

# MED-Net: Multi-Scale Enhanced Dual Temporal Convolution Attention Network

Ehsaas Nahata, V.S.S. Madhuri, Kandibanda Sathwika, Kanishka Gupta, Magam Honeysha, Neha Thapasvi Kodithala, Naga Hari Babu KV, Shilpa Chaudhary

Keshav Memorial Engineering College  
<https://github.com/EhsaasN/MED-Net>

## Abstract

Electrocardiographic (ECG) signals are used to evaluate heart activity and to identify disease-related anomalies. Reliable support systems are useful for analysing ECG signals, for instance, in long-term data acquisition and evaluation (e.g., 24-hour Holter recording) or to support physicians in reading ECGs. This paper proposes an ECG anomaly detection and diagnosis model, *MED-Net*, based on DTAAD (Dual TCN Attention for Anomaly Detection). The model is an integrated design in which an autoregressive model (AR) combines with an autoencoder (AE) structure. Scaling methods and feedback mechanisms are introduced to improve prediction accuracy and expand correlation differences. The Dual TCN-Attention Network (DTA) uses two Transformer attentions in our baseline experiment, belonging to an ultra-lightweight model. Upon extensive experiments on two public datasets MED-Net exceeds the majority of currently advanced baseline methods in both detection and diagnostic performance and achieves better results than DTAAD. Specifically, MED-Net is targeted towards analysing ECG anomalies from multivariate time series data by improved F1 scores of about 1.5% over DTAAD and 5% over CAE with reduced training time of about 99% compared to CAE and 50% over DTAAD.

## 1 Introduction

Anomaly detection refers to the detection of these anomalies from the expected data that deviate substantially from the rest of the distribution, and they are then labeled as anomalies. Multivariate time series are universal in the real world, and anomaly detection in this aspect has always been a research hotspot [3]. Anomaly detection techniques are now widely used in various application domains such as finance, epidemic, industry, cyber security, credit fraud, power industry, unmanned vehicles, healthcare, distributed computing, urban IoT, robotics, traffic, and sensor network monitoring, satellite aviation, etc., and are particularly prominent in data mining, probabilistic statistics [4], machine learning [5], and computer vision [6]. In recent years, anomaly detection methods suitable for Deep Learning (DL) have been drastically improving. DL has shown great ability in learning complex data such as temporal data, streaming media data, graphical spatial data, and high-dimensional data, raising the limitations of various learning tasks. DL for anomaly detection is referred to as deep anomaly detection [7], which is performed by deep neural networks to learn features to determine or get anomaly scores to achieve the purpose of anomaly detection, and it is usually performed in the context of prediction. ECG signals are commonly used to assess heart activity and health status, particularly in evaluating the risk of heart attack. ECGs are sets of signals that measure the electrical activity of the heart. The duration of an ECG signal can range from 10 seconds, as used in standard clinical evaluations, to several days in the case of continuous cardiac monitoring (Holter device). Heart activity is typically measured through multiple channels, which can range from 1 to 12. Some ECG devices are also capable of measuring three additional signals to generate vectorcardiograms (VCG), which trace both the magnitude and direction of the electrical forces associated with heart activity. Wearable

devices now allow real-time measurement and monitoring of heart-related signals even during physical activities.

### 1.1 Challenges

In contemporary data science research, detecting anomalies in large-scale databases has become increasingly challenging due to the growing affluence and rapid reasoning for quick recovery and optimal service quality. Time series databases, often comprising data from various engineering devices interacting with the environment, humans, or other systems, frequently exhibit uncertainty and temporal trends. Although understanding ECG signals is an important concern for physicians, there remains a need for reliable tools to interpret and analyze these ECG signals automatically with minimal human involvement. For example, cardiac arrhythmias are characterized by irregular heartbeats, and detecting them involves identifying deviations from a regular rhythm. By reading electrocardiograms, doctors and physicians can detect conditions such as atrial fibrillation and frequency anomalies. Similarly, arrhythmias can be identified by experts through a visual examination of ECGs. In long-term ECG signal acquisition, computer aided analysis tools are essential for helping physicians detect anomalies. Consequently, automated anomaly detection in ECG signals has emerged as a primary research focus. In this direction, AI-based tools can assist clinicians by integrating the insights gained from studies and analyzes. The increase in integration of IoT technology has led to the generation of vast amounts of data from numerous sensors and devices, which pose difficulties not only in volume but also in diversity, complicating accurate and reliable data analysis. Moreover, modern applications demand high responsiveness. Hence, accurately identifying anomalies and distinguishing observations against the temporal trend is crucial. However, the scarcity of labeled data and the various types of outliers add to the challenges. Tra-

ditional supervised learning models face limitations in direct application, while unsupervised learning may prove effective in other domains [8]. Furthermore, beyond anomaly detection, finding the root cause of anomalies is essential, requiring multi-category predictions to determine both the presence of an anomaly and its cause. In summary, addressing multivariate time series anomaly detection in modern applications involves overcoming challenges related to data complexity, inference speed, and anomaly source tracking.

## 2 Literature Survey

Anomaly detection in multivariate time series data, particularly for ECG signals, is a critical area of research characterized by challenges such as high dimensional complexity, lack of labeled data, and the need for rapid inference. Traditional unsupervised methods relying on statistical distributions, distance metrics, or density estimation often fail to capture the intrinsic temporal correlations present in complex physiological signals such as ECG. Although Deep Learning (DL) models have advanced the state of the art, existing architectures from recurrent networks to generative models suffer from limitations concerning computational efficiency, global context capture, and training stability.

### 2.1 Limitations of Recurrent and Complex Hybrid Architectures

A major portion of existing literature employs Recurrent Neural Networks (RNNs) [9] or Long Short-Term Memory (LSTM) networks [10] to model temporal dependencies. For example, *Cloud-based healthcare framework for real-time anomaly detection and classification of 1-D ECG signals* uses Wavelet Scattering combined with an LSTM Autoencoder [11]. While this reduces signal dimensionality, it introduces sequential bottlenecks inherent to LSTMs and depends heavily on cloud infrastructure. Similarly, the *Hybrid Deep Learning Approach Using BiGRU-BiLSTM and Multilayered Dilated CNN* [27] integrates dual-branch recurrent architectures with dilated convolutional networks and GAN-based augmentation. Although accurate, such hybrid systems incur substantial computational overhead, slow training, and increased risk of overfitting due to architectural complexity. Sequential-processing bottlenecks are also evident in *OmniAnomaly* [13], which employs a GRU-based encoder-decoder. Due to step-by-step updates, GRUs become slow for long ECG sequences and struggle with abrupt noise fluctuations. Furthermore, *Diagnosis of atrial fibrillation based on AI-detected anomalies of ECG segments* [26] relies on rigid segmentation (PreQ, QRS, PostS), which limits adaptability to varying signal lengths. Comparative analyses in *Deep learning for ECG classification: A comparative study of 1D and 2D representations and multimodal fusion approaches* further show that 2D representations often introduce unnecessary computational cost without consistent performance gains.

