

Joint angle prediction and joint-type classification in human gait analysis using explainable deep reinforcement learning

Deepak N. R.^{1,4}, Soumya Naik P. T.^{2,4}, Ambika P. R.^{2,4}, Shaik Sayeed Ahamed^{3,4}

¹Department of Information Science and Engineering, Atria Institute of Technology, Bangalore, India

²Department of Computer Science and Engineering, City Engineering College, Bangalore, India

³Department of Computer Science and Engineering (DS), Atria Institute of Technology, Bangalore, India

⁴Visvesvaraya Technological University, Belagavi, India

Article Info

Article history:

Received Mar 22, 2025

Revised Oct 15, 2025

Accepted Jan 11, 2026

Keywords:

Deep learning

Explainable artificial intelligence

Human gait analysis

Maximization

Q-learning and mutual

information

Rehabilitation

Reinforcement learning

ABSTRACT

Human gait analysis is a key component of rehabilitation, prosthetics, and sports science, especially for clinical evaluation and the development of adaptive assistive technologies. Accurate joint-angle estimation and dependable joint-type classification remain difficult because of the complex temporal behavior of gait signals and the limited interpretability of many deep learning (DL) approaches. While recent techniques have enhanced predictive accuracy, their clinical applicability is often limited by insufficient transparency and adaptability in learning mechanisms. To overcome these limitations, this work proposes an integrated framework that unifies DL, reinforcement learning (RL), and explainable artificial intelligence (XAI). Stochastic depth neural networks (SDNN) are applied for joint-angle regression, whereas deep feature factorization networks (DFFN) are used for multi-class joint-type classification. Optimization is achieved using Q-learning (QL) and mutual information maximization (MIM), ensuring stable convergence and improved learning efficiency. To improve interpretability, the framework incorporates counterfactual and contrastive explanations, feature ablation studies, and prediction probability analysis. Experimental findings show that the SDNN_MIM model attains an R^2 score of 0.9881, with RL rewards increasing from 0.997 to 0.999 during regression training. For joint-type classification, the DFFN_MIM model achieves an accuracy of 0.95, with reward values improving from 0.90 to 0.98. These results demonstrate the effectiveness of the proposed framework in delivering accurate and interpretable gait predictions, supporting its relevance to biomechanics, healthcare, personalized rehabilitation, and intelligent assistive systems.

This is an open access article under the [CC BY-SA](https://creativecommons.org/licenses/by-sa/4.0/) license.



Corresponding Author:

Shaik Sayeed Ahamed

Department of Computer Science and Engineering (DS), Atria Institute of Technology

Bangalore, Karnataka, 560064 India

Email: shaik.sayeedahamed1999@gmail.com

1. INTRODUCTION

Human gait analysis constitutes a core research domain in biomechanics, rehabilitation, prosthetics, and sports science, with significant relevance to clinical diagnosis, rehabilitation evaluation, and the development of intelligent assistive technologies. Accurate gait assessment supports early detection of movement impairments, facilitates improved prosthetic and orthotic design, and aids in injury prevention. Conventional gait analysis approaches primarily rely on motion-capture systems, force plates, and handcrafted biomechanical

models. While effective in controlled laboratory settings, these methods often face challenges related to high-dimensional data, inter-subject variability, limited generalization across diverse movement patterns, and extensive manual feature engineering requirements. In recent years, deep learning (DL) methods have been widely adopted to address these limitations by automatically learning hierarchical representations from gait data. Despite their effectiveness, DL-based gait models still exhibit notable limitations, including overfitting, restricted interpretability, and inefficient optimization, which constrain their clinical reliability. To mitigate these issues, explainable artificial intelligence (XAI) and deep reinforcement learning (DRL) have gained increasing attention in gait analysis research. In the context of clinical gait analysis (CGA), Slijepcevic *et al.* [1] categorized XAI techniques into data exploration, prediction explanation, and model explanation using approaches such as t-SNE and layer-wise relevance propagation (LRP). Although these methods enhanced interpretability, reinforcement learning (RL)-based optimization was not explored. Likewise, Madanu *et al.* [2] employed XAI for pain assessment, reducing subjectivity but without capturing the sequential and biomechanical complexity of gait. SHAP-based explanation techniques reported in [3], [4] improved clinical confidence; however, these studies were limited to supervised learning paradigms and lacked adaptive optimization strategies. RL has shown strong potential in sequential decision-making and continuous control tasks. The soft actor-critic (SAC) framework presented in [5], [6] enabled stable learning in continuous action spaces, and autonomous locomotion without predefined motion models was investigated in [7]. These studies mainly addressed robotic locomotion, where biomechanical constraints, safety requirements, and interpretability differ from human gait analysis. Guided SAC methods, such as [8], enhanced performance in partially observable environments; however, limited policy transparency restricts their clinical applicability. Model-based RL extensions incorporating uncertainty modeling and model predictive control (MPC) in [9] improved sample efficiency, yet their relevance to human gait remains constrained by safety and explainability concerns. Recent surveys and reviews [10]–[12] highlighted the promise of DRL for gait analysis and rehabilitation while identifying ongoing challenges, including small clinical datasets, dependence on simulated environments, and limited interpretability of learned policies. Explainable RL taxonomies in [12] and roadmap studies in [13] further emphasized the difficulty of explaining sequential decision-making processes in safety-critical applications. IMU-based gait investigations such as [14] demonstrated effective prediction of dynamic balance but did not incorporate reinforcement-driven optimization or biomechanical interpretability. Similarly, GRF-based gait classification in [15] utilized SHAP-based explanations without adaptive learning mechanisms. Beyond gait-focused research, XAI applications in healthcare and sports analytics [15], [16] reported challenges related to dataset quality, predictive performance, and generalization. Ethical transparency and accountability in machine learning were emphasized in [17], while sensitivity to dataset bias was discussed in [18] and [19]. Recent XAI-enabled gait decision-support studies [20], [21] applied LIME and SHAP to support clinical reasoning but encountered scalability and real-time interpretability limitations. Finally, [13] achieved strong foot-condition classification using handcrafted features and LIME explanations, yet lacked automated feature learning and reinforcement-based optimization. Overall, although prior studies demonstrate substantial progress in XAI and DRL for human movement analysis, a unified framework integrating deep neural networks, RL-driven optimization, and explainable mechanisms for accurate, adaptive, and clinically interpretable human gait prediction remains insufficiently investigated.

Despite the substantial progress achieved through deep learning in human gait analysis, several open challenges still restrict its clinical applicability. Most existing works depend on post-hoc interpretability methods applied to supervised learning models, which provide limited insight into model behavior and offer minimal explanation of sequential decision-making processes. As a result, the integration of XAI within RL-based gait analysis frameworks remains largely underexplored. Furthermore, current gait modeling strategies frequently face optimization challenges, including unstable training behavior, limited adaptability to time-varying gait patterns, and reduced generalization across subjects and movement conditions. Although RL approaches, such as SAC, have demonstrated strong performance in robotic locomotion, their effectiveness for modeling human gait dynamics—particularly for combined regression and multi-class classification tasks—has not been thoroughly examined. Addressing these gaps, this study proposes a unified deep learning framework augmented with RL and explainability components to enhance predictive accuracy, learning stability, and clinical interpretability in gait analysis. For joint-angle estimation, stochastic depth neural networks (SDNN) are adopted to improve generalization by dynamically bypassing network layers during training. To ensure stable and efficient optimization, QL and MIM are integrated into the learning process. For joint-type classification, deep feature factorization networks (DFFN) are employed to derive discriminative spatio-temporal gait representations, supporting robust multi-class decision-making. In addition, advanced XAI techniques—including counterfactual

and contrastive explanations, feature ablation analysis, and prediction confidence assessment—are incorporated to deliver clinically meaningful insights and enhance trust in model predictions. Overall, this work contributes a RL-driven and explainable gait analysis framework that unifies accurate prediction, adaptive learning, and transparent decision-making. The proposed methodology establishes a basis for reliable gait modeling applicable to intelligent assistive systems and future clinical deployment. The remainder of this paper is structured as follows: section 2 describes the dataset, preprocessing steps, problem formulation, model architectures, and the integration of RL and XAI strategies, section 3 presents the experimental results and interpretability analysis, and section 4 concludes the study with clinical implications and future research directions.

