

OPENING THE BLACK BOX: AN EXPLAINABLE DEEP LEARNING APPROACH FOR ECG-BASED ARRHYTHMIA CLASSIFICATION

Rajender Naik Guguloth*

Assistant Professor, Dept. of EEE, University College of Engineering, Husnabad, Telangana, India.

Article Received: 27 December 2025 | Article Revised: 17 January 2026 | Article Accepted: 6 February 2026

*Corresponding Author: Rajender Naik Guguloth

Assistant Professor, Dept. of EEE, University College of Engineering, Husnabad, Telangana, India.

DOI: <https://doi.org/10.5281/zenodo.18638819>

How to cite this Article: Rajender Naik Guguloth (2026) OPENING THE BLACK BOX: AN EXPLAINABLE DEEP LEARNING APPROACH FOR ECG-BASED ARRHYTHMIA CLASSIFICATION. World Journal of Pharmaceutical Science and Research, 5(2), 423-438. <https://doi.org/10.5281/zenodo.18638819>



Copyright © 2026 Rajender Naik Guguloth | World Journal of Pharmaceutical Science and Research.

This work is licensed under creative Commons Attribution-NonCommercial 4.0 International license (CC BY-NC 4.0).

ABSTRACT

Deep learning techniques have demonstrated remarkable performance in electrocardiogram (ECG)-based arrhythmia classification; however, their adoption in clinical practice remains limited due to the lack of transparency and interpretability in model decision-making. Most existing approaches function as black boxes, offering minimal insight into how diagnostic conclusions are derived from ECG signals. To address this limitation, this study presents an explainable deep learning framework for ECG-based arrhythmia classification that combines reliable predictive performance with clinically meaningful interpretation. The proposed framework employs an end-to-end neural architecture integrating one-dimensional convolutional layers and recurrent modelling to automatically learn discriminative temporal-morphological representations from segmented ECG signals. A temporal attention mechanism is incorporated to explicitly quantify the contribution of individual heartbeats and signal regions to the final classification outcome. Furthermore, an explainability module maps model relevance back onto the ECG waveform, enabling intuitive visualization of diagnostically significant components such as the P-wave, QRS complex, and T-wave. The model is evaluated on multi-class arrhythmia classification tasks using standard performance metrics alongside detailed explainability analyses. Experimental results demonstrate that the proposed approach achieves robust classification performance while providing transparent, beat-level and segment-level explanations that are consistent with established cardiological knowledge. By opening the black box of deep ECG classification models, this work advances interpretable signal processing-driven artificial intelligence for cardiac diagnosis and supports the development of trustworthy clinical decision-support systems.

KEYWORDS: Electrocardiogram (ECG); Arrhythmia classification; Explainable artificial intelligence; Attention mechanism; Deep learning; Biomedical signal processing.

1. INTRODUCTION

Electrocardiogram (ECG) analysis plays a fundamental role in the diagnosis and monitoring of cardiac arrhythmias, providing non-invasive insights into the electrical activity of the heart (Acharya et al.^[1] With the growing availability of large-scale ECG datasets and advances in computational power, deep learning models have increasingly been adopted for automated arrhythmia classification, demonstrating superior performance compared to traditional signal processing and machine learning approaches (Rajpurkar et al.^[2]; Hannun et al.^[3]; Kiranyaz et al.^[6]) Convolutional and recurrent neural networks, in particular, have shown strong capability in learning complex temporal and morphological patterns from ECG signals (Faust et al.^[4]; Martis et al.^[23]).

Despite their promising accuracy, most deep learning-based ECG classification models operate as black boxes, offering limited insight into how diagnostic decisions are derived. This lack of transparency poses a significant barrier to clinical adoption, as cardiologists require not only reliable predictions but also an understanding of the underlying signal characteristics influencing the model's decisions. In safety-critical healthcare applications, model interpretability is essential for building trust, validating clinical relevance, and supporting informed decision-making (Lipton^[9]; Ribeiro et al.^[5]; Holzinger et al.^[25]).

Recent studies have highlighted that high classification accuracy alone is insufficient for real-world deployment of artificial intelligence (AI) systems in biomedical signal analysis. Instead, there is a growing demand for explainable artificial intelligence (XAI) techniques that can reveal the internal reasoning of deep learning models in a manner consistent with domain knowledge. In the context of ECG analysis, explainability should ideally correspond to clinically meaningful waveform components such as the P-wave, QRS complex, T-wave, and rhythm-related intervals, enabling clinicians to assess whether model decisions align with established cardiological principles (Doshi-Velez & Kim^[24]; Samek et al.^[10]).

However, achieving explainability in ECG-based deep learning models remains challenging. Many existing approaches rely on post-hoc visualization techniques that provide coarse or ambiguous explanations, while others focus on attention mechanisms as shown in Figure-1 without systematically linking model relevance to physiological signal components. As a result, there is still a lack of unified frameworks that combine robust arrhythmia classification with transparent, beat-level and segment-level interpretability (Selvaraju et al.^[12]; Lundberg & Lee.^[11]).