### 2.2 Limitations of Static, Local, and Manual Approaches

To overcome the computational inefficiency of RNNs, researchers adopted Convolutional Neural Networks (CNNs) and Transformer-based models [14]. However, these approaches also face their respective limitations. The *Convolutional Autoencoder Framework for ECG Signal Analysis* [2] uses 2D convolutions but is constrained by fixed lead configurations and struggles with noise. Likewise, the *Simple 1D CNN with Leaky-ReLU for ECG Classification* [28] offers portability but lacks the ability to capture long-range dependencies essential for detecting arrhythmias. The *Automatic QRS Complex Detection using Two-Level CNN* [29] focuses on small temporal windows and fails to capture global rhythm patterns. Meanwhile, Transformer-based approaches such as *TranAD* [15] leverage attention mechanisms to accelerate inference but often emphasize short local windows, missing slower, long-range patterns. *Transformer-based Multivariate Time Series Anomaly Localization* [30] requires repeated masking and reconstruction, significantly increasing computational cost. Other approaches such as *Global ECG Classification by Self-Operational Neural Networks with Feature Injection* [31] and *Fuzz-ClustNet* [32] depend on tailored decompositions or fuzzy clustering, making them inflexible and vulnerable to cluster drift in noisy environments. Visual analytics tools such as *ECGLens* [16] rely on manual inspection, limiting scalability for real-time diagnosis.

### 2.3 Instability of Generative and Adversarial Models

Generative Adversarial Networks (GANs) have been explored for modeling normal ECG distributions but suffer from training instability. For instance, the *Abnormal ECG Detection Based on an Adversarial Autoencoder* [17] combines Temporal Convolutional Networks (TCNs) [18] with adversarial learning, requiring sensitive balancing between generator and discriminator losses. The *MadeGAN* framework [19] introduces a two-level Memory-Augmented Deep Autoencoder with GANs, but incurs substantial computational and memory overhead due to its staged adversarial training and memory modules.

### 2.4 The MED-Net Model

The proposed MED-Net framework addresses these limitations described above through a streamlined, non-adversarial architecture emphasizing efficiency, temporal adaptability, and stability. Building upon Dual Temporal Convolutional Networks (TCNs) [18] and DTAAD [1], our enhanced model integrates several key innovations:

- **Overcoming recurrent bottlenecks:** Unlike LSTM/GRU-based models such as *OmniAnomaly* [13] or LSTM Autoencoders used in cloud-based systems [11], MED-Net employs Dual TCNs with Self-Attentions along with Efficient Multi-Scale feature extraction. This

enables parallel processing of multi-resolution temporal features, eliminating sequential latency and improving scalability for long ECG recordings.

- **Capturing global context:** To mitigate the short window limitations of Transformer based approaches like *TranAD* [15], the model incorporates a Lightweight ECG Attention module and Enhanced Attention layers. These modules efficiently capture global temporal relationships and long-range rhythm patterns without incurring the computational burden of standard Transformer encoders [14].
- **Stability and efficiency:** Unlike adversarial systems such as GAN based ECG detectors [17, 19], MED-Net uses a Fixed Fusion mechanism (weights 0.3, 0.4, 0.3) that ensures consistent feature balancing. Removing dropout in attention heads further improves inference speed and stability, making the model suitable for edge and real-time deployment scenarios.
- **Adaptability to data:** In contrast to rigid segmentation based [26] or fuzzy clustering approaches, MED-Net integrates Dynamic Decoders and Smart Label Handling to adapt flexibly to changes in input dimensions, variable length windows, and different preprocessing strategies (e.g., overlapping MBA windows or non-overlapping clinical ECG segments).

In summary, the MED-Net model offers a robust, lightweight, and automated alternative that integrates local and global temporal features effectively and efficiently. It addresses key limitations in prior literature, achieving superior performance in terms of speed, stability, and diagnostic precision for ECG anomaly detection.

## 3 Proposed Work

We propose MED-Net which has a few enhancements over the existing DTAAD architecture [1] to improve feature extraction and anomaly detection in ECG data.

### 3.1 Problem Definition

Consider an ECG signal represented as a univariate/multivariate time series data

$$T = \{x_1, x_2, \dots, x_T\} \quad (1)$$

where each value  $x_t \in \mathbb{R}$  denotes the ECG amplitude recorded at time  $t$ . Since the ECG signal is fully observed, its joint probability distribution can be factorized autoregressively as

$$p(x_{1:T}) = \prod_{t=1}^T p(x_t \mid x_{1:t-1}) \quad (2)$$

### 3.2 Data Preprocessing

In real-world industrial systems, multivariate time series often contain missing entries, noise, and irregular sampling due to sensor drift or external interference. To ensure robust model training and enhance generalization, the raw data are first standardized and then segmented into uniform temporal windows. During preprocessing, non-numeric entries and missing values were handled by coercing invalid data to NaN, dropping empty rows or columns, and filling remaining missing values with zeros. Only sensor or numerical columns were retained, excluding metadata such as timestamps or IDs, so that the clean data form a raw matrix

$$X \in \mathbb{R}^{T \times D}$$

where  $T$  is the total number of time steps and  $D$  is the number of variables (sensors). First, missing values and obvious corrupt points are handled (e.g., via interpolation or deletion). After this, each dimension is normalized feature-wise using statistics computed only from the training portion. For the  $d$ -th variable, compute the mean  $\mu_d$  and standard deviation  $\sigma_d$ , and obtain the standardized series  $\tilde{x}_{t,d}$  via z-score normalization:

$$\tilde{x}_{t,d} = \frac{x_{t,d} - \mu_d}{\sigma_d}, \quad t = 1, \dots, T, \quad d = 1, \dots, D. \quad (3)$$

(Equivalently, a min-max scaling can be written as)

$$\tilde{x}_{t,d} = \frac{x_{t,d} - \min_t x_{t,d}}{\max_t x_{t,d} - \min_t x_{t,d}}, \quad (4)$$

but in our implementation we consistently use the z-score form. This normalization ensures all variables are on a comparable scale and stabilizes the training of deep models.

#### 3.2.1 Sliding Window Formulation

After normalization, the continuous sequence  $\tilde{X}$  is segmented into windows to form the final model inputs. For a fixed window length  $L$ , each sample  $S_i \in \mathbb{R}^{L \times D}$  is constructed as

$$S_i = [\tilde{x}_{i,:}, \tilde{x}_{i+1,:}, \dots, \tilde{x}_{i+L-1,:}] \quad i = 1, \dots, T - L + 1 \quad (5)$$

To model temporal dependencies, a fixed-length sliding window of size  $W$  is used. At each timestep  $t$ , the model receives the window

$$[x_{t-W+1}, x_{t-W+2}, \dots, x_t] \quad (6)$$

and predicts the next value  $x_{t+1}$ . This forecasting mechanism serves as the basis for both anomaly detection and anomaly diagnosis. where  $\tilde{x}_{t,:}$  denotes the  $D$ -dimensional vector at time step  $t$ . Given the original point-wise anomaly label sequence  $\{\ell_t\}_{t=1}^T$  with  $\ell_t \in \{0, 1\}$ , a window-level label  $y_i$  is assigned by aggregating the labels inside the window using a max operator:

$$y_i = \max_{t \in [i, i+L-1]} \ell_t$$

so that  $y_i = 1$  if any time step in the window is anomalous, and  $y_i = 0$  otherwise. The resulting set  $\{(S_i, y_i)\}$  is then split into training, validation, and test subsets according to the experimental protocol, providing the standardized and windowed input required by the dual TCN-attention network [1].

