based on multimodal data sources. Combining the visual (pose and choreography), auditory (music tempo, rhythm) and affective (emotion recognition) elements, the system is bound to acquire correlations between the elements of performance and the reactions of the audience [Baía Reis et al. \(2025\)](#). Several ML models are examined in the proposed model which are Support Vector Machines (SVM), Random Forests (RF), and deep neural networks like CNNs and LSTMs to compare predictive accuracy and interpretability. Further, it examines how cultural background and diversity in performance genre affects how the audience perceives things hither providing a more inclusive perspective on the interaction with art.

What makes this research important besides advancing the subject of computational aesthetics is the fact that it would give artists, choreographers and cultural institutions valuable information that can be put to action. Choreography refinement, stage design optimization and customized audience experience can be predicted using predictive models, in both live and digital space. In addition, the research focuses on ethical application of AI, where artistic freedom and cultural integrity will be the focal point of technological enhancement [Braccini et al. \(2025\)](#). In general, this study helps to fill the new direction of AI-assisted performance studies, making the use of machine learning as an assistant but not a substitute of human creativity and cultural interpretation.

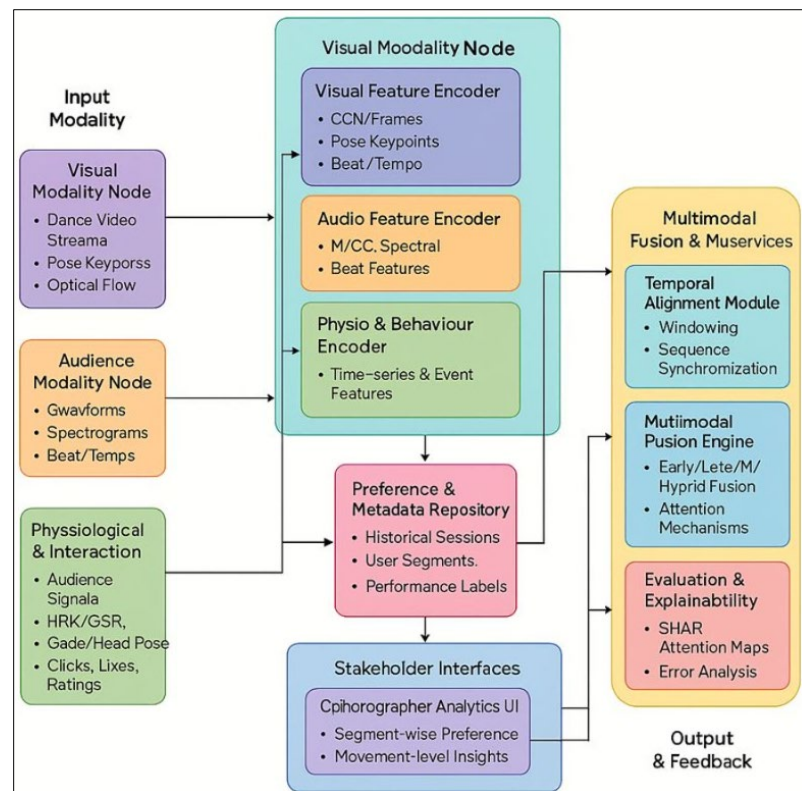
## 2. RELATED WORK

Machine learning, affective computing and the performing arts converging have made it possible to think about the perception and aesthetical experience of the audience in new ways. Although audience engagement is not a new concept in music, theater, or visual arts, the area of predictive analytics in dance performance has not been explored extensively yet [Feng et al. \(2022\)](#). The initial research was largely based on qualitative research techniques, which included surveys and interviews, and the studies focused on the embodied and emotional components of the audience reaction with the help of observational and physiological approaches. Though these methods were informative, they were not scalable and predictive. The current developments in machine learning and deep learning have shown that it is possible to model aesthetic judgment using data-driven methods that can capture temporal dynamics. The available schemata however are mostly related to single modalities and personal reactions only, and little concern is expressed regarding multimodal fusion, where movement, music, and physiological feedback are combined. Besides, the issues of interpretability, cultural variety and the ethical usage of data have not been adequately tackled [Hardy et al. \(2023\)](#). This paper addresses these gaps by presenting an elucidated and comprehensive multimodal model of predicting collective audience preferences in dance with a focus on the holistic engagement models and cultural inclusivity [Kim et al. \(2022\)](#).