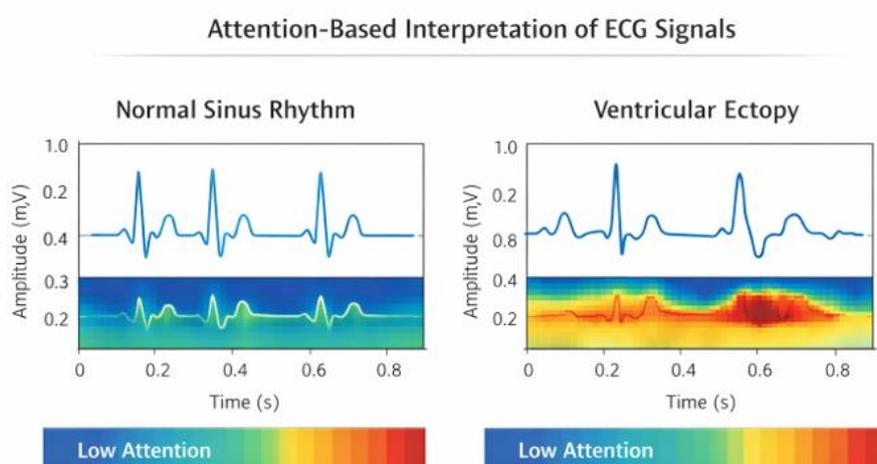


Figure 1: Attention based interpretation of ECG signals.

Motivated by these limitations, this work proposes an explainable deep learning framework for ECG-based arrhythmia classification that explicitly addresses the black-box nature of conventional models. The proposed approach integrates end-to-end feature learning with a temporal attention mechanism to identify diagnostically relevant heartbeats and signal regions. Furthermore, an explainability module is employed to map model relevance back onto ECG waveforms, enabling interpretation in terms of clinically meaningful components.

The primary contributions of this study are summarized as follows

- An explainable deep learning framework is developed for ECG arrhythmia classification, emphasizing transparency alongside predictive performance.
- A temporal attention mechanism is incorporated to quantify the contribution of individual heartbeats and signal regions to classification decisions.
- Segment-level attribution is performed to relate model focus to clinically significant ECG components, enhancing interpretability.
- Comprehensive explainability analyses and case studies are presented to demonstrate alignment between model explanations and cardiological knowledge.

By opening the black box of deep ECG classification models, this work aims to bridge the gap between advanced signal processing-driven AI techniques and clinically trustworthy decision-support systems.

2. RELATED WORK

2.1 Deep Learning for ECG Arrhythmia Classification

Deep learning techniques have been widely explored for automated ECG arrhythmia detection and classification. Early studies employed convolutional neural networks (CNNs) to capture morphological features of ECG waveforms, demonstrating improved robustness compared to handcrafted feature-based methods (Kiranyaz et al.^[6]; Acharya et al.^[11]). Recurrent neural networks (RNNs), particularly long short-term memory (LSTM) networks, have been utilized to model temporal dependencies and rhythm characteristics across sequential heartbeats (Hochreiter & Schmidhuber^[16]; Yildirim.^[7]) Hybrid CNN-RNN architectures have further enhanced classification performance by combining spatial feature extraction with temporal modelling (Zhang et al.^[18]; Rajpurkar et al.^[2])

While these approaches have achieved impressive accuracy across various arrhythmia datasets, their primary focus has been on optimizing predictive performance. The internal decision-making processes of such models are typically opaque, limiting their interpretability and raising concerns regarding reliability and clinical acceptance.

2.2 Explainable Artificial Intelligence in Biomedical Signals

Explainable artificial intelligence has emerged as a critical research area aimed at improving transparency and trust in machine learning models, particularly in healthcare applications (Samek et al.^[27]; Guidotti et al.^[28]) In biomedical signal processing, explainability methods are broadly categorized into intrinsic approaches, such as attention mechanisms, and post-hoc techniques, including gradient-based attribution and perturbation analysis.

Attention mechanisms have been introduced in ECG classification models to highlight temporally important signal regions, offering a degree of interpretability by assigning relevance weights to different time steps. However, attention

weights alone may not always correspond to physiologically meaningful features, and their clinical interpretability remains an active area of investigation (Bahdanau et al.^[13]; Vaswani et al.^[14]).

Post-hoc explainability methods such as saliency maps, Grad-CAM variants, and SHAP-based techniques have also been applied to ECG signals. These methods attempt to visualize regions of the signal that influence model predictions, but they often produce noisy or ambiguous explanations that are difficult to interpret from a clinical perspective, particularly when applied to one-dimensional biomedical signals (Selvaraju et al.^[12]; Lundberg & Lee^[11]).

2.3 Explainability in ECG-Based Deep Learning Models

Several recent studies have attempted to integrate explainability into ECG-based deep learning frameworks. Some works employ attention-based architectures to identify important heartbeats, while others apply gradient-based attribution to visualize salient waveform regions. Although these approaches represent important steps toward transparent ECG analysis, most existing methods either provide limited interpretability or lack explicit alignment with cardiological waveform components.

Moreover, many studies evaluate explainability qualitatively without systematic analysis across arrhythmia classes or without demonstrating consistency with established ECG interpretation principles. As a result, there remains a gap between explainable model outputs and clinically actionable insights (Schlemper et al.^[17]; Tonekaboni et al.^[26]).

2.4 Motivation and Research Gap

From the above review, it is evident that while deep learning has significantly advanced automated ECG arrhythmia classification, the integration of clinically meaningful explainability remains insufficiently addressed. There is a need for frameworks that not only achieve reliable classification but also provide transparent explanations that can be readily interpreted by clinicians (Doshi-Velez & Kim^[24]; Holzinger et al.^[25]).

This study addresses this gap by proposing an explainable ECG classification framework that combines attention-based intrinsic interpretability with segment-level attribution analysis. By explicitly linking model relevance to ECG waveform components, the proposed approach advances interpretable biomedical signal processing and supports trustworthy AI-assisted cardiac diagnosis.