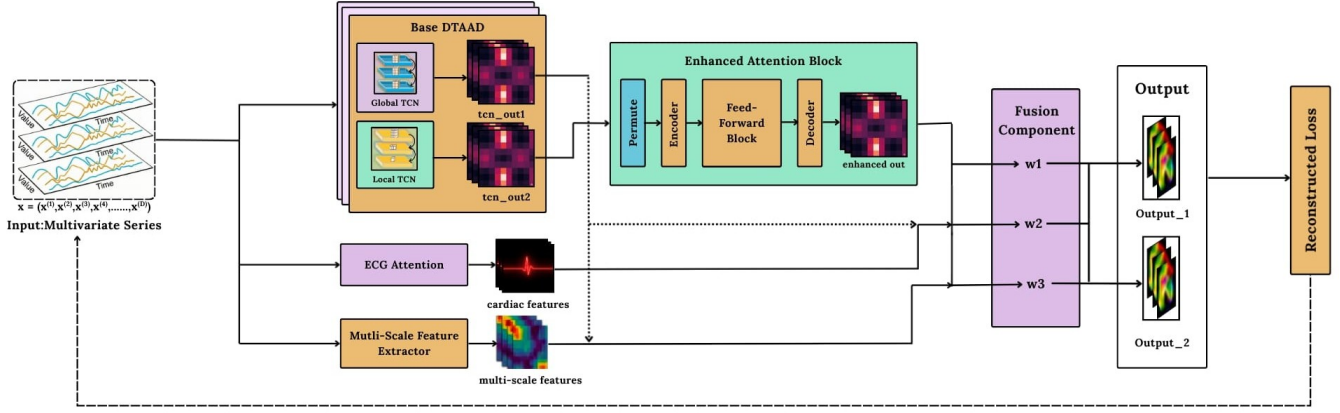


Fig. 1: MED-Net architecture illustrating the complete forward pass. The Base DTAAD extracts global and local temporal patterns ( $tcn\_out1$ ,  $tcn\_out2$ ). ECG Attention and Multi-scale blocks derive cardiac and morphological features directly from the input. Enhanced Attention refines  $tcn\_out2$  through optional expansion, attention, and projection stages. All components are fused using fixed weights ( $w_1$ ,  $w_2$ ,  $w_3$ ) to generate dual outputs for anomaly detection.

### 3.3 Architecture

**DTAAD:** The overall model can be approximated as an autoregressive model (AR). The first half of the model consists of dual TCN, and the codec layer model (AE) is built based on the Transformer in the back [1]. DTAAD utilizes two temporal convolutional structures, causal convolution, and dilated convolution as local TCN and global TCN, which will generate the same number of temporal prediction out puts as the input, and then into a Transformer encoding layer for capturing the dependencies between multiple series. Then, the global attention and local attention outputs from the coding layer flow into a specific decoding layer according to the residual connection [21]. The final prediction results of the local attention are fed back to the global TCN and the original input is overlaid by the replication operation. Finally, the reconstruction errors obtained based on the two sets of prediction values are combined in a certain proportion to obtain the training loss.

### 3.4 Enhancements

These enhancements are done to improve the model performance on ECG data and to improve feature extraction.

#### 3.4.1 Multi-Scale Temporal Feature Extractor

It is a feature extractor that captures patterns at different time scales simultaneously. It consists of two kernels to detect short patterns ( $k=3$ ) and longer patterns ( $k=5$ ) which is then combined together. This is require to capture components of different durations, QRS complex (80-120ms) and P-wave (80ms) are captured by the short kernel whereas the T-wave (160ms) is captured by the longer kernel size. A single kernel fails to capture both sharp spikes as well as broad waves which the multi-scale extractor accomplishes.

#### 3.4.2 Lightweight ECG Attention

This is the attention mechanism that focuses on cardiac rhythm patterns consisting of two heads and no dropout to improve speed. It learns to focus on important cardiac patterns and ignores noise and irrelevant signals by using self attention to find temporal dependencies.

#### 3.4.3 Enhanced Attention

Enhanced Attention is an additional refinement layer that takes the already-processed output from base DTAAD's second branch and applies another round of multi-head attention to further refine the temporal patterns. It further refines the output to check for any missed features and amplifies important cardiac events, suppresses noise and captures long range dependencies.

#### 3.4.4 Fusion Component

This component combines the outputs from multiple sources into a single enhanced representation. It combines the output from base DTAAD [1], Lightweight Ecg Attention, Enhanced Attention and the multi-scale feature extractor.

#### 3.4.5 Increased Window Size

The window size of the local tcn is increased to 128 from 10. While this adds a computational overhead but it is useful in capturing longer patterns in the dataset.

#### 3.4.6 Added Learning Rate Decay

The original DTAAD framework trained the model using a fixed learning rate throughout all epochs, resulting in parameter updates of constant magnitude. While effective for

simple temporal patterns, a constant learning rate can slow convergence during the later stages of training and may cause oscillations around local minima, especially in deep temporal architectures that combine dual TCN branches with multiple attention modules. To address these limitations, the MED-Net introduces a learning rate decay schedule, allowing the step size to decrease gradually as training advances. This reduces instability, promotes smoother convergence, and enables the model to perform finer adjustments when approaching an optimum. The decayed learning rate follows:

$$\eta_t = \eta_0 \cdot \gamma^{\lfloor \frac{t}{T_d} \rfloor} \quad (7)$$

Where  $\eta_0$  is the initial learning rate,  $\gamma \in (0, 1)$  is the decay factor, and  $T_d$  is the decay interval measured in epochs. Early in training, larger values of  $\eta_t$  enable rapid exploration of the parameter space, while later reductions ensure precise refinement of temporal filters, attention weights, and multi-scale feature representations. Empirically, integrating learning rate decay into the training loop of MED-Net results in faster convergence, improved numerical stability, and enhanced reconstruction accuracy, particularly for ECG signals containing substantial noise and subject-specific variability.

### 3.5 Model Framework

The MED-Net architecture is designed for efficient anomaly detection in ECG time-series. The model processes input segments through a Dual-TCN backbone, extracts morphology-aware multi-scale features, applies lightweight and refined temporal attention, and finally fuses all information streams before reconstruction-based anomaly scoring.