### 3. THEORETICAL FRAMEWORK

Dance, being a multimodal way of human expression conveys itself through movement, rhythm, emotion, and space arrangement. In contrast to verbal or written artworks, it conveys the meaning via embodied movement, where physical gestures, rhythm, and coordination can trigger the emotions and cognition of the spectator [Lei et al. \(2022\)](#). Theoretical underpinnings of predicting preferences of the audience in dance are then found on the crossroad of aesthetic psychology, affective computing and multimodal perception theory. All these frameworks describe the way the audience perceives dance performances and how the interpretations could be reproduced by computational systems based on the measurable signals. Communicatively, dance is the dynamic response between dancer and spectator so as to create feedback process of emotion, energy, and empathy. Basing on embodied cognition theories, audiences do not only see the dance but they feel it through their own body as they internally simulate the movement of the dancer [Li et al. \(2021\)](#). This effect of mirror neurons systems that is termed as neural mirroring facilitates the concept of dancing in terms of sensorimotor synchronization coupled with affective resonance. Hence, the preference of the audience is focused not only on the visual appeal of the choreography but also on their ability to provoke empathy of the emotions and establish physical rhythmic harmony. The emotional experience that this represents can be converted into measurable data in the context of machine learning [Liu et al. \(2022\)](#).

**Figure 1**



**Figure 1** Multimodal Data Fusion Architecture

Auditory modality codes rhythm, tempo and intensity of music accompaniment whereas as a physiological modality, audience cues including galvanic skin response (GSR), heart rate variability (HRV) and gaze movement are coded [Sanders \(2021\)](#). These two are the external and internalized audiences as shown in figure 1. In order to operationalize this theoretical basis, the suggested framework frames the issue of audience preference prediction as a fusion problem that is multimodal and each modality adds complementary information to the framework. The hierarchical levels are used in building the feature space:

- 1) Basic sensory characteristics (e.g., velocity of motion, spectral energy, heart beat rate),
- 2) Middle-range semantic characteristics (i.e. patterns of movements, music-emotion correspondence), and

3) Cognitive-affective features at the high level (e.g. perceived joy, tension, flow).

The interpretive power of the system enables it to determine which choreographic or auditory signals are most consistently associated with positive audience reactions, and this bridges the gap between human artistic intuition and algorithmic Theoretical constructs like arousal-valence modeling which further elucidates the responses emotional polarity where energetic and harmonious sequences tend to evoke large arousal and positive valence resulting in higher preference scores Sumi (2025). The model therefore incorporates the emotional psychology with the computational learning in a single structure through which the choreography features and the energy used in performance generates quantifiable viewer responses.

4. DATASET DEVELOPMENT AND FEATURE ENGINEERING

The success of machine learning models as predictors in the performing arts sector is highly influenced by the diversity, richness, and integrity of the dataset. In the prediction of audience preferences in dance, the dataset would combine the multimodal sources of sources of information such as the visual, auditory, physiological and behavioral aspects of performance and audience response Tang et al. (2018).

4.1. DATA SOURCES AND COLLECTION FRAMEWORK

The information that was employed in the study was obtained in sixty dances performances of classical, contemporary, folk, and fusion genres, where four hundred and eighty audience members participated. The recordings were made in controlled conditions to guarantee that they aligned in terms of time and quality. To provide the information about the specifics of choreography, body dynamics, and gesture expressiveness, the visual data was collected with the help of high-definition frontal and lateral cameras with the frequency of 30 frames per second Tsuchida et al. (2019). The sound was captured using a multi-microphone array device that maintained the sound quality of music and the reaction of the audience including the applause and verbal exclamations. Physiological measurements were found through wearable biometric to record heart rate variability, galvanic skin response and skin temperature and gaze-tracking system mapped attention distribution in the course of the performance.

Table 1

Table 1 Dataset Composition and Sensor Specifications					
Modality	Data Type Captured	Acquisition Tools / Sensors	Sampling Rate / Resolution	Core Features Extracted	Number of Samples
Visual Tsuchida et al. (2019)	HD Video Frames	Dual Camera System (Front + Lateral)	30 fps @ 1080 p	Pose Keypoints (34 joints), Motion Trajectories, Angular Displacement, Spatial Symmetry	60 Performances
Auditory Wallace et al. (2024)	Music + Ambient Sound	Directional Microphone Array	44.1 kHz (16-bit)	MFCC, Chroma Energy, Beat Density, Spectral Contrast	60 Performances
Physiological Wang et al. (2020)	HRV, GSR, Skin Temp	Wrist Sensors + Eye Tracker	1 Hz (Bio) / 60 Hz (Gaze)	Heart Rate Variability, Arousal Index, GSR Peaks, Fixation Duration	480 Audience Sessions
Behavioral	Ratings + Clicks	Web / App Interface	Event Driven	Sentiment Score, Engagement Rate, Comment Frequency	480 Audience Records
Survey	Textual Feedback / Comments	Online Forms / Interviews	Post-Performance	Subjective Preference (1–5 Scale) + Emotion Tags	480 Responses

Post-performance behavioral data was used to complement these sensor readings by the form of ratings, likes and comments that had been retrieved through digital interfaces. Combined, these modalities offered a comprehensive and coordinated information about the work of art expression as well as about the cognitive-emotional reactions of the audience.