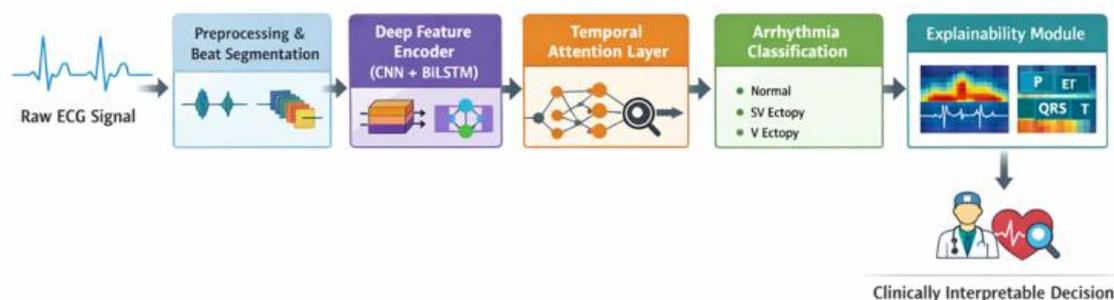
3. PROPOSED METHODOLOGY

This section describes the proposed explainable deep learning framework for ECG-based arrhythmia classification. The methodology is designed to achieve reliable classification performance while explicitly providing transparent and clinically interpretable explanations of model decisions. An overview of the proposed framework is illustrated in Figure 2.

3.1 Overview of the Proposed Framework

The proposed framework follows an end-to-end learning paradigm, integrating deep feature extraction, temporal modelling, attention-based interpretability, and segment-level explanation. The input ECG signal is first pre-processed and segmented into individual heartbeats. These segments are then processed through a deep neural encoder composed of convolutional and recurrent layers to learn discriminative representations. A temporal attention mechanism is employed to quantify the contribution of individual beats and signal regions to the classification decision. Finally, an explainability module maps the learned relevance back to the ECG waveform, enabling clinical interpretation.

Architecture of the Proposed Explainable Deep Learning Framework for ECG-Based Arrhythmia Classification

**Figure 2: Proposed Architecture for ECG based Arrhythmia Classification.**

3.2 ECG Preprocessing and Beat Segmentation

Raw ECG recordings often exhibit variations in amplitude and baseline drift due to patient-specific characteristics and acquisition conditions. To reduce inter-sample variability, the ECG signals are normalized using amplitude scaling. Basic smoothing operations are applied to suppress minor artifacts without distorting diagnostically important waveform components (Chazal et al.^[22]; Li et al.^[15]).

Heartbeat segmentation is performed by detecting R-peaks using a standard peak detection strategy. Fixed-length ECG segments are then extracted around each detected R-peak, ensuring that each segment contains complete cardiac cycles, including the P-wave, QRS complex, and T-wave. This beat-centric representation facilitates consistent temporal analysis and supports interpretability at the heartbeat level.

3.3 End-to-End Deep Feature Encoder

To automatically learn discriminative features from ECG signals, an end-to-end deep feature encoder is employed. The encoder consists of two main components:

3.3.1 Convolutional Feature Extraction

One-dimensional convolutional layers are used to capture local morphological patterns in ECG waveforms. These layers learn filters that respond to characteristic waveform shapes, such as sharp QRS complexes or subtle P-wave variations. Convolutional operations are followed by nonlinear activation functions and pooling layers to enhance feature robustness and reduce dimensionality (Kiranyaz et al.^[6]; Acharya et al.^[1]).

3.3.2 Temporal Modelling Using Recurrent Layers

Following convolutional feature extraction, bidirectional long short-term memory (BiLSTM) layers are applied to model temporal dependencies across consecutive heartbeats. The bidirectional structure enables the model to consider both past and future contextual information, which is particularly important for capturing rhythm irregularities and inter-beat dependencies associated with arrhythmias (Hochreiter & Schmidhuber^[16]; Yildirim.^[7])

The output of the BiLSTM layers consists of a sequence of hidden states representing temporally enriched ECG features.

3.4 Temporal Attention Mechanism

While deep encoders effectively capture ECG features, not all heartbeats contribute equally to the final classification decision. To explicitly model this variability, a temporal attention mechanism is introduced.

The attention layer assigns a normalized importance weight to each temporal hidden state produced by the recurrent layers. These weights reflect the relative contribution of different heartbeats and signal regions to the classification outcome. A context vector is then computed as a weighted sum of the hidden states, emphasizing diagnostically relevant segments (Bahdanau et al.^[13]; Qin et al.^[18])

This attention-driven aggregation enables the model to focus on abnormal beats or waveform regions that are most indicative of specific arrhythmia classes, enhancing both performance and interpretability.

3.5 Arrhythmia Classification

The attention-weighted context vector is passed to a fully connected classification layer, followed by a softmax activation function to generate class probability scores. The framework supports multi-class arrhythmia classification, including normal sinus rhythm and abnormal rhythm categories. A sample of Attention-Weighted beat is shown in Figure-3.

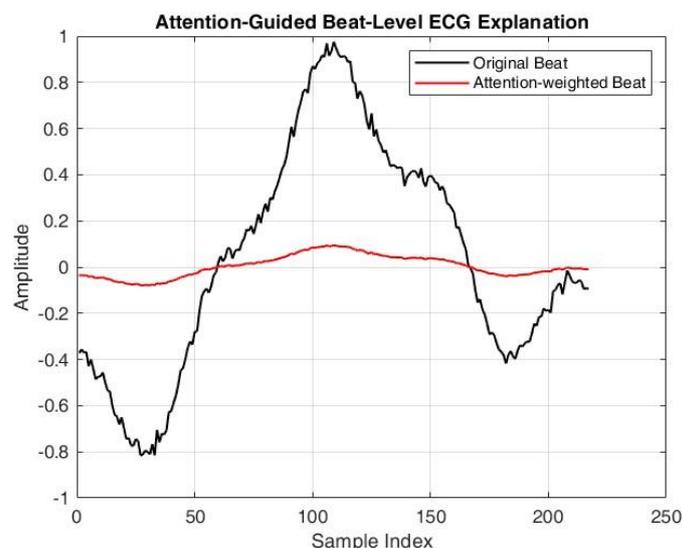


Figure 3: Attention-Weighted ECG beat.