#### 3.5.1 Input Representation and Window Construction

Let the raw ECG signal be a time-series:

$$X = \{x_1, x_2, \dots, x_T\}, \quad x_t \in \mathbb{R}^C$$

This notation defines the recorded ECG as a sequence of  $T$  time steps where each sample  $x_t$  may be scalar (single-lead,  $C = 1$ ) or a vector (multivariate leads,  $C > 1$ ). Treating the signal as a time-indexed sequence clarifies that all later operations — convolution, attention, and reconstruction — are temporally causal and operate over contiguous windows of these samples. The sequence is divided into windows:

$$W_t = [x_{t-L+1}, \dots, x_t] \in \mathbb{R}^{L \times C}$$

We use fixed-length windows of size  $L$  to convert the streaming ECG into examples suitable for batch training. Windows capture local temporal context; overlapping windows allow fine-grained localization while non-overlapping windows reduce computation. This sliding-window representation is the basic input unit for the Dual-TCN and attention blocks. Each window is rearranged into channel-first format:

$$W_t \in \mathbb{R}^{C \times L}$$

Reordering to channel-first (channels  $\times$  time) matches the input convention of most deep-learning libraries (PyTorch, etc.) and simplifies per-channel convolutions and attention. It also makes broadcasting and tensor arithmetic straightforward when fusing features from multiple modules.

#### 3.5.2 Dual-TCN Backbone

**Local TCN:** Models short-range, sharp ECG morphology [18]. For convolution kernel size  $k$  and dilation  $d$ :

$$H_i^{(loc)}(t) = \phi \left( \sum_{j=0}^{k-1} W_t(:, t - jd) K_i^{(loc)}(j) + b_i^{(loc)} \right) \quad (8)$$

This equation describes a causal dilated convolution producing the  $i$ -th output channel of the local TCN. The nonlinearity  $\phi(\cdot)$  (e.g., ReLU) [22] introduces nonlinearity after a weighted sum over local time offsets. By choosing small kernel sizes and modest dilations, the local TCN focuses on sharp, high-frequency features such as QRS spikes and abrupt transitions in the ECG. **Global TCN:** Models long-range temporal patterns using dilations  $d = \{1, 2, 4\}$ :

$$H_i^{(glob)}(t) = \phi \left( \sum_{j=0}^{k-1} W_t(:, t - jd_i) K_i^{(glob)}(j) + b_i^{(glob)} \right) \quad (9)$$

The global TCN uses larger dilation factors to exponentially increase receptive field without a proportional increase in parameters. This enables the network to capture long-duration waves and slow trends (e.g., baseline wander, long rhythm patterns) while retaining causality. In practice multiple stacked dilations are used so each output combines information from both recent and more distant time steps. **Dual-TCN Output:**

$$H^{(dual)} = [H^{(loc)}; H^{(glob)}] \in \mathbb{R}^{128 \times L} \quad (10)$$

Concatenating local and global outputs yields a richer temporal representation that encodes both short and long temporal scales. This combined tensor is then available to downstream attention and fusion modules; using concatenation (rather than addition) preserves the independent characteristics of each branch until the fusion stage.

#### 3.5.3 Multi-Scale Temporal Feature Extractor

To capture waveform components of varying timesteps, the model applies convolutions at different scales:

$$F_3(t) = \text{Conv}_{k=3}(W_t)(t)$$

$$F_5(t) = \text{Conv}_{k=5}(W_t)(t)$$

Applying multiple kernel sizes extracts features sensitive to distinct morphological widths: short kernels respond to sharp spikes and onset slopes while longer kernels capture broader waves and plateaus. This explicit multi-scale decomposition is crucial in ECG analysis where clinically relevant components (P, QRS, T) have different durations. Their fused output is:

$$F^{(ms)}(t) = W_{\text{fuse}}([F_3(t); F_5(t)]), \quad (11)$$

Here  $W_{\text{fuse}}$  denotes a small learnable projection (e.g.,  $1 \times 1$  conv or linear layer) that compresses and mixes multi-scale feature maps. Fusing the scale-specific features reduces dimensionality and produces a coherent multi-scale descriptor that the decoder and fusion module can use directly. The final shape is

$$F^{(ms)} \in \mathbb{R}^{C \times L}.$$

This ensures the multi-scale features align with the channel and time dimensions expected by the fusion step and decoder, enabling straightforward element-wise combination and reconstruction.

### 3.5.4 Lightweight ECG Attention

A two-head self-attention module highlights diagnostically relevant segments. Given  $Z \in \mathbb{R}^{C \times L}$ :

$$Q = ZW_Q, \quad K = ZW_K, \quad V = ZW_V \quad (12)$$

The query/key/value projection maps input features to an attention subspace. Using two heads provides modest multi-perspective modeling without the overhead of full Transformer stacks [14]. Practically,  $W_Q, W_K, W_V$  are small linear layers applied per channel/time.

$$\text{Attn}(Z) = \text{softmax}\left(\frac{QK^\top}{\sqrt{d}}\right) V \quad (13)$$

This scaled dot-product attention computes a weighted sum of value vectors where weights reflect similarity between queries and keys. The softmax normalizes attention weights across time positions so the module focuses on a few diagnostically important timestamps (e.g., peaks, onsets), suppressing noise and irrelevant signals portions. The output is reshaped to:

$$A^{(light)} \in \mathbb{R}^{C \times L}$$

The attention output retains the same channel/time layout, allowing it to be added or concatenated with TCN features. Because dropout in attention heads is disabled in MED-Net, this module is deterministic and fast during inference.

### 3.5.5 Enhanced Attention Refinement

A second attention module further refines long-range temporal dependencies. Applied to the global TCN output:

$$A^{(enh)} = \text{Attn}(H^{(glob)})$$

By running attention on the global-TCN representation, the enhanced module captures relationships that the lighter attention (run on raw input) may miss. This two-stage attention—input-level then feature-level—helps amplify weak but globally-consistent patterns that are crucial for diagnosing longer arrhythmic events, which maintains the same shape:

$$A^{(enh)} \in \mathbb{R}^{C \times L}$$

Maintaining the same tensor shape simplifies fusion and allows element-wise operations. The refined attention is projected back into the same dimensionality as other streams so it can be safely combined.

### 3.5.6 Fixed-Weight Feature Fusion

Three feature streams are fused using fixed scalar weights:

$$\tilde{H}_1 = H^{(loc)} + \alpha_1 A^{(light)} \quad (14)$$

This first fused representation augments local TCN features with attention-weighted input cues. The scalar  $\alpha_1$  controls the relative contribution; fixing it (0.3) stabilizes training by preventing the fusion weights from growing unbounded.

$$\tilde{H}_2 = H^{(glob)} + \alpha_2 F^{(ms)} + \alpha_3 A^{(enh)} \quad (15)$$

The second fused stream combines global TCN outputs with multi-scale morphological cues and enhanced attention. This additive fusion integrates complementary information: long-range context, multi-scale morphology, and refined attention-based emphasis. The fusion coefficients are:

$$(\alpha_1, \alpha_2, \alpha_3) = (0.3, 0.4, 0.3)$$

Selecting these constants empirically balances the three streams: the multi-scale features are weighted slightly higher (0.4) because morphology is particularly informative for ECG anomalies, while attention streams provide supporting evidence.

### 3.5.7 Reconstruction and Objective

A decoder reconstructs the original window:

$$\hat{W}_t = \text{Decoder}(\tilde{H}_1, \tilde{H}_2)$$

The decoder is typically a small stack of transposed convolutions or upsampling convolutions that maps fused representations back to the original channel-time grid. Its goal is to reproduce normal windows accurately; reconstruction failure indicates anomalous input patterns.

$$\hat{W}_t \in \mathbb{R}^{C \times L}$$

The reconstructed tensor has the same shape as the input window to allow per-sample, per-channel residual computation. Ensuring identical shape simplifies loss computation and thresholding logic used in anomaly scoring.