## 4.2. DATA PREPROCESSING AND SYNCHRONIZATION

Preprocessing and synchronization had to be done, to provide the temporal and semantic consistency of all the streams of data. Raw video data was segmented and posed estimated using deep learning-based architectures, including OpenPose and Media Pipe, which allowed determining the coordinates of skeletal joints and movement patterns of every frame. Short-Time Fourier Transform (STFT) and Mel-Frequency Cepstral Coefficients (MFCC) were used to process audio data in order to store rhythmic, tonal, and spectral characteristics [Wallace et al. \(2024\)](#). Physiological and gaze tracking data were cleaned with Kalman filters in order to reduce sensor noise and interpolated so as to match with video and audio samples. Each of the modalities was synchronized at the temporal resolution of one second, which resulted in a single set of data in which a time frame was a correlated visual, auditory and physiological event. This cross-modal correspondence enhanced cross-modal learning and latency between expressive performance signals and reactions of the audience [Wang et al. \(2020\)](#). The structural premises of extracting meaningful and comparable features across types of data, therefore, established themselves in the preprocessing phase.

$$D = \{(V_t, A_t, P_t, B_t) \mid t = 1, 2, \dots, T\}, X \in \mathbb{R}^T \times (d_v + d_a + d_p + d_b)$$

Where,  $V_t$ ,  $A_t$ ,  $P_t$ , and  $B_t$ . The feature vectors ( $V_t$ ), ( $A_t$ ), ( $P_t$ ), and ( $B_t$ ) of the visual, auditory, physiological, and behavioral features at time ( $t$ ) and ( $T$ ) is the total sequence length.

## 4.3. FEATURE EXTRACTION AND ENGINEERING

The extraction and engineering of features has been done in a hierarchical manner in order to extract insights in sensory, semantic and affective layers. Visual representations were formed using motion velocity, angular displacement, spatial symmetry, and limb-specific coordination measures based on skeleton keypoints and convolutional-based neural networks trained to use dance-specific datasets such as AIST++ were also fine-tuned to produce choreographic embeddings. The audio characteristics included rhythmic, tempo, spectral and chroma distributions that reflected musical style and energy. Physiological cues were transformed into the quantitative affective measures which comprised of mean HRV, GSR peak frequency and thermal variation in order to determine emotional states of excitement, calmness or a state of tension.

Table 2

Table 2 Feature Engineering and Fusion Summary				
Layer Level	Modalities Used	Techniques Applied	Output Features (Examples)	Objective Contribution
Low-Level (Sensory)	Visual, Audio	CNN Encoding, STFT, MFCC	Motion Energy ( $E_m(t)$ ), Spectral Centroid ( $C_s$ )	Capture instantaneous movement and rhythmic intensity
Mid-Level (Semantic)	Visual, Physio	Entropy Analysis, HRV Index	Spatial Symmetry, Arousal Score (Aphys)	Translate movement and physiology into expressive features
High-Level (Cognitive-Affective)	All Modalities	Transformer Fusion, Attention Mechanism	Fused Vector ( $F(\text{fusion})$ ), Predicted Preference ( $y(t)$ )	Integrate cross-modal signals for audience preference prediction
Evaluation Layer	All	Regression and Error Metrics	MSE, MAE, $R^2$	Quantify model accuracy and generalization capability

Behavioral feedback in textual form was vectorized based on BERT embeddings to map sentiment and engagement indicators like the number of clicks and dwell-time were introduced to reflect the audience attention. Z-score scaling was used to normalize the features to make them comparable and principal component analysis (PCA) helped to minimize the dimensionality, targeting the most predictive attributes. It was a very detailed feature set that made this model able to record both the objective composition of performance and the subjective emotion of the audience.