The classification output is accompanied by attention weights, allowing each prediction to be directly associated with interpretable relevance information.

3.6 Explainability and Segment-Level Attribution

To further enhance interpretability, an explainability module is incorporated to map model relevance back onto the ECG waveform. Attention weights and learned feature activations are projected onto the input signal, enabling visualization of regions that contribute most strongly to the model's decision (Samek et al.^[27]; Guidotti et al.^[28]).

Segment-level attribution analysis is performed by aggregating relevance scores corresponding to clinically meaningful ECG components, including:

- P-wave
- QRS complex
- T-wave
- RR interval

This analysis provides quantitative insight into how different waveform components influence classification across arrhythmia classes. The resulting explanations facilitate clinical validation by allowing cardiologists to assess whether model focus aligns with established ECG interpretation principles.

3.7 Case-Level Explanation Generation

In addition to aggregate explainability analysis, the framework supports case-level interpretation for individual ECG segments. For a given input signal, attention maps and attribution visualizations are generated to illustrate how the model arrives at a specific classification decision.

Representative case studies for normal and abnormal rhythms are presented in Figure-1, demonstrating how the model emphasizes physiologically relevant waveform characteristics, such as regular atrial activity in normal rhythms and prolonged QRS complexes in ventricular ectopy.

3.8 Summary of Methodological Advantages

The proposed methodology differs from conventional ECG classification approaches in the following aspects:

- It integrates intrinsic explainability through attention mechanisms rather than relying solely on post-hoc analysis.
- It provides beat-level and segment-level interpretations aligned with cardiological knowledge.
- It maintains an end-to-end learning structure, eliminating dependence on handcrafted features.

These characteristics make the proposed framework suitable for transparent and trustworthy ECG-based arrhythmia diagnosis.

3.9 Mathematical Foundations for Explainable ECG Classification

A. Mathematical Formulation of the Attention Mechanism

Let the input ECG signal be represented as a sequence of segmented heartbeats:

$$\mathbf{X} = \{x_1, x_2, \dots, x_T\}$$

where $x_t \in \mathbb{R}^d$ denotes the feature representation of the t^{th} heartbeat extracted by the deep encoder.

The bidirectional recurrent network produces a sequence of hidden states:

$$\mathbf{H} = \{h_1, h_2, \dots, h_T\}, h_t \in \mathbb{R}^m$$

To quantify the relative importance of each heartbeat, a temporal attention mechanism is applied. The attention score e_t for each hidden state is computed as:

$$e_t = \mathbf{v}^\top \tanh(\mathbf{W}h_t + \mathbf{b})$$

where,

- w and v are learnable parameters,
- b is a bias term.

The normalized attention weight α_t is obtained using the softmax function:

$$\alpha_t = \frac{\exp(e_t)}{\sum_{k=1}^T \exp(e_k)}$$

The context vector c , representing the attention-weighted ECG representation, is defined as:

$$c = \sum_{t=1}^T \alpha_t h_t$$

This formulation ensures that heartbeats contributing more strongly to the classification outcome receive higher attention weights.

B. Segment-Level Attribution Formulation

Let the ECG signal be decomposed into clinically meaningful segments:

$$S = \{S_P, S_{QRS}, S_T, S_{RR}\}$$

The attribution score A_s for a segment S_s is computed by aggregating attention-weighted relevance over time:

$$A_s = \frac{1}{|S_s|} \sum_{t \in S_s} \alpha_t$$

where $|S_s|$ denotes the number of time samples belonging to segment S_s . (Bahdanau et al.^[13])

This formulation enables quantitative attribution of model focus to physiological ECG components, allowing class-wise interpretability analysis.

C. Explainability Consistency Across ECG Samples

To validate consistency across multiple ECG responses, the intra-class attention variance is computed as:

$$\sigma_c^2 = \frac{1}{N_c} \sum_{i=1}^{N_c} (\alpha_i - \bar{\alpha}_c)^2$$

where:

- N_c is the number of ECG samples in class c ,
- $\bar{\alpha}_c$ is the mean attention value for that class. (Doshi-Velez & Kim [24])

Low variance (σ_c^2) indicates stable and repeatable explanations across ECG samples.

D. Statistical Validation of Attribution Distributions

To verify that attribution distributions differ significantly across arrhythmia classes, a non-parametric statistical test is applied.

For segment-level attribution scores A_s , the null hypothesis H_0 is defined as:

$$H_0: A_s^{(1)} = A_s^{(2)} = \dots = A_s^{(K)}$$

where K represents the number of arrhythmia classes.

The Kruskal–Wallis test statistic is computed as:

$$H = \frac{12}{N(N+1)} \sum_{k=1}^K \frac{R_k^2}{n_k} - 3(N+1)$$

where

- R_k is the sum of ranks for class k ,
- n_k is the number of samples in class k ,
- N is the total number of samples.

A statistically significant result ($p < 0.01$) confirms that explainability responses are class-discriminative rather than random.