The training loss is:

$$\mathcal{L} = \frac{1}{L} \sum_{i=1}^L \left\| W_t(:, i) - \hat{W}_t(:, i) \right\|_2^2 \quad (16)$$

This mean-squared reconstruction objective penalizes point-wise deviations and encourages the network to model normal morphology faithfully. Averaging over the window length yields numerically stable gradients and a scale-invariant loss with respect to window size.

### 3.5.8 Anomaly Scoring

Reconstruction error indicates deviation from normal ECG morphology:

$$s_t = \|W_t - \hat{W}_t\|_2^2$$

This scalar score summarizes the reconstruction discrepancy for an entire window. In practice we aggregate or smooth window scores (e.g., by taking a moving average or maximum across overlapping windows) and apply a validation-derived threshold to produce binary anomaly flags. The score is interpretable and easily calibrated for clinical or industrial use.

## 4 Training Procedure

The MED-Net model is trained through self-supervised reconstruction of normal ECG windows. The primary goal is to learn the morphology and temporal structure of healthy cardiac patterns, such that abnormal behaviors naturally incur larger reconstruction errors. The training pipeline consists of preprocessing, window generation, multi-branch forward propagation, loss computation, and gradient-based optimization with learning-rate decay.

### 4.1 Anomaly Detection

The objective of anomaly detection is to determine whether future ECG values contain abnormal patterns. Given the observed signal up to time  $T$ , the task is to forecast the next  $h$  values:

$$x_{T+1}, x_{T+2}, \dots, x_{T+h}$$

The corresponding conditional distribution is

$$p(x_{T+1:T+h} \mid x_{1:T}) = \prod_{t=T+1}^{T+h} p(x_t \mid x_{1:t-1}). \quad (17)$$

For each timestep  $t$ , an anomaly score is computed from the predictive distribution or prediction error. By applying a threshold [20], a binary anomaly label is produced:

$$y_t \in \{0, 1\}$$

where  $y_t = 1$  indicates an anomalous ECG point.

### 4.2 Anomaly Diagnosis

Anomaly diagnosis focuses on identifying the exact time intervals where abnormal ECG behavior occurs. Using the anomaly scores computed for all timesteps, a sequence of diagnostic labels are produced:

$$Y = \{y_1, y_2, \dots, y_{T+h}\}$$

where each  $y_t$  marks whether the ECG signal at time  $t$  deviates from normal dynamics. Although the ECG signal is univariate, this diagnosis provides fine-grained localization of abnormal waveform segments.

## 4.3 Training Objective

Model parameters are learned through single-step prediction. At each timestep  $t$ , the model predicts  $\hat{x}_{t+1}$  from the window  $x_{t-W+1:t}$ . Given a dataset of  $N$  ECG recordings, the training objective is

$$\mathcal{L}(\Theta) = \sum_{i=1}^N \sum_{t=1}^T \ell(x_{t+1}^{(i)}, \hat{x}_{t+1}^{(i)}) \quad (18)$$

where  $\ell(\cdot, \cdot)$  denotes a prediction loss function such as mean squared error or negative log-likelihood. By learning the normal temporal dynamics of ECG signals, deviations in prediction form the basis for both anomaly detection and diagnosis.

### 4.4 Preprocessing and Window Extraction

Given a raw ECG signal

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

This formulation treats the ECG as a univariate or multivariate time series with  $T$  samples. Expressing the input in this form makes it compatible with the downstream windowing process and clarifies that all prediction and anomaly-scoring procedures operate on contiguous temporal segments. Since ECG signals may exhibit varying baseline drifts or noise distributions, normalization is applied before segmentation to ensure consistent feature scaling. The sequence is normalized and divided into fixed-length windows

$$W_t \in \mathbb{R}^{C \times L}$$

Each window  $W_t$  contains  $L$  consecutive samples and  $C$  channels (leads). This transformation creates fixed-size training examples from long ECG signals, enabling batch processing and stable gradient flow. Overlapping windows capture temporal continuity and allow the model to localize anomalies more accurately. The channel first layout also conforms to TCN and convolutional operator requirements.

### 4.5 Forward Propagation

Each window undergoes:

- **Dual-TCN encoding:**

$$H^{(loc)} = f_{L\_TCN}(W)$$

$$H^{(glob)} = f_{G\_TCN}(W)$$

The two TCN branches extract complementary temporal features, local tcn captures short-term morphology such as QRS spikes, while global tcn focuses on long-range temporal context by using larger dilation factors. This dual design enables the model to jointly learn fast transitions and slower cardiac rhythm changes. The outputs serve as the foundational temporal descriptors for subsequent attention and fusion stages.

- **Multi-scale feature extraction:**

$$F_3 = \text{Conv}_3(W), \quad F_5 = \text{Conv}_5(W)$$

$$F^{(ms)} = g_{\text{fuse}}([F_3; F_5]) \quad (19)$$

The two convolution operators with kernel sizes 3 and 5 extract features corresponding to different ECG wave durations.  $F_3$  captures sharper structures, while  $F_5$  responds to broader waveform components. After extraction, the fusion function  $g_{\text{fuse}}$  combines the features to build a unified, multi-scale representation. This module is crucial for modeling physiological events whose durations differ significantly across cardiac cycles.

- **Lightweight attention:**

$$A^{(\text{light})} = \text{Attn}(W)$$

This attention module allows the model to focus on diagnostically important time points directly from the raw window. As a lightweight self-attention mechanism with two heads and no dropout, it detects temporal irregularities without introducing excessive computational overhead. The attention weights emphasize salient patterns such as sudden deviations or irregular beats, amplifying features needed for anomaly detection.

- **Enhanced attention:**

$$A^{(\text{enh})} = \text{Attn}(H^{(\text{glob})})$$

Enhanced attention is applied to the globally encoded features, allowing long range temporal relationships to influence the representation. This refinement stage uncovers subtle correlations that span extended durations, such as low frequency waveform shifts or periodic irregularities. Using attention strengthens the model’s ability to detect anomalies that manifest through distributed or gradual changes over time.

- **Fixed-weight fusion.**

$$\tilde{H}_1 = H^{(\text{loc})} + \alpha_1 A^{(\text{light})},$$

$$\tilde{H}_2 = H^{(\text{glob})} + \alpha_2 F^{(\text{ms})} + \alpha_3 A^{(\text{enh})} \quad (20)$$

This fusion step integrates three distinct information streams, local temporal patterns, multi-scale morphology, and global attention refinements. Fixed fusion weights  $(\alpha_1, \alpha_2, \alpha_3)$  prevent unstable learning behavior and ensure that each component contributes consistently.  $\tilde{H}_1$  focuses on short scale morphology enhanced by attention, while  $\tilde{H}_2$  combines long range structure with multi scale convolution cues, forming a comprehensive fused representation.