2. METHOD

2.1. Research design

The increasing demand for data-driven and clinically dependable human movement analysis highlights the challenge of accurately modeling complex gait dynamics. This study concentrates on developing a unified framework capable of performing joint-angle regression and multi-class joint-type classification while maintaining robustness, learning stability, and clinical interpretability. To accomplish this, the proposed approach integrates deep neural architectures with RL and mutual information-based optimization, forming a cohesive pipeline illustrated in Figures 1–4. For joint-angle estimation, SDNN are employed to capture temporal joint trajectories. As shown in Figure 1, SDNN utilizes a probabilistic layer-skipping strategy in which each network block (P0–P3) is assigned a survival probability. Shallower layers remain active during training, while deeper layers are selectively bypassed. When a layer is skipped, its output is substituted with a shortcut connection from the preceding layer, enabling uninterrupted forward propagation. This architecture mitigates overfitting, enhances generalization, and promotes stable learning from noisy and variable gait signals by learning hierarchical temporal representations. For multi-class joint-type classification, deep feature factorization (DFF), depicted in Figure 2, is applied to enable structured feature decomposition and dimensionality reduction. Raw gait signals are initially processed through feature extraction and reshaped into matrix form, which is subsequently factorized into basis and activation matrices. Methods such as singular value decomposition, non-negative matrix factorization, or principal component analysis produce compact yet informative representations that preserve essential spatio-temporal characteristics while reducing redundancy, thereby improving both discriminative capability and computational efficiency. To support adaptive optimization, RL is incorporated through a QL mechanism, as illustrated in Figure 3. In this configuration, the model functions as an agent that receives reward feedback based on prediction performance. Incorrect predictions generate corrective rewards, directing iterative Q-value updates and policy refinement. Through continuous interaction and feedback, the model progressively improves learning stability and classification accuracy. Complementing this process, MIM, shown in Figure 4, is employed to reinforce feature relevance across modalities. By maximizing shared information among complementary feature representations, MIM ensures that retained features remain informative and non-redundant, ultimately improving representation quality and downstream performance.

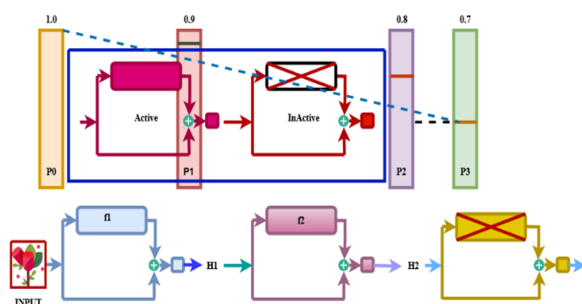


Figure 1. Flow diagram of SDNN

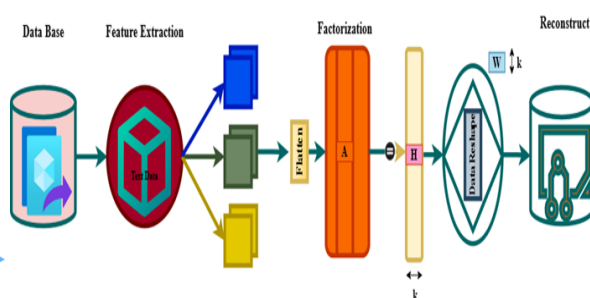


Figure 2. Flow diagram of DFF

2.2. Data sources and preprocessing

This study utilizes a multivariate human gait dataset sourced from the UCI machine learning repository, released on December 14, 2022. The dataset consists of 181,800 time-series samples acquired from 10 healthy participants performing gait under three experimental conditions: unbraced, knee-braced, and ankle-braced. Under each condition, participants completed 10 gait cycles, with joint-angle trajectories captured at 101 discrete time points corresponding to a complete gait cycle. Each sample is characterized by seven attributes, including subject identifier, walking condition, replication index, leg side, joint type (ankle, knee, or hip), time step, and joint angle expressed in degrees. Data acquisition was conducted at the Human Dynamics and Controls Laboratory, University of Illinois at Urbana–Champaign [22]–[24], and the dataset contains no missing entries. The balanced distribution across subjects, walking conditions, limbs, and joint types makes the dataset appropriate for both regression and classification tasks in biomechanical gait analysis. For the joint-angle regression task, the objective was to estimate continuous joint-angle values using subject-specific and gait-related features. Data preprocessing involved loading the dataset with Pandas, encoding categorical variables, and normalizing numerical features using MinMaxScaler. A feature matrix comprising 29 predictors was formed, with joint angle designated as the regression target. The dataset was subsequently split into training and testing subsets and reshaped into sequential formats compatible with the SDNN-based regression architecture. For multi-class joint-type classification, the aim was to identify joint categories using the same input attributes. Joint labels were one-hot encoded, numerical features were normalized, and a dataset containing 27 input features was constructed using the identical train–test split. The classification data were then arranged into structured sequences suitable for the DFFN-based architecture. Overall, these preprocessing procedures produced clean, balanced, and well-organized datasets, establishing a reliable basis for accurate and interpretable gait analysis across varying walking conditions.

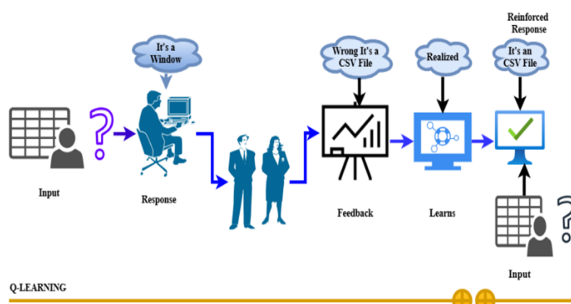


Figure 3. Flow diagram of QL

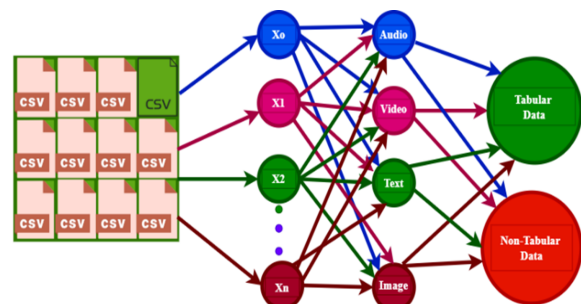


Figure 4. Flow diagram of MIM

2.3. Model architecture and justification

This study proposes a unified framework that integrates neural networks (NN), RL, and XAI to address joint-angle regression and multi-class joint-type classification in human gait analysis. The overall workflow starts with dataset preparation, where noise and outliers are managed, categorical variables are encoded, and numerical features are normalized using Min–Max scaling. The processed data are then partitioned into training and testing sets to enable balanced and unbiased evaluation. For joint-angle regression, two variants of the SDNN are developed. The SDNN_QL model incorporates QL to support policy-driven optimization during training, while the SDNN_MIM model applies MIM to enhance feature representation and generalization performance. Both variants are designed to effectively capture temporal gait dynamics while minimizing prediction error in joint-angle estimation. Regression performance is assessed using mean squared error (MSE), mean absolute error (MAE), and the coefficient of determination (R^2), complemented by residual and performance plots that assist in validating learning stability and predictive reliability. For multi-class joint-type classification, two DFFN variants are utilized. The DFFN_QL model integrates QL to optimize action-selection behavior during classification, whereas the DFFN_MIM model employs MIM to reinforce learned feature embeddings. These models are trained to discriminate among ankle, knee, and hip joint categories. Classification performance is measured using accuracy, precision, recall, F1 score, and prediction probability distributions, with additional insights derived from confusion matrices, ROC curves, and precision–recall plots. To enhance