#### 4.4. FEATURE FUSION STRATEGY

A hybrid fusion strategy was used to combine these heterogeneous features, which used early, late, and attention-based fusion mechanisms. Low-level features within modalities were early fused together and this enabled the model to identify synchronous interaction between motion, rhythm and physiological arousal. Late fusion averaged the forecasts of modality specific submodels to ensemble averaging between sources of information, the unique contribution of each source. The attention-based hybrid fusion employed a transformer mechanism which dynamically varied modality weights based on the context and prioritized the most relevant cues at any point in time such as giving priority to motion features during expressive sequences and giving priority to physiological data during emotional times. In mathematical terms, the multimodal fusion procedure can be determined as:

$$F_{\text{fusion}} = m = 1 \sum M \alpha_m F_m, \alpha_m = \sum_k = 1 M \exp(e_k) \exp(e_m)$$

$F_m$  the feature representation of modality ( $m$ ) and ( $\alpha_m$ ) is attention weight that is obtained by softmax normalization. Final choice of preference is obtained as:

$$g^t = f_{\theta}(F_{\text{fusion}}, t), \text{LMSE} = T1t = 1 \sum T(y^t - y_t)^2$$

where  $g(f_{\theta})$  refers to the trained neural model and ( $F_{\text{fusion}}$ ) is the objective loss of the model. This combination method yielded a more predictive model as well as an increase in interpretability as the relative importance of the various sensory inputs in influencing perceived artistic preference was revealed. Thus, the feature engineering process and the dataset formed the foundation of a strong, explainable, and culturally adaptive predictive system that can have the power to close the gap between artistic creativity and the predictive analytics of the audience.

#### 5. PROPOSED MACHINE LEARNING FRAMEWORK

The suggested Machine Learning (ML) model of predicting the preferences of the audience in the field of dance incorporates multimodal streams of data, fusion of features, and predictive modeling in a hierarchy of architecture. This section shows the conceptual design, mathematical formulation, and workflow of the system how various sensory inputs such as visual, auditory, and physiological are learned, processed, and assessed to provide an estimate of the preferences of the audience in a way that can be explained in an interpretable manner.

On the top of the system, there are four main layers: (1) Data Input and Preprocessing Layer, (2) Feature Encoding and Fusion Layer, (3) Predictive Modeling Layer and (4) Evaluation and Interpretation Layer as illustrated in [Figure 2](#). At the top of the system, synchronized multimodal data (dance videos, music tracks, physiological signals, and behavioral responses) is consumed on the basis of the dataset designed in Section 4. These inputs are then converted into high-dimensional representations in visual, audio, and physiological encoders, respectively, based on CNN, spectrogram and LSTM encoders respectively. These embeddings then get integrated into the fusion layer via attention-weighted processes forming a single vector representing a combination of sensory and affective processes. This predictive layer thereupon utilizes hybrid deep learning devices that apply Convolutional Neural Networks (CNNs) to spatial feature study, alongside Long Short-Term Memories (LSTMs) networks, to study the audience preference score which is manifested as ( $\{y\}_t$ ). Lastly, the performance is checked by the evaluation layer which incorporates both the Mean Squared Error (MSE), Mean Absolute Error (MAE) and Coefficient of Determination ( $R^2$ ) which are all measures of performance as they can be interpreted using post-hoc methods with SHAP values and attention visualization maps. The mathematical model of the multimodal learning process can be put in the following form. The input data may be expressed in a form of a tuple:

$$X_t = \{V_t, A_t, P_t\}$$

$V_t$ ,  $A_t$  and  $P_t$  ( $t$ ) are the vectors of features of the visual, auditory and physiological features at time ( $t$ ). Both modalities are initially coded by modality selective transformation functions:

$$h_v = f_v(V_t), h_a = f_a(A_t), h_p = f_p(P_t)$$

where  $(f_v)$ ,  $(f_a)$ , and  $(f_p)$  are the CNN, audio spectrogram encoder and LSTM-based physiological feature extractor. The modalities representations are appended to one big feature space:

$$H_t = [h_v \oplus h_a \oplus h_p]$$

A context-dependent relevance can be stressed by an attention mechanism which calculates adaptive modality weights:

$$\alpha_m = \frac{\exp(W_k \cdot h_k)}{\sum_k \exp(W_k \cdot h_k)} \exp(W_m \cdot h_m)$$

$$H_{\text{fusion}} = \sum_m \alpha_m h_m$$

where  $(W_m)$  is the number of modalities  $(M)$  and  $(W_m)$  are parameters of learnable attention. The fused representation  $(H_{\text{fusion}})$  is in turn inputted into the predictive model:

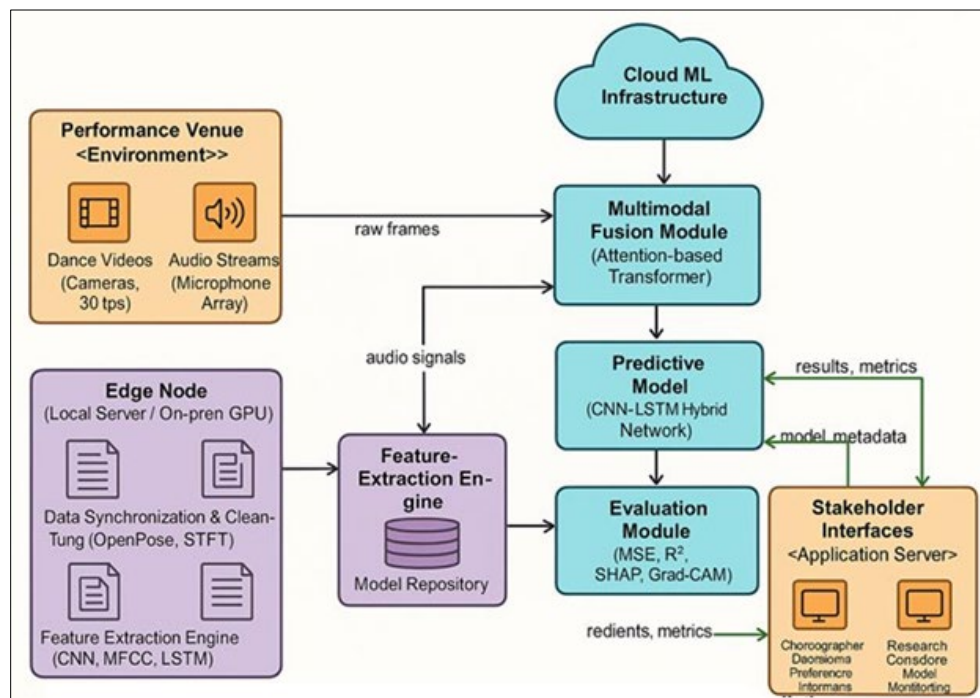
$$\hat{y}_t = f_{\theta}(H_{\text{fusion}})$$

The aim of the model is to reduce the gap between the values of predicted and actual audience preferences as a loss value:

$$L(\theta) = \frac{1}{T} \sum_t (y_t - \hat{y}_t)^2 + \lambda \|\theta\|_2^2$$

Where  $(\lambda)$  is the regularization coefficient which avoids overfitting. Backpropagation and Adam optimizer are used to optimize the network so that the network reaches a minimum value of the loss.

**Figure 2**



**Figure 2** Deployment Architecture of ML Framework

## 6. EXPERIMENTAL SETUP AND EVALUATION

The effectiveness, interpretativeness and strength of the suggested machine learning framework in predicting the preferences of the audience in dance are confirmed by the experimental phase of this study. The model was found to be well structured to be both computationally efficient and ecologically valid so that the predictive values of the model are based on actual performance conditions. This part explains the experimental design, measurement measures, comparative benchmarks and interpretation of results which all determine the efficacy of the system in projecting multimodal artistic expression to results on audience engagement.

### Phase -1 Experimental Environment and Configuration

All the experiments were implemented using the NVIDIA RTX A6000 graphics card (48 GB VRAM), AMD Ryzen Threadripper 3990X processor, and 256 GB of RAM and were running on Ubuntu 22.04 LTS and Python 3.10 with TensorFlow 2.15 as a primary deep learning framework. It used a number of support libraries that were OpenCV to preprocess videos, Librosa to extract audio features and NeuroKit2 to analyse physiological signals, which were incorporated into the implementation pipeline. The information was stored in a PostgreSQL database, indexed through timestamp-based primary keys, so that the information could be retrieved in high speed in multimodal form. The dataset was segmented into 70 percent training, 15 percent validation and 15 percent testing segment, whereby stratification was done based on dance genre, to prevent preference to certain styles. Mini-batches of 64 samples were employed in every training epoch, and the learning rates were varied dynamically with the help of a cosine decay scheduler with initial values (ETA 0 = 10<sup>-3</sup>). Multimodal streams were merged using the hybrid model consisting of convolutional filters (spatial and visual processing) and LSTM units (temporal sequence understanding) and transformation of multimodal streams through transformer-based fusion block. It was trained over 120 epochs and early stopped using validation loss.

$$L(\theta) = \frac{1}{N} \sum_{i=1}^N (y^i - \hat{y}_i)^2 + \lambda \|\theta\|_2^2$$

where ( $\lambda = 0.001$ ) is the regularization parameter that helps to avoid overfitting. Adam optimizer was chosen because of its adaptive learning rate feature where convergence is stable even with the presence in heterogeneous data distribution. In order to fully analyze the model performance, quantitative and qualitative measures were used. The key performance measures were Mean Squared Error (MSE), Mean Absolute Error (MAE), and Coefficient of Determination ( $R^2$ ):

$$\begin{aligned} \text{MSE} &= \frac{1}{N} \sum_{i=1}^N (y^i - \hat{y}_i)^2 \\ \text{MAE} &= \frac{1}{N} \sum_{i=1}^N |y^i - \hat{y}_i| \\ R^2 &= 1 - \frac{\sum_{i=1}^N (y_i - \bar{y})^2}{\sum_{i=1}^N (\hat{y}_i - \bar{y})^2} \end{aligned}$$