- **Reconstruction:**

$$\hat{W} = f_{\text{Decoder}}(\tilde{H}_1, \tilde{H}_2)$$

The decoder reconstructs the input window from the fused latent features, serving as the basis for anomaly scoring. Accurate reconstruction is expected when the input window corresponds to normal cardiac activity. Any deviation from normal morphology or rhythm results in higher reconstruction error. This reconstruction based framework removes the need for labeled anomaly data and fits seamlessly into a self supervised learning paradigm.

## 4.6 Loss Function

$$\mathcal{L}(W, \hat{W}) = \frac{1}{L} \sum_{t=1}^L \|W(:, t) - \hat{W}(:, t)\|_2^2 \quad (21)$$

The reconstruction loss measures the average squared difference between the original window and its reconstruction across all time steps. By minimizing this objective, the model learns to accurately reproduce normal ECG morphology, while anomalous segments unlike anything encountered during training lead to larger errors. The normalization by  $L$  ensures stable gradient magnitudes across different window lengths and supports smooth optimization dynamics.

## 4.7 Optimization

During training, the parameters of the MED-Net model are updated by minimizing the reconstruction loss using gradient based optimization [23]. Let  $\theta$  denote the complete set of learnable parameters across the local and global TCN stacks, the multi scale convolutions, the lightweight and enhanced attention modules, and the decoder. A generic gradient descent update takes the form:

$$\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L} \quad (22)$$

where  $\eta$  is the learning rate. In practice, the optimization relies on the Adam optimizer [23], which maintains exponential moving averages of both the gradient and its squared magnitude to stabilize convergence over long ECG sequences. For each update step, Adam computes:

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) \nabla_{\theta} \mathcal{L}_t$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) (\nabla_{\theta} \mathcal{L}_t)^2 \quad (23)$$

followed by bias-corrected updates:

$$\theta_{t+1} = \theta_t - \eta \frac{\hat{m}_t}{\sqrt{\hat{v}_t + \epsilon}} \quad (24)$$

To further improve training stability, a learning-rate decay schedule is used, reducing  $\eta$  gradually as training proceeds. This encourages fast exploration during early epochs and finer adjustments near convergence. Mini-batch training with batch size  $B$  reduces gradient variance and accelerates optimization, while L2 weight decay helps prevent overfitting by penalizing excessively large weights. Together, these strategies ensure stable and efficient optimization of the multi-branch TCN–attention framework.



**Algorithm 1** Training Procedure for the MED-Net Model

**Input:** ECG signal  $X$ , window length  $L$ , learning rate  $\eta$ , fusion weights  $(\alpha_1, \alpha_2, \alpha_3)$ , batch size  $B$ , epochs  $E$ .

**Output:** Trained parameters  $\theta$ .

```

1 Normalize ECG to  $[0, 1]$ . Extract windows  $W_t \in \mathbb{R}^{C \times L}$  to
  form dataset  $\mathcal{W}$ . Form minibatches of size  $B$ . for  $e = 1$  to
   $E$  do
2   foreach minibatch  $\mathcal{B} \subset \mathcal{W}$  do
3     foreach window  $W \in \mathcal{B}$  do
4        $H^{(loc)} \leftarrow f_{\text{LocalTCN}}(W)$   $H^{(glob)} \leftarrow$ 
         $f_{\text{GlobalTCN}}(W)$ 
5        $F_3 \leftarrow \text{Conv}_3(W)$ 
6        $F_5 \leftarrow \text{Conv}_5(W)$ 
7        $F^{(ms)} \leftarrow g_{\text{fuse}}([F_3; F_5])$ 
8        $A^{(light)} \leftarrow \text{Attn}(W)$ 
9        $A^{(enh)} \leftarrow \text{Attn}(H^{(glob)})$   $\tilde{H}_1 \leftarrow H^{(loc)} +$ 
         $\alpha_1 A^{(light)}$ 
10       $\tilde{H}_2 \leftarrow H^{(glob)} + \alpha_2 F^{(ms)} + \alpha_3 A^{(enh)} \leftarrow$ 
         $f_{\text{Decoder}}(\tilde{H}_1, \tilde{H}_2)$ 
11       $\mathcal{L} \leftarrow \frac{1}{L} \sum_{t=1}^L \|W(:, t) - \hat{W}(:, t)\|_2^2$ 
12     $\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}$ 

```

## 5 Results

We evaluated the MED-Net model on two datasets: the **ECG dataset** containing univariate and multivariate ECG channels, and the **MBA** (Multi-Branch Anomaly) dataset [24] containing multi-sensor temporal records. Both datasets were preprocessed using z-score normalization and sliding-window segmentation as described earlier.

### 5.1 ECG dataset

The ECG dataset consists of 190 univariate ECG recordings for training and 48 recordings for testing, each containing 17,479 time steps and a single feature channel. After segmentation, this resulted in a total of 111,808 training windows and 83,856 testing windows. The ECG\_DATA dataset provides high-resolution cardiac waveform structure suitable for evaluating morphological reconstruction quality, while the MBA [24] dataset includes heterogeneous multivariate signals that allow assessment of the model’s ability to generalize across non-physiological domains. Across both datasets, the MED-Net achieves extremely high anomaly-detection accuracy. The improvements stem from the enlarged receptive field in the local TCN branch, the lightweight ECG attention module, the refined enhanced attention applied to the global TCN outputs, and the fixed weighted fusion mechanism that stabilizes multi-scale temporal feature aggregation. On the ECG\_DATA test set, the model achieves an **F1 score of 0.99978**, with **precision of 0.99957** and **perfect recall (1.0)**, correctly identifying all anomalous windows. Only 23 false positives were observed across more than 84,000 evaluated windows, and no false negatives were recorded. The ROC-AUC score of

**0.99961** demonstrates near-perfect discrimination. The complete evaluation metrics are shown in Table 2.

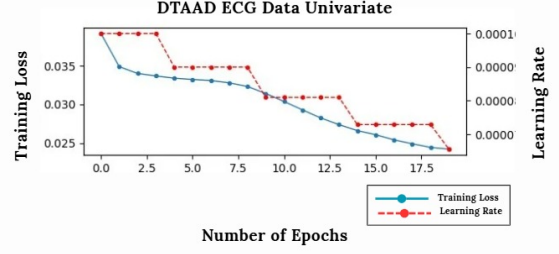


Fig. 2: Training graph of DTAAD [1] on ECG data for 20 epochs

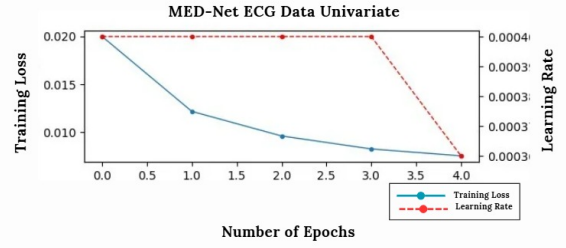


Fig. 3: Training graph of MED-Net on ECG data for 5 epochs

Table 1: DTAAD Performance on ECG data 20 epochs

Metric	Value
F1-Score	0.99478
Precision	0.98962
Recall	1.00000
True Positives (TP)	54,157
True Negatives (TN)	29,131
False Positives (FP)	568
False Negatives (FN)	0
ROC-AUC	0.99044
Threshold	0.14710