transparency and clinical interpretability, the framework incorporates multiple XAI techniques. Counterfactual explanations identify minimal changes in input features required to modify predictions, while contrastive explanations highlight differences between predicted outcomes and alternative classes. Feature ablation analysis evaluates the contribution of individual input variables, and prediction probability analysis demonstrates model confidence across both regression and classification tasks. These interpretability findings are presented through visual and textual representations to support clear understanding of model decision-making. The complete architecture is shown in Figure 5, where Figure 5(a) illustrates the SDNN_QL_MIM regression model and Figure 5(b) displays the DFFN_QL_MIM classification model.

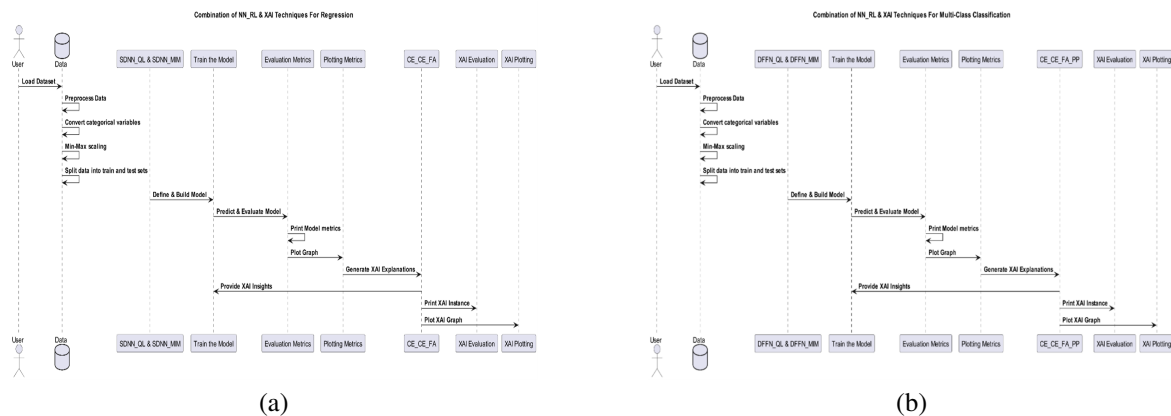


Figure 5. Model architectures (a) SDNN_QL_MIM for regression and (b) DFFN_QL_MIM for multi-class classification

2.4. Performance metrics

The proposed gait analysis framework is assessed using standard performance metrics suitable for both joint-angle regression and multi-class joint-type classification. These metrics are selected to capture prediction accuracy, learning stability, and generalization capability, which are critical for clinically dependable evaluation using the SDNN_QL_MIM and DFFN_QL_MIM models. For joint-angle regression, model performance is evaluated using MSE, MAE, and the coefficient of determination (R^2). MSE places greater emphasis on larger discrepancies between predicted and actual joint-angle values, whereas MAE offers a more intuitive measure of average prediction error. The R^2 metric reflects how effectively the model explains variance in joint-angle data, enabling meaningful comparison across different regression models and optimization strategies. For multi-class joint-type classification, evaluation concentrates on accuracy, precision, recall, F1 score, and prediction probability distributions. Accuracy represents overall classification effectiveness, while precision and recall characterize class-specific reliability and sensitivity. The F1 score balances these measures to provide a unified performance indicator. To further analyze class-level behavior and decision boundaries, confusion matrices, receiver operating characteristic (ROC) curves, and precision-recall plots are utilized. Collectively, these metrics provide a comprehensive evaluation of the robustness and effectiveness of the proposed gait prediction framework.

2.5. Integration of XAI techniques

The proposed framework incorporates multiple XAI techniques to improve transparency and confidence in black-box learning models applied to human gait analysis. When interpretability is needed, input data are preprocessed and forwarded through the trained model to obtain predictions. Counterfactual explanations are subsequently generated by identifying minimal and plausible modifications in the input that result in different prediction outcomes, ensuring clinical relevance. In parallel, contrastive explanations are utilized to compare the predicted outcome with alternative scenarios, thereby emphasizing the key features that drive model decisions. To further examine feature relevance, feature ablation is performed by systematically removing or perturbing individual input variables and analyzing the resulting variations in model outputs. This procedure enables a quantitative evaluation of feature importance. In the multi-class classification setting, prediction probability analysis is applied to assess class-wise confidence levels and determine the features that most strongly

influence the predicted joint category. For instance, when a sample is classified as Joint Class 2, the associated probability scores reflect the relative contribution of the corresponding input features (X variables). Collectively, these XAI techniques deliver clear and actionable insights into model behavior. When combined with RL-based decision refinement and mutual information-guided feature optimization, the framework enables accurate, interpretable, and clinically meaningful joint-angle prediction and joint-type classification.

2.6. Real-world implications

The proposed framework, integrating deep learning with RL and XAI, demonstrates strong practical relevance for biomechanics, rehabilitation engineering, prosthetics, and CGA. Accurate joint-angle prediction and joint-type classification can support clinicians in the early detection of movement disorders, enable personalized rehabilitation strategies, and contribute to the design of more effective prosthetic and assistive devices. The incorporation of RL allows the models to adapt to evolving gait patterns and sustain stable performance across varying walking conditions. Moreover, the inclusion of XAI techniques—such as counterfactual and contrastive explanations, feature ablation, and prediction probability analysis—enhances transparency by enabling clinicians and domain experts to interpret and validate model predictions. This level of interpretability addresses common concerns related to black-box learning models and promotes responsible clinical deployment. By unifying adaptive learning with explainable decision-making, the proposed framework provides a practical basis for implementing intelligent gait analysis systems in real-world environments. As data-driven human movement analysis continues to advance, such adaptive and explainable approaches are expected to play a significant role in the development of assistive technologies and evidence-based healthcare solutions.

2.7. Mathematical formulation

This section presents concise mathematical formulations of the XAI techniques used in this study, namely counterfactual explanations, contrastive explanations, and feature ablation. These formulations describe how minimal input perturbations influence model predictions and enable transparent interpretation for both regression and multi-class classification tasks.

2.8. Counterfactual explanations

Counterfactual explanations identify the minimal modification to an input instance that changes the model's prediction. Input features are normalized using Min-Max scaling is defined as (1):

$$x_{\text{norm}} = \frac{x - x_{\min}}{x_{\max} - x_{\min}} \quad (1)$$

The counterfactual objective is defined by minimizing a loss function that shifts the prediction from the original output to a target outcome is defined as (2):

$$L(x) = -P(y_{\text{target}} | x) + P(y_{\text{orig}} | x) \quad (2)$$

The optimal counterfactual instance is obtained as (3):

$$x^* = \arg \min_x L(x) \quad (3)$$

2.9. Contrastive explanations

Contrastive explanations analyze how small perturbations in the input alter the model's prediction. A contrastive instance is generated by adding bounded Gaussian noise as (4):

$$x_{\text{con}} = \text{clip}(x + \mathcal{N}(0, \sigma^2), 0, 1) \quad (4)$$

The model prediction for both original and perturbed inputs is given by (5):

$$\hat{y} = f(x) \quad (5)$$

Differences between these predictions highlight features that most strongly influence decision boundaries.