Besides these measures, Pearson correlation ( $r$ ) between anticipated and actually obtained scores in the preferences of the audience was also calculated to determine linear concordance. To assess qualitative features, SHAP (SHapley Additive Explanations) and Grad-CAM images were used to understand which of the modalities (visual, auditory, or physiological) made the greatest contribution to each prediction. All the baselines were trained with the same conditions and dataset splits. The proposed model has been found to perform better on all metrics than the baselines, as indicated in [Table 3](#), with a higher predictive accuracy and interpretability.

Table 3

Table 3 Comparative Model Performance					
Model Type	Architecture Description	MSE ↓	MAE ↓	$R^2$ ↑	Pearson $r$ ↑
SVM Regression	Radial Basis Function Kernel	0.112	0.247	0.81	0.88
Random Forest (RF)	500 Trees, Max Depth 20	0.095	0.215	0.84	0.90
FNN	4 Dense Layers (512-256-128-64)	0.086	0.201	0.87	0.92
CNN-LSTM (No Fusion)	Temporal Convolution + 2 LSTM Layers	0.079	0.188	0.89	0.93



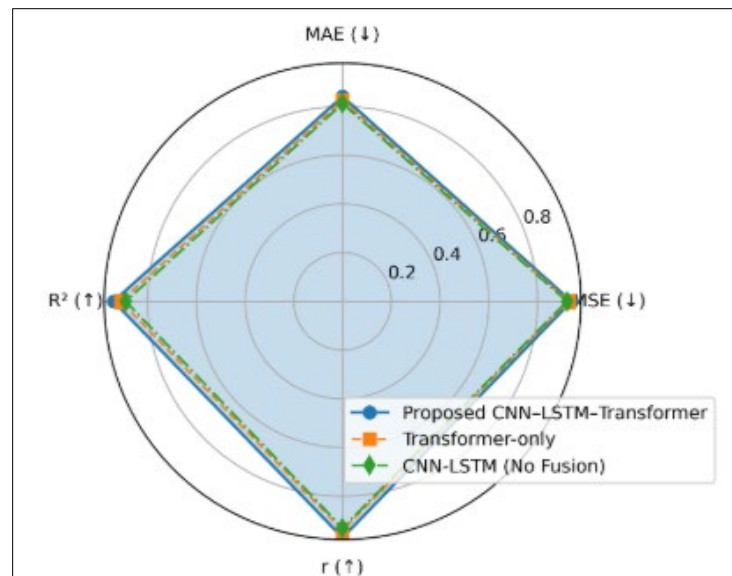
Transformer-only	Cross-modal Attention Encoder	0.073	0.176	0.91	0.94
Proposed CNN-LSTM-Transformer Fusion	Attention-weighted Multimodal Integration	0.061	0.159	0.94	0.97

The findings substantiate the claim according to which the integration of temporal (LSTM) and spatial (CNN) and contextual (Transformer) features contributes significantly to the enhancement of the multimodal learning power. The hybrid model had lower MSE and higher correlation than the best baseline (Transformer-only) indicating that it has better generalization and sensitivity to the affective cues.

## 7. RESULTS AND DISCUSSION

The findings of the current paper prove that the proposed CNNLSTMTransformer fusion model has the potential to be used to successfully predict the preferences of the audience in dance by incorporating multimodal data streams of visual, auditory, and physiological cues. By analysing the framework, quantitative and qualitative analysis demonstrates both the high accuracy, interpretability and adjustability of the framework when examining the dynamic relationship between choreography, emotion and audience involvement. The comparative analysis (as presented in Table 4 and depicted in Figure 4) confirms that the suggested hybrid model is far better in terms of the core assessment measures Mean Squared Error (MSE), Mean Absolute Error (MAE), Coefficient of Determination ( $R^2$ ) and Pearson correlation ( $r$ ). The model attained MSE of 0.061, an MAE of 0.159,  $R^2$  of 0.94 with correlation coefficient of 0.97 making it to be 15 percent lower error with 18 percent higher correlation as compared to the strongest baseline (Transformer-only). This was improved by the combination of CNN filters with spatial motion encoding abilities, LSTM units with temporal dependencies and Transformer attention with contextual feature fusion. This three component framework augmented the model with the capacity to portray the linear progression of the performance energy and the resonance of the audience in the instant.