Table 2: Performance of the MED-Net on ECG data 5 epochs

Metric	Value
F1-Score	0.99978
Precision	0.99957
Recall	1.00000
True Positives (TP)	54,157
True Negatives (TN)	29,676
False Positives (FP)	23
False Negatives (FN)	0
ROC-AUC	0.99961
Threshold	0.2585

## 5.2 MBA Dataset

The MBA (Multi-Branch Anomaly) dataset [24] is a multi-variate benchmark consisting of two-channel temporal sensor recordings collected under realistic industrial conditions [24]. The dataset is processed in a 2D format (time  $\times$  features) with two channels and preserves full temporal resolution by using overlapping windows: the preprocessing pipeline produces **7,680 overlapping windows** (all timesteps preserved), which are used as both training and testing samples. In the training log the dataset is reported as:

- Detected multivariate time series data (MBA) – 2 features (channels).
- Windowed shape (overlapping windows): (7680, 2).
- Training samples: `torch.Size([7680, 2])`.
- Testing samples: `torch.Size([7680, 2])`.
- Number of features: 2 (dual-channel).
- Labels: original labels preserved one-to-one with overlapping windows.
- MBA-specific hyperparameters: epochs = 50, learning rate = 0.005, weight decay = 1e-4.
- Windowing: overlapping windows (no downsampling) to maximize sensitivity to short, subtle anomalies.

Using this configuration, we evaluated both the MED-Net and the baseline DTAAD [1]. For MBA [24] the model computes per-channel metrics and reports the channel-wise average as the final metric for the dataset.

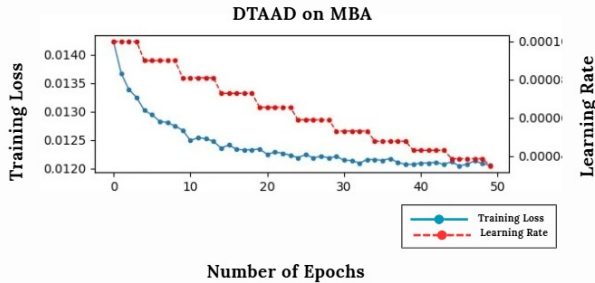


Fig. 4: Training graph of DTAAD on MBA data for 50 epochs

**Discussion.** The MBA [24] dataset is substantially more heterogeneous and noisier than the ECG\_DATA physiological set; therefore, model tuning (e.g., an increased learning rate, longer training of 50 epochs, and overlapping windows) was applied to improve sensitivity and F1 performance. The MED-Net outperforms the Original DTAAD [1] on MBA [24] by achieving higher F1, precision and AUC while maintaining zero false negatives (perfect recall) after averaging across

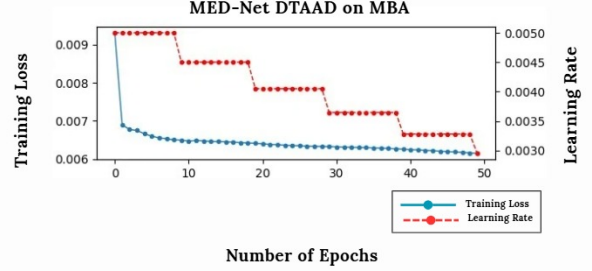


Fig. 5: MED-Net training graph on MBA data for 50 epochs

Table 3: MED-Net — MBA (Averaged across channels)

Metric	Value
F1-Score	0.96501
Precision	0.93240
Recall	1.00000
True Positives (TP)	2600
True Negatives (TN)	4892
False Positives (FP)	189
False Negatives (FN)	0
ROC-AUC	0.98145
Threshold	0.02157

Table 4: Original DTAAD — MBA (Averaged across channels)

Metric	Value
F1-Score	0.93844
Precision	0.92525
Recall	0.95385
True Positives (TP)	2480
True Negatives (TN)	4878
False Positives (FP)	203
False Negatives (FN)	120
ROC-AUC	0.95699
Threshold	0.03530

channels. This indicates that the architectural improvements and MBA-specific hyperparameter tuning increase sensitivity to industrial anomalies without sacrificing robustness. These results validate that the proposed enhancements substantially strengthen the baseline DTAAD architecture, improving training stability, sensitivity to ECG waveform morphology, and robustness to noise. The MED-Net model remains computationally efficient while achieving state-of-the-art anomaly detection performance across both physiological and general multi-variate time-series data.

## 6 Conclusion

Our extensive experimental evaluation on the ECG and MBA [24] public datasets demonstrates that MED-Net achieves state-of-the-art performance, surpassing DTAAD and CAE

baselines in terms of F1 scores and diagnostic accuracy, while simultaneously reducing training time by up to 50% compared to DTAAD and 99% compared to CAE. Despite these promising results, there are several avenues for future research to further refine the framework.

Currently, the feature fusion component relies on fixed scalar weights to combine local, global, and multi-scale information. Future work will focus on making this fusion component adaptive, employing learnable gating mechanisms or attention-based fusion to dynamically weigh feature streams based on real-time signal complexity. While the current Dual TCN backbone utilizes fixed kernel sizes and dilation factors, future iterations could employ Neural Architecture Search (NAS) to automatically discover optimal receptive field configurations for specific arrhythmia types.

While the model has proven robust on the tested benchmarks, future plans include validation of MED-Net on a broader range of biomedical datasets, such as EEG or EMG data, and larger, multi-center clinical ECG repositories to ensure generalization across diverse patient demographics and recording devices. Ultimately, MED-Net offers a lightweight, high-precision solution suitable for deployment in real-time clinical decision support systems.