2.10. Feature ablation

Feature ablation evaluates the importance of individual features by measuring prediction changes after feature removal. For a given feature j , the perturbed input is defined as (6):

$$X' = X \text{ with } X[:, j] = 0 \quad (6)$$

The impact of the ablated feature is quantified by the absolute prediction difference as (7):

$$L_j = |f(X) - f(X')| \quad (7)$$

To enable fair comparison across features, the ablation scores are normalized as (8):

$$L_j^{\text{norm}} = \frac{L_j}{\sqrt{\sum_{j=1}^n L_j^2}} \quad (8)$$

Higher normalized scores indicate greater influence of the corresponding feature on the model's output.

2.11. Hyperparameter tuning strategy

Hyperparameter tuning was conducted independently for the joint-angle regression and multi-class joint-type classification tasks to ensure stable convergence and dependable model performance. For the regression task, the SDNN model was trained using a test split of 0.3 and a fixed random seed of 42 to guarantee reproducibility. The network architecture comprised five stochastic depth layers with a survival probability of 0.8. Each hidden layer included 32 neurons with ReLU activation, while a linear activation function was employed at the output layer to enable continuous joint-angle prediction. Model optimization was carried out using the Adam optimizer, which supported training stability and reduced overfitting. For the multi-class joint-type classification task, the DFFN model defined target variables as features beginning with `joint_`, with 30% of the dataset allocated for testing and the same random seed of 42. The architecture incorporated a feature factorization layer with 512 neurons, followed by interaction layers consisting of 256, 128, and 64 neurons. Additional non-linear transformation layers with 128 and 64 neurons were included, and a dropout rate of 0.4 was applied to enhance generalization. The final softmax layer contained three neurons corresponding to the joint-type classes. Training was performed using the Adam optimizer with an initial learning rate of 0.0001, exponential decay steps of 10,000, a decay rate of 0.8, staircase decay enabled, and categorical cross-entropy as the loss function to ensure stable and reliable classification.

Tables 1 and 2 summarize the experimental configurations applied for the regression and multi-class classification tasks, respectively. Across all experiments, deep learning and RL parameters were maintained consistently to ensure fair comparison across different XAI techniques. To support interpretability, XAI explanations were generated for both the initial and final predictions.

Table 1. RL parameter settings for regression (QL vs. MIM)

Parameters	QL-regression	MIM-regression
Total training epochs for RL model	30	30
Batch size for training	64	64
Initial exploration rate (ϵ)	0.5	0.5
Exploration decay rate	0.99	0.99
Discount factor (γ)	0.95	0.95
Frequency of updating target model	10	5
Target model for RL updates	–	Clone of main model
Possible learning rate values	[0.00001, 0.00005, 0.0001, 0.0005, 0.001]	–
Possible dropout rate values	[0.2, 0.3, 0.3, 0.4, 0.5]	–
Possible action values	–	[(0.00001, 0.2, 0.6), (0.00005, 0.3, 0.7), (0.0001, 0.3, 0.8), (0.0005, 0.4, 0.9), (0.001, 0.5, 0.9)]
Learning rate for Q-table updates	0.5	–
Number of features in training set	–	$X_{\text{train}}.\text{shape}[1]$
Counter for successful episodes	–	0
Reward function	$1/(1 + \text{MSE})$	$1/(1 + \text{MSE})$
Maximum reward value	1.0	1.0
Reward threshold for success count	0.8	0.8
Verbosity level	0	0

Table 2. RL parameter settings for multi-class classification (QL vs. MIM)

Parameters	QL-multi class	MIM-multi class
Number of training epochs	50	50
Batch size for training	64	64
Initial exploration rate (ϵ)	0.9	0.9
Exploration decay rate	0.98	0.98
Discount factor (γ)	0.99	0.99
Frequency of updating target model	10 epochs	10
Possible action values	[(0.00001, 0.3, 128), (0.00005, 0.4, 256), (0.0001, 0.4, 512), (0.0005, 0.5, 1024), (0.001, 0.6, 2048)]	[(0.00001, 0.3, 128), (0.00005, 0.4, 256), (0.0001, 0.4, 512), (0.0005, 0.5, 1024), (0.001, 0.6, 2048)]
Learning rate for Q-table updates	0.9	–
Number of features in dataset	–	$X_{\text{train}}.\text{shape}[1]$
Scaling factor for intrinsic reward	–	0.5
Dropout rate in hidden layers	–	0.4
Number of neurons in interaction layers	–	128, 256, 512, 1024, 2048

3. RESULTS AND DISCUSSION

3.1. Experimental setup

The experimental setup utilizes advanced DL, RL, and XAI techniques to support efficient and robust gait analysis. Data preprocessing and performance evaluation were performed using the scikit-learn library, while deep neural network architectures were designed and trained with TensorFlow/Keras. RL components were incorporated to enable adaptive optimization during training, and XAI techniques were integrated to improve transparency and interpretability. This unified setup facilitates reliable joint-angle regression and multi-class joint-type classification with clinically meaningful insights.

3.2. Exploratory data analysis and feature insights

Figure 6 illustrates a lollipop chart summarizing the mean values of all input features. The time feature shows the highest mean value (approximately 50), followed by the angle feature (approximately 12.15), indicating their dominant numerical magnitude within the dataset. In contrast, features such as subject, condition, replication, leg, and joint exhibit lower mean values (ranging from 1 to 5), reflecting their categorical or discrete nature. Figure 7 presents a line plot with error bars representing the mean and standard deviation of each feature. The time feature demonstrates both the highest mean and the greatest variability, whereas angle shows moderate variation. The remaining features display shorter error bars, indicating limited variability consistent with categorical attributes. Figure 8 depicts a hexbin plot visualizing the joint distribution of Class and Hypertension, where color intensity denotes data density. This visualization emphasizes dominant class–hypertension combinations while minimizing visual clutter from individual data points. Finally, the correlation matrix in Figure 9 indicates generally weak linear relationships among features, with a modest positive correlation (0.22) identified between time and angle. The overall low linear dependency supports the application of nonlinear and multivariate modeling approaches to capture complex gait dynamics.