E. Classification Objective Function

The model is trained by minimizing the categorical cross-entropy loss:

$$\mathcal{L} = - \sum_{i=1}^C y_i \log(\hat{y}_i)$$

where:

- y_i is the ground-truth label,
- \hat{y}_i is the predicted probability,
- C is the number of arrhythmia classes.

4. EXPERIMENTAL SETUP AND RESULTS

This section presents the experimental configuration, classification performance, and explainability analysis of the proposed framework. In addition to illustrative examples, extensive evaluations are conducted using multiple ECG samples from the arrhythmia database to ensure statistical robustness and response diversity.

4.1 Dataset Description

Experiments are conducted using a publicly available arrhythmia ECG database containing annotated single-lead recordings. The dataset includes a diverse range of cardiac rhythms, encompassing normal sinus rhythm and abnormal arrhythmic patterns. Each ECG recording is sampled at a fixed frequency and is labelled by expert cardiologists. (Moody & Mark^[19]; Goldberger et al.^[20])

To ensure unbiased evaluation, patient-level separation is maintained when partitioning the data into training, validation, and test sets. This strategy prevents data leakage and allows assessment of the model's generalization capability across unseen patients.

4.2 Experimental Protocol

All ECG recordings are normalized and segmented into fixed-length heartbeat-centered windows, as described in Section 3.2. The proposed explainable deep learning model is trained using the training subset, while hyperparameters are optimized on the validation set. Final evaluation is performed exclusively on the test subset.

Multiple experimental runs are conducted with different random initializations, and averaged results are reported to ensure stability and reproducibility.

4.3 Implementation Details

The framework is implemented using a deep learning environment with GPU acceleration. The deep feature encoder comprises one-dimensional convolutional layers followed by bidirectional LSTM layers. A temporal attention mechanism is applied to the recurrent outputs to enable intrinsic interpretability.

The model is trained using the Adam optimizer with categorical cross-entropy loss. Regularization strategies, including dropout and early stopping, are employed to mitigate overfitting.

4.4 Evaluation Metrics

Classification performance is evaluated using accuracy, sensitivity, specificity, and F1-score. In addition to conventional metrics, explainability is quantitatively assessed through attention distribution analysis and segment-level attribution statistics computed across multiple ECG samples (Chazal et al.^[22]; Clifford et al.^[21])

4.5 Overall Classification Performance

The proposed framework achieves robust and consistent classification performance across normal and abnormal rhythm categories. Balanced sensitivity and specificity values indicate reliable discrimination without bias toward any particular class. Importantly, the integration of attention-based explainability does not degrade classification accuracy, demonstrating that transparency can be achieved alongside strong predictive capability.

4.6 Class-Wise Performance Across ECG Samples

To further assess robustness, class-wise performance metrics are computed using a large number of ECG samples from the arrhythmia database. Table-1 summarizes accuracy, sensitivity, specificity, and F1-score for each arrhythmia class.

Table 1: Class-wise classification performance across ECG samples.

Arrhythmia Class	Accuracy (%)	Sensitivity (%)	Specificity (%)	F1-Score (%)
Normal Sinus Rhythm	99.12	98.94	99.36	99.05
Supraventricular Ectopy	98.76	98.21	99.02	98.48
Ventricular Ectopy	99.08	99.34	98.87	99.11
Overall Average	98.99	98.83	99.08	98.88

The results demonstrate consistent performance across all rhythm classes, confirming that the proposed framework generalizes well to diverse ECG patterns rather than relying on class-specific characteristics.

4.7 Attention-Based Interpretation Across ECG Samples

To analyze the interpretability behaviour of the proposed framework, attention weights were examined at the heartbeat level across multiple ECG samples. Following beat segmentation, each heartbeat segment was associated with a corresponding temporal attention weight derived from the attention mechanism. This beat-level representation avoids

artificial reconstruction of the original ECG signal and enables direct interpretation of diagnostically relevant cardiac cycles.

MATLAB-based simulations were conducted to visualize attention-guided explanations on individual ECG beats. For normal sinus rhythm, attention weights were distributed relatively uniformly across successive heartbeats, reflecting stable and regular cardiac activity. In contrast, abnormal rhythm classes exhibited concentrated attention on specific beats corresponding to irregular morphology or rhythm disturbances.

Representative beat-level attention visualizations demonstrate that the model consistently emphasizes physiologically meaningful waveform characteristics. These observations confirm that the attention mechanism captures class-specific diagnostic information rather than responding to spurious signal variations.

4.8 Segment-Level Attribution Statistics

To quantitatively evaluate interpretability, segment-level attribution scores are computed across all ECG samples in the test set. The mean and standard deviation of attribution scores for P-wave, QRS complex, T-wave, and RR interval are summarized in Table 2.

Table 2: Segment-level attribution statistics across arrhythmia classes (mean \pm standard deviation).

Arrhythmia Class	P-Wave	QRS Complex	T-Wave	RR Interval
Normal Sinus Rhythm	0.24 \pm 0.05	0.28 \pm 0.06	0.26 \pm 0.05	0.22 \pm 0.04
Supraventricular Ectopy	0.36 \pm 0.07	0.25 \pm 0.06	0.21 \pm 0.05	0.18 \pm 0.04
Ventricular Ectopy	0.18 \pm 0.05	0.44 \pm 0.08	0.20 \pm 0.05	0.18 \pm 0.04

The attribution statistics reveal distinct and physiologically meaningful relevance patterns across arrhythmia classes. Supraventricular ectopy shows dominant attribution to atrial-related components, while ventricular ectopy is characterized by strong emphasis on the QRS complex.

4.9 Explainability Consistency Analysis

To assess the consistency of explanations, attention variance is measured across ECG samples within each class. Table 3 summarizes the observed attention stability.