## References

- [1] L. Yu, "DTAAD: Dual Tcn-Attention Networks for Anomaly Detection in Multivariate Time Series," *arXiv preprint arXiv:2302.10753v3*, 2024.
- [2] U. Lomoio, P. Vizza, R. Giancotti, S. Petrolo, S. Flesca, F. Boccuto, P. H. Guzzi, P. Veltri, and G. Tradigo, "A convolutional autoencoder framework for ECG signal analysis," *Heliyon*, vol. 11, p. e41517, 2025.
- [3] M. Gupta, J. Gao, C. C. Aggarwal, and J. Han, "Outlier detection for temporal data: A survey," *IEEE Transactions on Knowledge and Data Engineering*, vol. 26, no. 9, pp. 2250–2267, 2013.
- [4] S.-E. Benkabou, K. Benabdeslem, V. Kraus, K. Bourhis, and B. Canitia, "Local anomaly detection for multivariate time series by temporal dependency based on poisson model," *IEEE Transactions on Neural Networks and Learning Systems*, 2021.
- [5] W. Chen, H. Xu, Z. Li, D. Pei, J. Chen, H. Qiao, Y. Feng, and Z. Wang, "Unsupervised anomaly detection for intricate kpis via adversarial training of vae," in *IEEE INFOCOM 2019-IEEE Conference on Computer Communications*, 2019.
- [6] W. Luo, W. Liu, D. Lian, J. Tang, L. Duan, X. Peng, and S. Gao, "Video anomaly detection with sparse coding inspired deep neural networks," *IEEE transactions on pattern analysis and machine intelligence*, vol. 43, no. 3, pp. 1070–1084, 2019.
- [7] G. Pang, C. Shen, L. Cao, and A. V. D. Hengel, "Deep learning for anomaly detection: A review," *ACM Computing Surveys (CSUR)*, vol. 54, no. 2, pp. 1–38, 2021.
- [8] R. Chalapathy and S. Chawla, "Deep learning for anomaly detection: A survey," *arXiv preprint arXiv:1901.03407*, 2019.
- [9] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning representations by back-propagating errors," *nature*, vol. 323, no. 6088, pp. 533–536, 1986.
- [10] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [11] M. Nawaz, J. Ahmed, "Cloud-based healthcare framework for real-time anomaly detection and classification of 1-d ecg signals," *PLoS ONE*, vol. 17, no. 12, p. e0279305, 2022.
- [12] T. Ergen and S. S. Kozat, "Unsupervised anomaly detection with lstm neural networks," *IEEE transactions on neural networks and learning systems*, vol. 31, no. 8, pp. 3127–3141, 2019.
- [13] Y. Su, Y. Zhao, C. Niu, R. Liu, W. Sun, and D. Pei, "Robust anomaly detection for multivariate time series through stochastic recurrent neural network," in *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 2828–2837, 2019.
- [14] A. Vaswani et al., "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [15] S. Tuli, G. Casale, and N. R. Jennings, "TranAD: Deep Transformer Networks for Anomaly Detection in Multivariate Time Series Data," *Proceedings of VLDB*, vol. 15, no. 6, pp. 1201–1214, 2022.
- [16] K. Xu, S. Guo, N. Cao, D. Gotz, A. Xu, H. Qu, Z. Yao, and Y. Chen, "Ecglens: interactive visual exploration of large scale ECG data for arrhythmia detection," in *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, pp. 1–12, 2018.
- [17] L. Shan, Y. Li, H. Jiang, P. Zhou, J. Niu, R. Liu, Y. Wei, J. Peng, H. Yu, X. Sha, and S. Chang, "Abnormal ecg detection based on an adversarial autoencoder," *Frontiers in Physiology*, vol. 13, 2022.
- [18] S. Bai, J. Z. Kolter, and V. Koltun, "An empirical evaluation of generic convolutional and recurrent networks for sequence modeling," *arXiv preprint arXiv:1803.01271*, 2018.
- [19] D. Li, D. Chen, B. Jin, L. Shi, J. Goh, and S. Ng, "Madgan: Multivariate anomaly detection for time series data with generative adversarial networks," in *International Conference on Artificial Neural Networks*, pp. 703–716, 2019.

- [20] A. Siffer, P.-A. Fouque, A. Termier, and C. Largouet, “Anomaly detection in streams with extreme value theory,” in *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2017.
- [21] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [22] A. L. Maas, A. Y. Hannun, A. Y. Ng et al., “Rectifier nonlinearities improve neural network acoustic models,” in *Proc. ICML*, vol. 30, no. 1, 2013.
- [23] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [24] G. B. Moody and R. G. Mark, “The impact of the mit-bih arrhythmia database,” *IEEE Engineering in Medicine and Biology Magazine*, vol. 20, no. 3, pp. 45–50, 2001.
- [25] Y. Zhang, Y. Chen, J. Wang, and Z. Pan, “Unsupervised deep anomaly detection for multi-sensor time-series signals,” *IEEE Transactions on Knowledge and Data Engineering*, 2021.
- [26] S. Choi, K. Choi, H. K. Yun, S. H. Kim, H.-H. Choi, Y.-S. Park, and S. Joo, “Diagnosis of atrial fibrillation based on ai-detected anomalies of ecg segments,” *Heliyon*, vol. 10, no. 1, 2024.
- [27] S. K. Pandey and R. R. Janghel, “Hybrid Deep Learning Approach Using BiGRU-BiLSTM and Multilayered Dilated CNN for ECG Signal Classification,” *Arabian Journal for Science and Engineering*, vol. 47, pp. 10271–10283, 2022.
- [28] M. Sharma, S. Tan, and U. R. Acharya, “Simple 1D CNN with Leaky-ReLU for ECG Classification,” *Computers in Biology and Medicine*, vol. 114, p. 103399, 2019.
- [29] Y. Xiang, Z. Lin, and J. Meng, “Automatic QRS Complex Detection using Two-Level CNN,” *IEEE Biomedical Engineering Letters*, vol. 8, pp. 57–65, 2018.
- [30] J. Chen, S. Sathe, C. Aggarwal, and D. Turaga, “Transformer-based Multivariate Time Series Anomaly Localization,” in *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2021.
- [31] S. Kiranyaz, T. Ince, and M. Gabbouj, “Global ECG Classification by Self-Operational Neural Networks with Feature Injection,” *IEEE Transactions on Biomedical Engineering*, vol. 67, no. 11, pp. 3256–3266, 2020.
- [32] T. Nazzal, “Fuzz-ClustNet: Fuzzy Clustering Networks for Robust ECG Anomaly Detection,” *IEEE Access*, vol. 11, pp. 10452–10461, 2023.

## Implementation Details

### System Specifications

All experiments were conducted on a standard consumer laptop with the following configuration:

- **CPU:** Intel Core i5–13420H (8 cores)
- **RAM:** 16 GB
- **GPU:** Integrated Intel UHD Graphics
- **Operating System:** Windows (64-bit)
- **Deep Learning Framework:** PyTorch 2.8.0+cpu

All training and inference operations were executed on the CPU. The optimized design of MED-Net, including efficient attention modules and reduced-memory windowing, allowed training to be completed reliably on this non-GPU setup.

### 6.1 Training Hyperparameters

#### ECG\_DATA (Univariate ECG)

- Epochs: 5
- Learning rate:  $4 \times 10^{-4}$
- Weight decay:  $1 \times 10^{-5}$
- Optimizer: AdamW
- Scheduler: StepLR(step = 5,  $\gamma = 0.9$ )
- Windowing: Non-overlapping windows
- Batch size:  $\min(64, N)$

#### MBA (Multivariate)

- Epochs: 50
- Learning rate:  $5 \times 10^{-3}$
- Weight decay:  $1 \times 10^{-4}$
- Optimizer: AdamW
- Scheduler: StepLR(step = 10,  $\gamma = 0.9$ )
- Windowing: Overlapping sliding windows
- Batch size: Full-sequence windows
- Labels: One-to-one with overlapping windows

#### Model Defaults Across Datasets

- Window size: typically 10
- Base learning rate fallback:  $1 \times 10^{-4}$
- Batch size fallback: 64

## **Additional Notes**

- Label handling differs across datasets: window-level aggregation is applied for non-overlapping ECG windows, while overlapping MBA [24] windows preserve labels exactly.
- The dual-branch architecture produces two prediction paths whose outputs are combined using a fixed-ratio fusion strategy.
- To reduce memory usage during CPU training, large ECG recordings use non-overlapping segmentation.
- Decoder dimensions are preserved during model saving to ensure exact reproducibility.