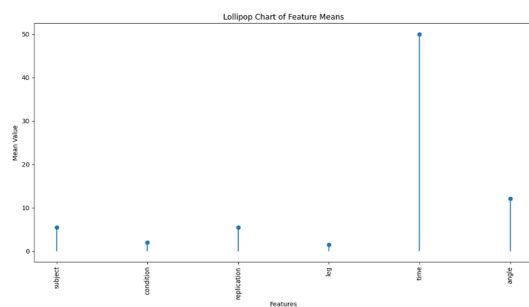


Figure 6. Lollipop chart of feature means

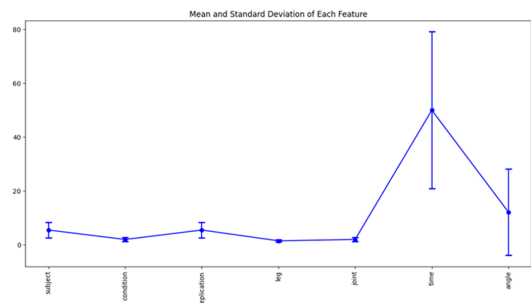


Figure 7. Mean and standard deviation for each feature

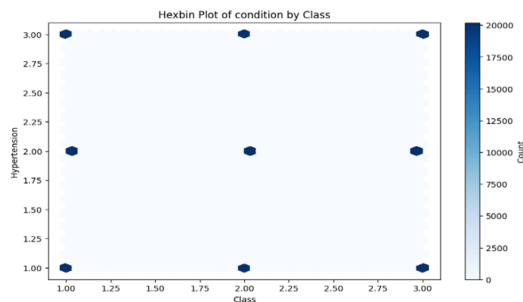


Figure 8. Hexbin plot

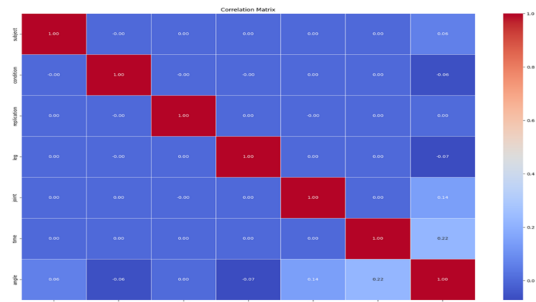


Figure 9. Correlation matrix

3.3. Regression performance analysis

Figure 10 provides a comparative evaluation of integrated NN and RL-based regression models, where the SDNN framework is optimized using QL and MIM. In Figure 10(a), the QL-based model displays a gradual rise in reward values from approximately 0.992 to 0.998 over 30 epochs, indicating steady performance improvement with minor fluctuations. In contrast, Figure 10(b) illustrates that the MIM-based model converges more quickly, increasing from about 0.997 to nearly 0.999 within the same epoch range. Overall, although both optimization strategies demonstrate effective learning behavior, SDNN_MIM achieves faster convergence and marginally higher reward values than SDNN_QL, indicating superior optimization efficiency for joint-angle regression tasks.

Figure 11 presents a comparative assessment of integrated NN and RL-based regression models for joint-angle prediction, specifically SDNN_QL and SDNN_MIM. Performance is evaluated using MSE, MAE, and R^2 . The SDNN_MIM model records lower errors (MSE = 0.0003, MAE = 0.0125) compared to SDNN_QL (MSE = 0.0006, MAE = 0.0183) and achieves a higher R^2 score (0.9881 vs. 0.9750), reflecting improved variance explanation and model fit. These findings suggest that MIM strengthens feature learning and regression accuracy, whereas QL is relatively less effective. Overall, SDNN_MIM emerges as the most effective model for joint-angle regression, while maintaining strong interpretability.

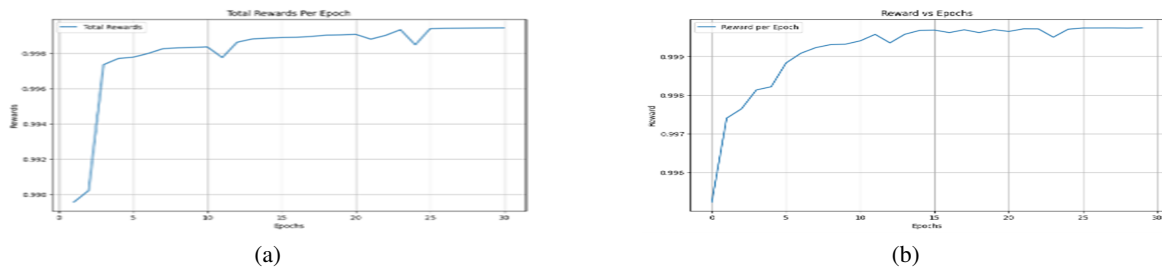


Figure 10. Model performance analysis (a) SDNN_QL_MIM for regression and (b) SDNN_QL_MIM for regression

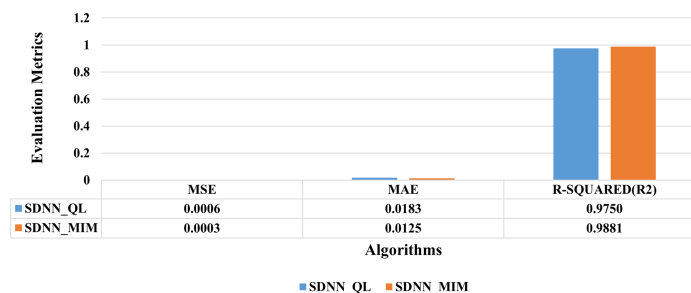


Figure 11. Comparative analysis of combined NN and RL-based regression models

Figure 12 integrates regression performance analysis with XAI-based explanations. In Figure 12(a), counterfactual explanations of the SDNN_QL_MIM model analyze feature contributions across the first, second, last, and last-to-first predictions under QL and MIM. For QL, the initial prediction is mainly driven by *time*, *joint_2*, *leg_2*, *replication_1*, *condition_3*, and *subject_9*, while the final prediction shifts toward *joint_3*, *replication_3*, *condition_3*, and *subject_4*. Under MIM, the second prediction emphasizes *time*, *joint_3*, *leg_2*, *replication_8*, *condition_3*, and *subject_10*, whereas the last-to-first prediction highlights *time*, *joint_2*, *leg_1*, *replication_4*, *condition_3*, and *subject_9*. Across all predictions, *time* and *condition_3* consistently emerge as dominant features. Figure 12(b) presents contrastive explanations that further examine feature variations across prediction stages. For QL, the first prediction is influenced by *time*, *joint_2*, *leg_2*, *replication_1*, *condition_3*, and *subject_9*, while the final prediction shifts toward *joint_3*, *replication_3*, *condition_3*, and *subject_4*. Under MIM, the second prediction highlights *time*, *joint_3*, *leg_2*, *replication_8*, *condition_3*, and *subject_10*, whereas the last-to-first prediction emphasizes *time*, *joint_2*, *leg_1*, *replication_4*, *condition_3*, and *subject_9*. These findings indicate stable temporal and condition-related features, with other variables adapting based on the learning strategy. In Figure 12(c), feature ablation analysis assesses feature importance through sensitivity comparisons across predictions. For QL, the initial prediction is affected by *time*, *joint_2*, *leg_2*, *replication_1*, *condition_3*, and *subject_9*, while the final prediction shifts toward *joint_3*, *replication_3*, *condition_3*, and *subject_4*. For MIM, the second prediction is influenced by *time*, *joint_3*, *leg_2*, *replication_8*, *condition_3*, and *subject_10*, whereas the last-to-first prediction highlights *time*, *joint_2*, *leg_1*, *replication_4*, *condition_3*, and *subject_9*. Across all XAI techniques, *time* and *condition_3* consistently emerge as the most dominant and stable features influencing the target variable (*angle*). Overall, temporal and condition-related factors govern prediction stability, while joint, leg, replication, and subject identifiers contribute adaptively to model refinement in gait joint-angle regression.

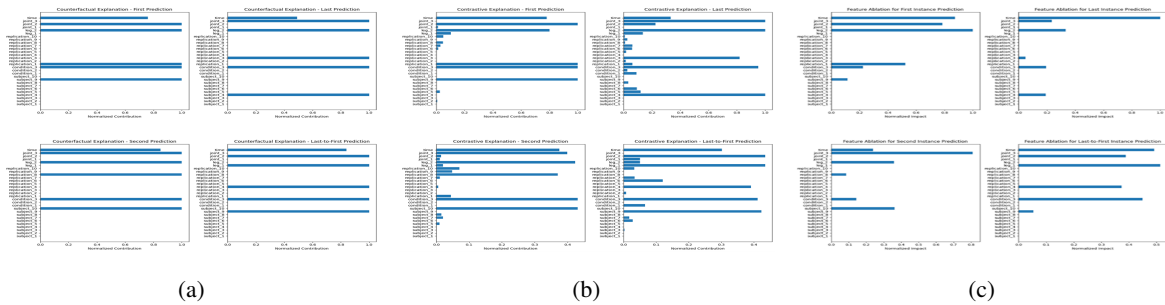


Figure 12. Regression performance analysis with XAI-based explanations: (a) counterfactual explanations, (b) contrastive explanations, and (c) feature ablation

3.4. Multi-class classification performance analysis

Figure 13 presents a comparative analysis of integrated NN- and RL-based multi-class classification models using DFNN optimized with QL and MIM across 50 epochs. In Figure 13(a), the QL-based model exhibits a gradual and oscillatory increase in reward values from approximately 0.70 to 0.96, indicating slower and less stable convergence. In contrast, Figure 13(b) shows that the MIM-based model rapidly exceeds 0.90 within the first 10 epochs and stabilizes around 0.98 by epoch 50. Overall, while both optimization strategies demonstrate effective learning behavior, MIM achieves faster convergence and greater learning stability, making it a more efficient optimization strategy than QL for multi-class joint-type classification.