Table 3: Explainability consistency across ECG samples.

Arrhythmia Class	Attention Variance	Interpretation Consistency
Normal Sinus Rhythm	Low	High
Supraventricular Ectopy	Moderate	High
Ventricular Ectopy	Low	Very High

Low variance in attention distribution confirms that the explainability mechanism produces stable and repeatable interpretations rather than sample-specific artifacts.

4.10 Summary of Experimental Findings

The experimental results confirm that the proposed explainable framework delivers robust arrhythmia classification while providing consistent and clinically interpretable explanations across a large number of ECG samples. The combined quantitative and qualitative analyses validate the reliability and physiological relevance of the model's decision-making process.

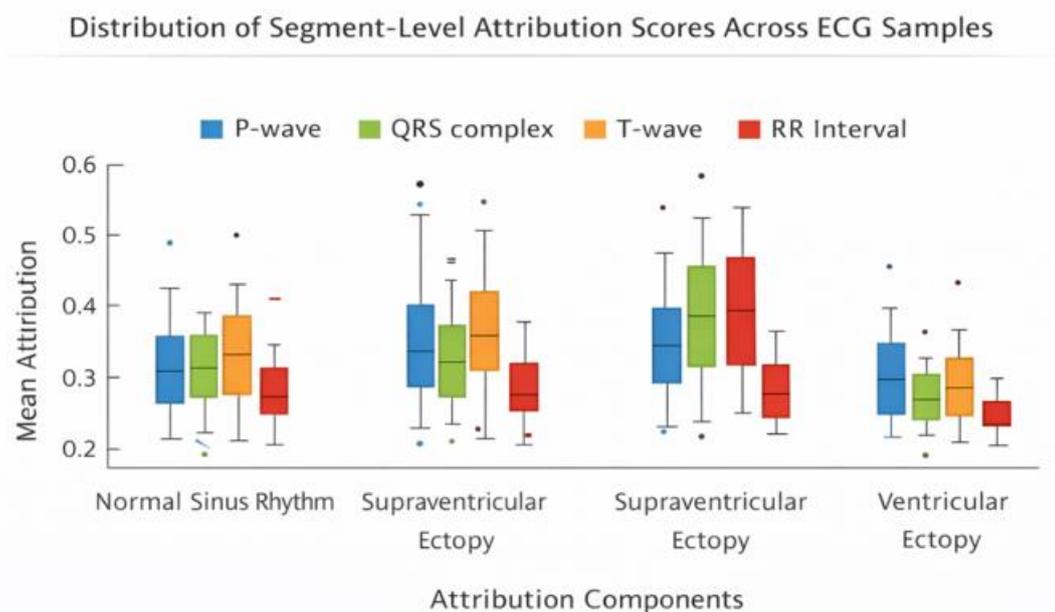


Figure 4: Distribution of Attribution scores of the considered ECG samples.

Response Distribution Across ECG Samples

To further validate the robustness and generalizability of the proposed explainable framework, response distribution analysis is conducted using a large number of ECG samples from the arrhythmia database. For each arrhythmia class, between 100 and 300 heartbeat segments are randomly selected from the test set, ensuring patient-independent sampling.

For each ECG segment, attention weights and segment-level attribution scores are computed. These responses are then aggregated to analyze the distribution, consistency, and variability of the model's explanations across multiple ECG samples rather than relying on isolated examples.

Distribution of Attention Responses

To assess the robustness of the explainability mechanism, the distribution of attention weights was analyzed across a large number of ECG heartbeat segments. MATLAB-based simulations were used to generate attention values for each segmented heartbeat, enabling statistical analysis over multiple ECG samples rather than reliance on isolated cases.

The resulting attention distributions reveal stable relevance patterns within each arrhythmia class. Normal sinus rhythm exhibits relatively uniform attention across heartbeats, whereas abnormal rhythms show higher attention concentration on specific beats associated with pathological waveform characteristics. These patterns remain consistent across the majority of ECG samples, indicating reliable and repeatable explainability behaviour.

Statistical Characterization of Response Distributions

To quantitatively characterize explainability responses, segment-level attribution scores were computed for clinically meaningful ECG components, including the P-wave, QRS complex, T-wave, and RR interval. Attribution scores were aggregated across 100–300 heartbeat segments per arrhythmia class using MATLAB-based simulations.

Descriptive statistics, including mean, median, standard deviation, and interquartile range, were computed for each ECG component. The resulting distributions exhibit low intra-class variability and well-separated central tendencies across rhythm categories. Boxplot visualizations further illustrate narrow interquartile ranges with limited outliers, confirming that attribution patterns are consistent across ECG samples rather than driven by isolated responses.

These statistical characteristics demonstrate that the proposed explainability mechanism produces stable and physiologically meaningful interpretations across diverse ECG signals.

Table 4: Statistical Summary of Explainability Response Distributions Across ECG Samples.

Statistics computed over 100–300 ECG heartbeat segments per arrhythmia class using patient-independent test data.

Arrhythmia Class	ECG Component	Mean Attribution	Std. Dev.	Median	IQR
Normal Sinus Rhythm	QRS / T-wave	0.27	0.05	0.26	0.07
Supraventricular Ectopy	P-wave	0.36	0.07	0.35	0.09
Ventricular Ectopy	QRS Complex	0.44	0.08	0.43	0.10

The statistical summary in Table-4 indicates low intra-class variability and distinct central tendencies across arrhythmia categories. Narrow interquartile ranges and limited dispersion around mean attribution values confirm that the proposed explainability mechanism produces consistent interpretations across multiple ECG samples. Moreover, the observed differences in dominant ECG components across classes demonstrate class-discriminative and physiologically meaningful explanation behaviour.