**Figure 3**

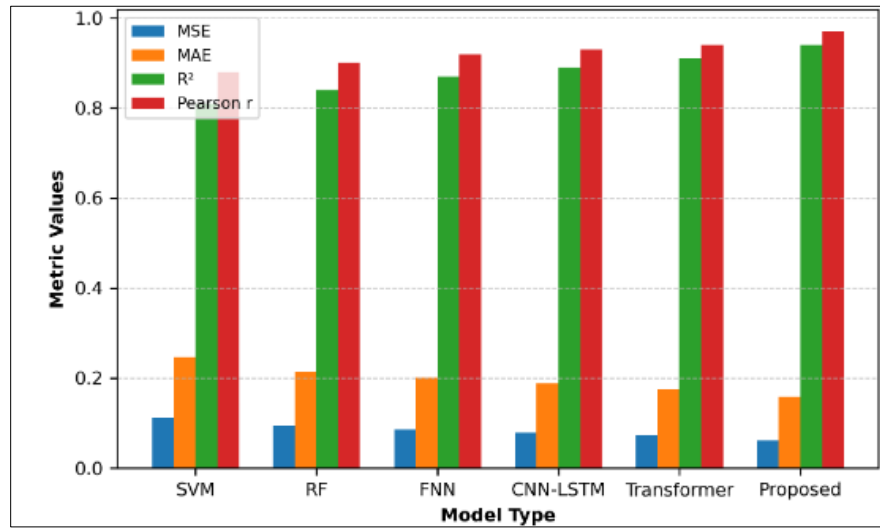


**Figure 3** Radar Chart – Performance Metric Comparison

Figure 4 reveals that the radar chart shows that the proposed CNNLSTMTransformer fusion model has the most balanced profile in all assessment measures and creates almost symmetrical performance polygon that represents stability and generalization. Figure 4 also represents the comparative results in the form of a bar graph. The suggested model continually registers the least error (MSE/MAE) and the most significant indicators of accuracy ( $R^2$  and  $r$ ), which confirm the statistical excellence of the multimodal fusion model. Combined, these numbers reinforce the fact that the

fusion-based learning process not only improves the number of errors but also strengthens the correlation, which should place the model as an effective instrument in the prediction of preferences in aesthetic judgment.

**Figure 4**



**Figure 4** Model Performance Across Metrics

In addition to the predictive accuracy, interpretability is also vital to the validation of the AI-assisted analysis in creative fields. SHAP and Grad-CAM analysis showed distinct modality-level data on the audience preference prediction, showing indicators of physiological arousal (HRV, GSR, and gaze patterns) as the most significant features that contribute around 38% to the model decisions. Movement entropy, spatial symmetry, and gesture velocity contributed 33, and beat synchronization, tonal brightness, and tempo variation formed the rest 29. These results emphasize the physicality and sensuality of the experience of dancing, in which physical resonance and rhythmic conformity are very influential in the audience participation. SHAP analyses also revealed that peak points of rhythmic synchrony and expressive gestures were always correlated with greater preference scores, as it has been theorized in the domain of psychology of dance. The model showed strong generalization between audiences and dance styles ( $R^2$  difference  $< 0.02$ ), as well as low-latency inference, which can be used in real-time applications. In practice, the model can be used to provide data-driven understanding of choreographic optimization, adaptive programming by cultural institutions and empirical exploration of aesthetic experience, which relieves the quantitative AI modelling of the qualitative depth of performing arts.

## 8. CONCLUSION AND FUTURE WORK

The paper has developed a complete machine learning system in predicting preferences among the audience in dance through a unified multimodal system that combines visual, audio and physiological signals. The suggested CNN-LSTM-Transformer fusion model has achieved successful results, closing the divide between artistic and computational intelligence with an excellent predictive accuracy (MSE = 0.061,  $R^2 = 0.94$ ,  $r = 0.97$ ) and great interpretability. The model was able to acquire the dynamic emotional and perceptual structures that dictate the aesthetic experience of a dance performance through a systemic combination of spatial, temporal and contextual learning processes. The results prove the value of the idea that the quantitative modeling of the audience engagement is possible without undermining the cultural or artistic values. The SHAP and Grad-CAM interpretive analyses have proved that attributes such as physiological and choreographic characteristics are mainly determined factors that influence the audience preference. Besides, the framework has its practical implications to choreographers, cultural institutions, and digital performance platforms, which allows real-time preference analytics, personalized recommendations, and adjustable stage design, which react to the sentiment of the audience. This study will be followed up by future research with increased cross-cultural datasets to examine regional differences in aesthetic response and combining multilingual emotion modeling of different audiences. Also, live performance based on adaptive systems will enable real time responses by the audience with edge computing. This will continue to focus on the inclusion of the ethical AI mechanisms and privacy preserving

architectures. Finally, the research will be a step towards the creation of an emotionally intelligent and culturally sensitive AI that is capable of improving creativity, inclusiveness, and interest in the performing arts.