Figure 13. Model performance analysis: (a) DFNN_QL_MIM for multi-class classification and (b) DFNN_MIM_MIM for multi-class classification

Figure 14 presents a comparative assessment of combined NN and RL-based classification models for multi-class joint-type classification, specifically DFFN_QL and DFFN_MIM, evaluated using accuracy, precision, recall, F1 score, and Cohen's Kappa. The DFFN_MIM model attains higher accuracy (0.95 vs. 0.94), recall (0.95 vs. 0.94), and F1 score (0.95 vs. 0.94), while both models achieve identical precision (0.95) and Cohen's Kappa (0.92). These findings suggest that MIM improves feature representation and class discrimination relative to QL. Overall, DFFN_MIM exhibits superior and more consistent performance for multi-class joint-type classification, while preserving strong interpretability.

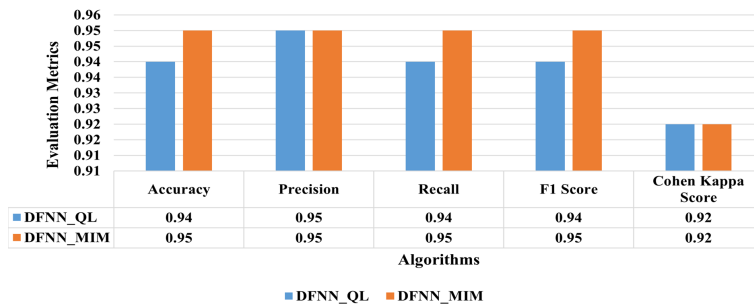


Figure 14. Comparative analysis of combined NN and RL-based multi-class classification models

Figure 15 illustrates the integration of multi-class classification performance with XAI-based explanations. In Figure 15(a), counterfactual explanations for DFFN_QL and DFFN_MIM examine feature contributions across four prediction stages. For DFFN_QL, the first prediction is influenced by *angle*, *time*, *leg_2*, *replication_1*, *condition_3*, and *subject_9*, whereas the final prediction shifts toward *replication_3*, *condition_3*, and *subject_4*, with *angle* and *time* remaining dominant. In DFFN_MIM, the second prediction emphasizes *angle*, *time*, *leg_2*, *replication_8*, and *subject_10*, while the last-to-first prediction highlights *leg_1*, *replication_4*, and *subject_9*. Across all predictions, *angle*, *time*, and *condition_3* consistently emerge as the most influential features. Figure 15(b) presents contrastive explanations that further analyze feature variations across identical prediction stages. For DFFN_QL, the first prediction is affected by *angle*, *time*, *leg_2*, *replication_1*, *condition_3*, and *subject_9*, whereas the final prediction shifts toward *replication_3*, *condition_3*, and *subject_4*. Under DFFN_MIM, the second prediction emphasizes *angle*, *time*, *leg_2*, *replication_8*, and *subject_10*, while the last-to-first prediction highlights *leg_1*, *replication_4*, and *subject_9*. These findings indicate that biomechanical and temporal features remain stable, whereas other variables adapt across learning strategies. In Figure 15(c), feature ablation analysis assesses feature importance through sensitivity comparisons across prediction stages. For DFFN_QL, the first prediction is most sensitive to the removal of *angle*, while the final prediction is strongly influenced by *time*, *leg_2*, *replication_3*, and *subject_4*. In DFFN_MIM, the second prediction is affected by *angle*, *leg_2*, *condition_3*, and *subject_10*, whereas the last-to-first prediction is dominated by *angle*. Across all XAI techniques, *time* and *condition_3* consistently emerge as the most dominant and stable features influencing the target variable (*joint*). Overall, temporal and condition-related factors drive prediction stability, while joint, leg, replication, and subject identifiers contribute adaptively and dynamically to model refinement in gait multi-class joint-type classification.

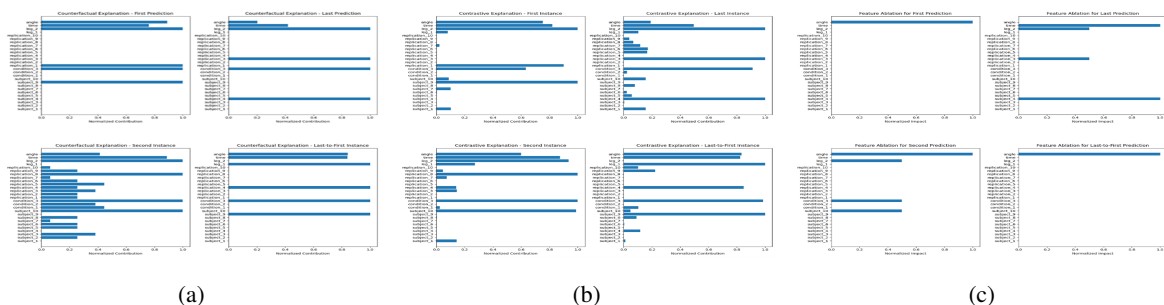


Figure 15. Multi-class classification performance analysis with XAI-based explanations: (a) counterfactual explanations, (b) contrastive explanations, and (c) feature ablation

Figure 16 integrates the multi-class classification performance analysis with prediction probability analysis for the DFFN_QL_MIM model, demonstrating high confidence across all four test samples in the multi-class classification task as shown in Figures 16(a) and (b). For the first and last-to-first samples, the model assigns a probability of 1.0 to Class 1, whereas for the second and last samples it assigns a probability of 1.0 to Class 2, with all remaining classes receiving zero probability. These results indicate highly decisive and unambiguous class predictions for each instance. Overall, Class 1 emerges as the most frequently predicted class, highlighting its dominant contribution to the target variable joint in gait joint-type classification.