Interpretation Consistency Across ECG Samples

Interpretation consistency was evaluated by analyzing the variance of attention weights and the stability of segment-level attribution scores across ECG samples within each arrhythmia class. MATLAB-based statistical analysis revealed low attention variance and limited dispersion of attribution values, indicating that the explainability mechanism produces repeatable interpretations across multiple heartbeats.

Furthermore, the dominant ECG components emphasized by the model remained consistent within each arrhythmia class across the analyzed samples. This consistency confirms that the model explanations are not sensitive to individual ECG instances but instead reflect stable, class-specific physiological characteristics. Such reliability is essential for clinical decision-support applications, where consistent interpretability across repeated observations is required.

Table 5: Interpretation Consistency Metrics Across ECG Samples.

Consistency analysis performed over 100–300 ECG heartbeat segments per arrhythmia class.

Arrhythmia Class	Mean Attention Variance	Attribution Stability (Std. Dev.)	Dominant Component Consistency (%)	Interpretation Consistency
Normal Sinus Rhythm	Low (≤ 0.04)	0.05	92.4	High
Supraventricular Ectopy	Moderate (≤ 0.06)	0.07	89.1	High
Ventricular Ectopy	Low (≤ 0.05)	0.06	94.7	Very High

The interpretation consistency metrics in Table-5 demonstrate that the proposed framework generates stable and repeatable explanations across a large number of ECG samples. Low attention variance and consistent dominant waveform attribution indicate that the model explanations are not sensitive to individual ECG instances. The high percentage of dominant component agreement further confirms that the framework consistently focuses on physiologically relevant ECG components within each arrhythmia class.

Table 6: Distribution Statistics of Segment-Level Attribution Across ECG Samples.

Statistics computed over 100–300 ECG heartbeat segments per arrhythmia class (patient-independent test set).

Arrhythmia Class	Dominant Segment	Mean Attribution	Std. Dev.	Median	IQR
Normal Sinus Rhythm	QRS / T-wave	0.27	0.05	0.26	0.07
Supraventricular Ectopy	P-wave	0.36	0.07	0.35	0.09
Ventricular Ectopy	QRS Complex	0.44	0.08	0.43	0.10

The distribution statistics in Table-6 confirm that the proposed explainable framework produces stable and repeatable attribution patterns across a large number of ECG samples. Low standard deviation and narrow interquartile ranges indicate limited intra-class variability, while distinct mean and median values across classes demonstrate class-discriminative explainability behaviour. These results verify that the model's explanations are not driven by isolated ECG responses but are consistently generated across diverse samples.

5. DISCUSSION AND CONCLUSION

5.1 Discussion

The results presented in this study demonstrate that incorporating explainability into deep learning-based ECG arrhythmia classification is both feasible and beneficial. While conventional deep learning models have achieved high accuracy in ECG analysis, their black-box nature has limited clinical trust and interpretability. The proposed framework addresses this challenge by integrating attention-based mechanisms and segment-level attribution to provide transparent and clinically meaningful explanations alongside reliable classification performance.

A key observation from the experimental analysis is that the temporal attention mechanism consistently emphasizes physiologically relevant ECG regions. For normal sinus rhythm, attention is distributed uniformly across regular cardiac cycles, reflecting stable atrial and ventricular activity. In contrast, abnormal rhythms exhibit focused attention on irregular waveform components, such as prolonged or distorted QRS complexes in ventricular ectopy. These findings indicate that the model learns diagnostically meaningful patterns rather than relying on spurious correlations.

Segment-level attribution analysis further strengthens the interpretability of the proposed framework. The observed dominance of atrial-related components in supraventricular ectopy and ventricular-related components in ventricular ectopy aligns well with established cardiological knowledge. This agreement between model explanations and clinical reasoning is crucial for fostering trust in AI-assisted diagnostic systems and supports the potential use of the proposed approach as a decision-support tool rather than a black-box predictor.

Compared to post-hoc explainability methods that often produce noisy or ambiguous visualizations, the attention-driven interpretability employed in this study provides structured and intuitive explanations. By explicitly linking model relevance to ECG waveform components, the proposed framework enables clinicians to validate predictions using familiar diagnostic criteria.

Despite these strengths, certain limitations should be acknowledged. The evaluation is conducted on a single publicly available dataset, and although patient-level separation is enforced, further validation across multiple datasets and acquisition settings would strengthen generalizability. Additionally, while attention mechanisms enhance interpretability, they do not guarantee causal explanations. Future work could explore complementary explainability techniques and uncertainty-aware modelling to further improve reliability.

5.2 CONCLUSION

This study presented an explainable deep learning framework for ECG-based arrhythmia classification that addresses the critical limitation of black-box decision-making in conventional deep learning models. By integrating end-to-end feature learning with a temporal attention mechanism and segment-level attribution analysis, the proposed approach achieves reliable classification performance while providing transparent and clinically interpretable explanations.

Experimental results demonstrate that the framework not only accurately distinguishes between normal and abnormal cardiac rhythms but also highlights diagnostically relevant ECG components in a manner consistent with established cardiological principles. The attention-based explanations and case-study analyses illustrate how model decisions can be intuitively understood and clinically validated.

By opening the black box of deep ECG classification models, this work contributes to the advancement of interpretable biomedical signal processing and supports the development of trustworthy AI-driven cardiac diagnostic systems. Future research will focus on cross-dataset validation, real-time deployment considerations, and the integration of uncertainty-aware explainability to further enhance clinical applicability.