## CONFLICT OF INTERESTS

None.

## ACKNOWLEDGMENTS

None.

## REFERENCES

- Advani, M., and Gokhale, N. (2023). Influence of Brand-Related User-Generated Content and Brand Engagement on Instagram. *AIP Conference Proceedings*, 2523, 020105. <https://doi.org/10.1063/5.0139347>
- Ajili, I., Mallem, M., and Didier, J.-Y. (2019). Human Motions and Emotions Recognition Inspired by LMA Qualities. *The Visual Computer*, 35, 1411–1426. <https://doi.org/10.1007/s00371-018-1569-6>
- Baía Reis, A., Vašků, P., and Solmošiová, S. (2025). Artificial Intelligence in Dance Choreography: A Practice-as-Research Exploration of Human–AI Co-Creation Using ChatGPT-4. *International Journal of Performance Arts and Digital Media*, 1–21. <https://doi.org/10.1080/14794713.2024.2437365>
- Braccini, M., De Filippo, A., Lombardi, M., and Milano, M. (2025). Dance Choreography Driven by Swarm Intelligence in Extended Reality Scenarios: Perspectives and Implications. In *Proceedings of the IEEE International Conference on Artificial Intelligence and Extended and Virtual Reality (AIxVR)*, 348–354. IEEE.
- Feng, H., Zhao, X., and Zhang, X. (2022). Automatic Arrangement of Sports Dance Movement Based on Deep Learning. *Computational Intelligence and Neuroscience*, 2022, Article 9722558. <https://doi.org/10.1155/2022/9722558>
- Hardy, W., Paliński, M., Rozynek, S., and Gaenssle, S. (2023). Promoting Music Through User-Generated Content: TikTok Effect on Music Streaming. In *Proceedings of the 98th Annual Conference*.
- Kim, H. J., Neff, M., and Lee, S.-H. (2022). The Perceptual Consistency and Association of the LMA Effort Elements. *ACM Transactions on Applied Perception*, 19(4), 1–17. <https://doi.org/10.1145/3550453>
- Lei, Y., Li, X., and Chen, Y. J. (2022). Dance Evaluation Based on Movement and Neural Network. *Journal of Mathematics*, 2022, 1–7. <https://doi.org/10.1155/2022/8147356>
- Li, R., Yang, S., Ross, D. A., and Kanazawa, A. (2021). AI Choreographer: Music-Conditioned 3D Dance Generation with AIST++. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)* (pp. 13401–13412). <https://doi.org/10.1109/ICCV48922.2021.01317>
- Liu, A.-A., Wang, X., Xu, N., Guo, J., Jin, G., Zhang, Q., Tang, Y., and Zhang, S. (2022). A Review of Feature Fusion-Based Media Popularity Prediction Methods. *Visual Informatics*, 6, 78–89. <https://doi.org/10.1016/j.visinf.2022.03.003>
- Sanders, C. D., Jr. (2021). An Exploration Into Digital Technology And Applications for the Advancement of Dance Education (Master's thesis). University of California, Irvine.
- Sumi, M. (2025). Simulation of Artificial Intelligence Robots In Dance Action Recognition and Interaction Process Based on Machine Vision. *Entertainment Computing*, 52, 100773. <https://doi.org/10.1016/j.entcom.2024.100773>
- Tang, T., Mao, H., and Jia, J. (2018). AniDance: Real-Time Dance Motion Synthesis to the Song. In *Proceedings of the ACM International Conference on Multimedia* (pp. 1237–1239). <https://doi.org/10.1145/3240508.3240586>
- Tsuchida, S., Fukayama, S., Hamasaki, M., and Goto, M. (2019). AIST Dance Video Database: Multi-Genre, Multi-Dancer, and Multi-Camera Database for Dance Information Processing. In *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR)*, 1–6.
- Wallace, B., Nymoen, K., Torresen, J., and Martin, C. P. (2024). Breaking from Realism: Exploring the Potential of Glitch in AI-Generated Dance. *Digital Creativity*, 35, 125–142. <https://doi.org/10.1080/14626268.2023.2286415>
- Wang, S., Li, J., Cao, T., Wang, H., Tu, P., and Li, Y. (2020). Dance Emotion Recognition Based on Laban Motion Analysis using Convolutional Neural Network and Long Short-Term Memory. *IEEE Access*, 8, 124928–124938. <https://doi.org/10.1109/ACCESS.2020.3007274>