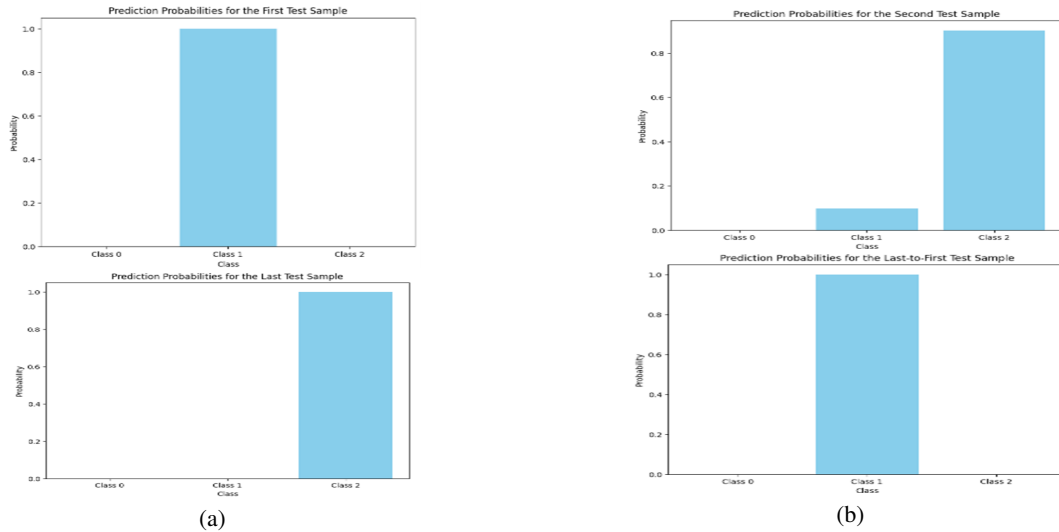


Figure 16. Prediction probabilities analysis (a) DFFN_QL for multi-class classification and (b) DFFN_MIM for multi-class classification

Table 3 indicates that the proposed hybrid framework exceeds existing studies in both regression and multi-class classification tasks. By integrating NN with RL and XAI, the proposed method delivers higher predictive accuracy and enhanced model transparency compared to previous approaches based on classical machine learning, deep learning without explainability, or XAI without RL, in line with the findings reported in [1], [11], [14], [25].

Table 3. Comparative summary of performance metrics across gait analysis studies

Study	Hybrid approach	MSE (regression)	MAE (regression)	R^2 (regression)	Accuracy (binary)	Accuracy (multi)	Reward (regression)	Reward (multi)
[1]	No	–	–	–	–	0.92	–	–
[25]	Yes	3.31	–	0.99	0.97	0.94	–	–
[11]	Yes	–	–	–	0.94	–	–	–
[14]	Yes	–	3.8	–	–	–	–	–
Proposed study	Yes	0.0003	0.0125	0.9881	–	0.95	0.999	0.98

3.5. Limitations and future work

Although the proposed framework demonstrates strong predictive performance and enhanced interpretability, its generalizability is affected by the reliance on gait data collected from healthy participants under controlled experimental conditions. Furthermore, the NN, RL, and XAI components are assessed in an offline setting, necessitating additional validation in real-time and clinically diverse environments. While counterfactual, contrastive, feature ablation, and prediction probability analyses improve transparency, their direct influence on clinical decision-making has not yet been empirically evaluated. Future research will therefore concentrate on exploring more adaptive RL strategies, extending validation to clinical populations with gait impairments, and enabling real-time deployment through wearable sensing technologies. These efforts are expected to enhance robustness, learning stability, and personalization across diverse gait conditions, thereby strengthening the practical and clinical relevance of the proposed framework.

4. CONCLUSION

This study investigates the integration of deep learning, RL, and XAI for joint-angle prediction and joint-type classification in human gait analysis. Experimental findings show that models optimized through MIM exhibit more consistent performance than QL-based models in both regression and multi-class classification tasks, particularly in modeling complex gait patterns. The incorporation of XAI techniques, including counterfactual and contrastive explanations, feature ablation, and prediction probability analysis, enhances model transparency by clarifying feature-level and class-level influences on prediction outcomes. These findings are significant for applications in rehabilitation, prosthetic control, and human motion analysis, where accurate predictions and interpretable models are essential for clinical decision-making. Although the present work relies on controlled experimental data collected from healthy individuals, it establishes a valuable foundation for interpretable gait modeling. Future research will concentrate on extending the proposed framework to clinical populations, assessing performance in real-world environments, exploring more adaptive RL strategies, and enabling real-time deployment using wearable sensing systems. Overall, this study contributes toward the development of more accurate, interpretable, and practically applicable gait analysis methods for personalized healthcare solutions.

FUNDING INFORMATION

Authors state no funding involved

AUTHOR CONTRIBUTIONS STATEMENT

This journal uses the Contributor Roles Taxonomy (CRediT) to recognize individual author contributions, reduce authorship disputes, and facilitate collaboration.

Name of Author	C	M	So	Va	Fo	I	R	D	O	E	Vi	Su	P	Fu
Deepak N. R.	✓		✓	✓		✓		✓	✓				✓	
Soumya Naik P. T.		✓						✓		✓		✓		
Ambika P. R.	✓		✓	✓		✓			✓		✓		✓	
Shaik Sayeed Ahamed	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

C : **C**onceptualization

M : **M**ethodology

So : **S**oftware

Va : **V**alidation

Fo : **F**ormal Analysis

I : **I**nvestigation

R : **R**esources

D : **D**ata Curation

O : Writing - **O**riginal Draft

E : Writing - Review & **E**ditting

Vi : **V**isualization

Su : **S**upervision

P : **P**roject Administration

Fu : **F**unding Acquisition

CONFLICT OF INTEREST STATEMENT

Authors state no conflict of interest.

DATA AVAILABILITY

Data supporting the findings of this study are available from the corresponding author, Shaik Sayeed Ahamed, upon reasonable request. The dataset is based on publicly available biomechanical gait recordings cited in [22]–[24], with preprocessing applied for FL experimentation as detailed in Section 2.2.

REFERENCES

- [1] D. Slijepcevic, F. Horst, S. Lapuschkin, B. Horsak, A. M. Raberger, A. Kranzl, and M. Zeppelzauer, "Explaining machine learning models for clinical gait analysis," *ACM Transactions on Computing for Healthcare (HEALTH)*, vol. 3, no. 2, pp. 1–27, 2021, doi: 10.1145/3474121.
- [2] R. Madanu, M. F. Abbod, F. J. Hsiao, W. T. Chen, and J. S. Shieh, "Explainable AI applied in machine learning for pain modeling: a review," *Technologies*, vol. 10, no. 3, p. 74, 2022, doi: 10.3390/technologies10030074.
- [3] P. Khera and N. Kumar, "Role of machine learning in gait analysis: A review," *Journal of Medical Engineering and Technology*, vol. 44, no. 8, pp. 441–467, 2020, doi: 10.1080/03091902.2020.1822940.