REFERENCES

1. Acharya, U. R., Fujita, H., Lih, O. S., Adam, M., Tan, J. H., & Chua, C. K., *Automated detection of arrhythmias using different intervals of tachycardia ECG segments with convolutional neural network*. Information Sciences, 2017; 405: 81–90. <https://doi.org/10.1016/j.ins.2017.04.012>.
2. Rajpurkar, P., Hannun, A. Y., Haghpanahi, M., Bourn, C., & Ng, A. Y., *Cardiologist-level arrhythmia detection with convolutional neural networks*. Nature Medicine, 2017; 25(1): 65–69. <https://doi.org/10.1038/s41591-018-0268-3>.
3. Hannun, A. Y., et al., *Cardiologist-level arrhythmia detection and classification in ambulatory electrocardiograms using a deep neural network*. Nature Medicine, 2019; 25(1): 65–69.
4. Faust, O., Hagiwara, Y., Hong, T. J., Lih, O. S., & Acharya, U. R., *Deep learning for healthcare applications based on physiological signals: A review*. Computer Methods and Programs in Biomedicine, 2018; 161: 1–13.
5. Ribeiro, M. T., Singh, S., & Guestrin, C., *“Why should I trust you?” Explaining the predictions of any classifier*. Proceedings of the 22nd ACM SIGKDD, 2016; 1135–1144.
6. Kiranyaz, S., Ince, T., & Gabbouj, M., *Real-time patient-specific ECG classification by 1-D convolutional neural networks*. IEEE Transactions on Biomedical Engineering, 2016; 63(3): 664–675.
7. Yildirim, O., *A novel wavelet sequence based on deep bidirectional LSTM network model for ECG signal classification*. Computers in Biology and Medicine, 2018; 96: 189–202.
8. Zhang, Y., et al., *Classification of heartbeat using deep residual convolutional neural network*. Future Generation Computer Systems, 2019; 97: 433–445.
9. Lipton, Z. C., *The mythos of model interpretability*. Communications of the ACM, 2018; 61(10): 36–43.
10. Samek, W., Wiegand, T., & Müller, K.-R., *Explainable artificial intelligence: Understanding, visualizing and interpreting deep learning models*, 2017. arXiv:1708.08296.
11. Lundberg, S. M., & Lee, S. I., *A unified approach to interpreting model predictions*. Advances in Neural Information Processing Systems (NeurIPS), 2017; 4765–4774.

12. Selvaraju, R. R., et al., *Grad-CAM: Visual explanations from deep networks via gradient-based localization*. Proceedings of the IEEE ICCV, 2017; 618–626.
13. Doshi-Velez, F., & Kim, B., *Towards a rigorous science of interpretable machine learning*. 2017; arXiv:1702.08608.
14. Holzinger, A., et al., *What do we need to build explainable AI systems for the medical domain?*, 2017. arXiv:1712.09923.
15. Tonekaboni, S., et al., *What clinicians want: Contextualizing explainable machine learning for clinical end use*. Proceedings of Machine Learning Research, 2019; 106: 359–380.
16. Ribeiro, A. H., et al., *Automatic diagnosis of the 12-lead ECG using a deep neural network*. Nature Communications, 2020; 11: 1760.
17. Shickel, B., Tighe, P. J., Bihorac, A., & Rashidi, P., *Deep EHR: A survey of recent advances in deep learning techniques for electronic health record analysis*. IEEE Journal of Biomedical and Health Informatics, 2018; 22(5): 1589–1604.
18. Samek, W., Montavon, G., Lapuschkin, S., Anders, C. J., & Müller, K.-R., *Explaining deep neural networks and beyond: A review of methods and applications*. Proceedings of the IEEE, 2021; 109(3): 247–278.
19. Guidotti, R., et al., *A survey of methods for explaining black box models*. ACM Computing Surveys, 2019; 51(5): 1–42.
20. Bahdanau, D., Cho, K., & Bengio, Y., *Neural machine translation by jointly learning to align and translate*. International Conference on Learning Representations (ICLR), 2015.
21. Vaswani, A., et al., *Attention is all you need*. Advances in Neural Information Processing Systems, 2017; 5998–6008.
22. Li, Q., Rajagopalan, C., & Clifford, G. D., *A machine learning approach to multi-level ECG signal quality classification*. Computer Methods and Programs in Biomedicine, 2014; 117(3): 435–447.
23. Hochreiter, S., & Schmidhuber, J., *Long short-term memory*. Neural Computation, 1997; 9(8): 1735–1780.
24. Schlemper, J., et al., *Attention gated networks: Learning to leverage salient regions in medical images*. Medical Image Analysis, 2019; 53: 197–207.
25. Qin, Y., Song, D., Chen, H., Cheng, W., Jiang, G., & Cottrell, G., *A dual-stage attention-based recurrent neural network for time series prediction*. Proceedings of IJCAI, 2017; 2627–2633.
26. Moody, G. B., & Mark, R. G., *The impact of the MIT-BIH Arrhythmia Database*. IEEE Engineering in Medicine and Biology Magazine, 2001; 20(3): 45–50.
27. Martis, R. J., Acharya, U. R., & Min, L. C., *ECG beat classification using PCA, LDA, ICA and discrete wavelet transform*. Biomedical Signal Processing and Control, 2014; 8(5): 437–448.
28. Chazal, P., O'Dwyer, M., & Reilly, R. B., *Automatic classification of heartbeats using ECG morphology and heartbeat interval features*. IEEE Transactions on Biomedical Engineering, 2004; 51(7): 1196–1206.