- [4] L. Lei, Y. Tan, K. Zheng, S. Liu, K. Zhang, and X. Shen, "Deep reinforcement learning for autonomous internet of things: models, applications, and challenges," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 3, pp. 1722–1760, 2020, doi: 10.1109/COMST.2020.2988367.
- [5] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, and S. Levine, "Soft actor-critic algorithms and applications," *arXiv preprint arXiv:1812.05905*, 2018.
- [6] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: off-policy maximum entropy deep reinforcement learning with a stochastic actor," in *Proceedings of the International Conference on Machine Learning (ICML)*, 2018, pp. 1861–1870.
- [7] T. Haarnoja, S. Ha, A. Zhou, J. Tan, G. Tucker, and S. Levine, "Learning to walk via deep reinforcement learning," *arXiv preprint arXiv:1812.11103*, 2018.
- [8] M. S. Islam, I. Hussain, M. M. Rahman, S. J. Park, and M. A. Hossain, "Explainable artificial intelligence model for stroke prediction using EEG signal," *Sensors*, vol. 22, no. 24, p. 9859, 2022, doi: 10.3390/s22249859.
- [9] A. Plaat, W. Kusters, and M. Preuss, "High-accuracy model-based reinforcement learning: a survey," *Artificial Intelligence Review*, vol. 56, no. 9, pp. 9541–9573, 2023, doi: 10.1007/s10462-022-10335-w.
- [10] K. M. Hossain and T. Oates, "Ten-Guard: tensor decomposition for backdoor attack detection in deep neural networks," in *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, IEEE, 2024.
- [11] A. Mishra, A. Shetkar, G. M. Bapat, R. Ojha, and T. T. Verlekar, "XAI-based gait analysis of patients walking with knee-ankle-foot orthosis using video cameras," *arXiv preprint arXiv:2402.16175*, 2024.
- [12] M. Haklidiir and H. Temeltaş, "Guided soft actor critic: a guided deep reinforcement learning approach for partially observable Markov decision processes," *IEEE Access*, vol. 9, pp. 159672–159683, 2021, doi: 10.1109/ACCESS.2021.3131772.
- [13] M. E. Özates, A. Yaman, F. Salami, S. Campos, S. I. Wolf, and U. Schneider, "Identification and interpretation of gait analysis features and foot conditions by explainable AI," *Scientific Reports*, vol. 14, no. 1, p. 5998, 2024.
- [14] J. B. Ko, J. S. Hong, Y. S. Shin, and K. B. Kim, "Machine learning-based predicted age of the elderly using instrumented timed up and go and six-minute walk tests," *Sensors*, vol. 22, no. 16, p. 5957, 2022, doi: 10.3390/s22165957.
- [15] G. Nicora, S. Pe, G. Santangelo, L. Billeci, I. G. Aprile, M. Germanotta, and S. Quaglini, "A systematic review of machine learning in robotics-assisted rehabilitation," 2024.
- [16] H. Van Eetvelde, L. D. Mendonça, C. Ley, R. Seil, and T. Tischer, "Machine learning methods in sport injury prediction and prevention: a systematic review," *Journal of Experimental Orthopaedics*, vol. 8, pp. 1–15, 2021.
- [17] D. Casacuberta, A. Guersenzvaig, and C. Moyano-Fernández, "Justificatory explanations in machine learning: increasing transparency through documented design decisions," *AI & Society*, vol. 39, no. 1, pp. 279–293, 2024.
- [18] F. Rita, L. De Santis, and I. F. Felice, "XAI for supporting gait analysis of patients with schizophrenia," 2024.
- [19] S. Fleischmann, S. Dietz, J. Shanbhag, A. Wuensch, M. Nitschke, J. Miehl, and A. D. Koelewijn, "Exploring dataset bias and scaling techniques in multi-source gait biomechanics: an explainable machine learning approach," *ACM Transactions on Intelligent Systems and Technology*, vol. 16, no. 1, pp. 1–19, 2024.
- [20] D. Slijepcevic, *Human gait analysis: machine learning-based classification of gait disorders*, Ph.D. dissertation, Technische Universität Wien, 2024.
- [21] S. Milani, N. Topin, M. Veloso, and F. Fang, "Explainable reinforcement learning: a survey and comparative review," *ACM Computing Surveys*, vol. 56, no. 7, pp. 1–36, 2024, doi: 10.1145/3616864.
- [22] K. A. Shorter, J. D. Polk, K. S. Rosengren, and E. T. Hsiao-Weckslar, "A new approach to detecting asymmetries in gait," *Clinical Biomechanics*, vol. 23, no. 4, pp. 459–467, 2008, doi: 10.1016/j.clinbiomech.2007.11.009.
- [23] N. E. Helwig, S. Hong, E. T. Hsiao-Weckslar, and J. D. Polk, "Methods to temporally align gait cycle data," *Journal of Biomechanics*, vol. 44, no. 3, pp. 561–566, 2011, doi: 10.1016/j.jbiomech.2010.09.015.
- [24] N. E. Helwig, K. A. Shorter, P. Ma, and E. T. Hsiao-Weckslar, "Smoothing spline analysis of variance models: a new tool for the analysis of cyclic biomechanical data," *Journal of Biomechanics*, vol. 49, no. 14, pp. 3216–3222, 2016, doi: 10.1016/j.jbiomech.2016.07.035.
- [25] S. S. Ahamed, A. Pasha, S. Rahman, and P. Kumar D. N., "Interpretable federated deep learning models for predicting gait dynamics in biomechanics," *Indonesian Journal of Electrical Engineering and Computer Science (IJECS)*, vol. 40, no. 2, pp. 1087–1099, 2025, doi: 10.11591/ijeecs.v40.i2.pp1087-1099.

BIOGRAPHIES OF AUTHORS







Deepak N. R. is a Professor in the Department of Computer Science and Engineering at Atria Institute of Technology, Bengaluru, India. He holds a B.E., M.Tech., and Ph.D. in Computer and Information Sciences from Visvesvaraya Technological University (VTU). He is a Fellow Certified in Human Resource Development (FC-HRD) and a member of ISTE. Dr. Deepak has 17 years of academic and research experience, focusing on wireless networks, deep learning, and blockchain. He currently supervises four research scholars pursuing their Ph.D. His research spans theoretical foundations to design and implementation, collaborating with researchers across multiple domains of computer science. He has actively participated in organizing and serving on five conference committees and multiple workshops, FDPs, and SDP programs. Additionally, he has served as a Program Chair for various academic events. Dr. Deepak is a passionate educator, dedicated to fostering academic and personal excellence among students. He strives to create an engaging and challenging learning environment that nurtures lifelong learners. His approach to education is deeply invested in administrative service, committee contributions, and an achievement-oriented teaching methodology. He can be contacted at email: deepaknrgowda@gmail.com.







Sowmya Naik P. T.     is a professor and head of the Department of Computer Science and Engineering at City Engineering College, Bengaluru, affiliated to Visvesvaraya Technological University (VTU), Belgaum, Karnataka, India. She completed her B.E. in Computer Science and Engineering in 2007 and M.Tech. in Computer Science and Engineering in 2012. She obtained her Ph.D. degree in Computer Science and Engineering from Visvesvaraya Technological University (VTU), Belgaum, Karnataka. Dr. Sowmya Naik is a member of ISTE and MIE. Her areas of interest include wireless sensor networks, cloud computing, big data, and machine learning. She can be contacted at email: sowmya.vturesearch@gmail.com.



Ambika P. R.     is a professor in the Department of Computer Science and Engineering at City Engineering College, Bengaluru, Karnataka, India. She completed her B.E. and M.Tech. in Computer Science and Engineering from Visvesvaraya Technological University (VTU), Belgaum, Karnataka, and earned her Ph.D. in Computer Science and Engineering from VTU, Belagavi, Karnataka. She has more than 14 years of teaching experience, including six years of research experience, and has also worked as a Technical Associate at SAP Labs, Bengaluru, for two years. She has been actively involved in organizing conferences, Faculty Development Programs (FDPs), workshops, technical seminars, and project exhibitions. She has guided three KSCST-funded student projects and served as Co-Coordinator for an ATAL (AICTE Training and Learning) Academy-sponsored FDP in 2023. She is the author of the book internet of things and has published over 15 research papers in Scopus-indexed journals and national and international conferences. She is a member of ISTE and IAENG. Her research interests include data mining, data science, internet of things, artificial intelligence, and machine learning. She can be contacted at email: ambikatanaji@gmail.com.



Shaik Sayeed Ahamed     is an assistant professor in the Department of Computer Science and Engineering (Data Science) at Atria Institute of Technology, Bengaluru, Karnataka, India, and is currently pursuing a Ph.D. at Visvesvaraya Technological University (VTU), Belagavi. He obtained his B.Tech. and M.Tech. degrees in Computer Science and Engineering from Presidency University and REVA University, Bengaluru, respectively. He has one year of teaching experience along with more than three years of research experience. His research contributions include six publications, comprising two journal articles, one international conference paper, and three book chapters published by Springer and other reputed outlets. He has also been actively involved in organizing and participating in Skill Development Programs (SDPs), Faculty Development Programs (FDPs), workshops, and technical seminars. His research interests span data science, artificial intelligence and machine learning, computer networks, and web development. He can be contacted at email: shaik.sayeedahamed1999@gmail.com.