

# Advancements in AI for Cardiovascular Diagnostics: A Synthesis of Deep Learning Models

Ikteder Akhand Udoj

November 13, 2024

# Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
<b>2</b>	<b>Background</b>	<b>5</b>
2.1	Echocardiogram Images . . . . .	6
2.2	Convolutional Neural Networks (CNNs) . . . . .	7
2.2.1	Convolutional Layers . . . . .	7
2.2.2	Pooling Layers . . . . .	7
2.2.3	Fully Connected Layers . . . . .	8
2.2.4	Training CNNs . . . . .	8
2.2.5	Applications in Medical Imaging . . . . .	8
2.2.6	Recent Advances . . . . .	9
2.3	Recurrent Neural Networks (RNN) and Long Short-Term Memory (LSTM) Networks . . . . .	9
2.4	ResNet Architecture . . . . .	10
<b>3</b>	<b>Review of Automatic Detection of Congestive Heart Failure Based on a Hybrid Deep Learning Algorithm in the Internet of Medical Things</b>	<b>12</b>
3.1	Hybrid Model Architecture . . . . .	12
3.2	Data Processing Scheme . . . . .	12
3.2.1	Preprocessing . . . . .	13
3.2.2	Feature Extraction . . . . .	13
3.3	Results and Discussion . . . . .	14
3.3.1	5-Minute ECG Analysis . . . . .	14
3.3.2	1-Minute ECG Analysis . . . . .	15
3.3.3	Comparative Analysis . . . . .	15
3.4	Discussion . . . . .	15
<b>4</b>	<b>Review of ECG Heartbeat Arrhythmia Classification Using Time-Series Augmented Signals and Deep Learning Approach</b>	<b>15</b>
4.1	Model Architecture . . . . .	15
4.2	Data Processing Scheme . . . . .	16
4.3	Results and Discussion . . . . .	17
<b>5</b>	<b>Review of High-Throughput Precision Phenotyping of Left Ventricular Hypertrophy With Cardiovascular Deep Learning</b>	<b>18</b>
5.1	Model Architecture . . . . .	18
5.2	Data Processing Scheme . . . . .	18
5.3	Results and Discussion . . . . .	19
<b>6</b>	<b>Review of Echonet-Dynamic study</b>	<b>21</b>
6.1	Model Architecture . . . . .	21
6.2	Data Processing Scheme . . . . .	22
6.3	Results and Discussion . . . . .	24
<b>7</b>	<b>Challenges and Gaps in Current Research</b>	<b>25</b>
<b>8</b>	<b>Opportunities for Future Research</b>	<b>26</b>

<b>9 Computer Artifact</b>	<b>27</b>
9.1 Model Architecture . . . . .	27
9.2 Datasets: PneumoniaMNIST and BreastMNIST . . . . .	28
9.3 Model Implementations and Training . . . . .	28
9.4 Results . . . . .	28
9.5 Discussion of Model Performance . . . . .	34
<b>10 Conclusion</b>	<b>34</b>

## Abstract

Cardiovascular disease (CVD) remains a leading cause of mortality worldwide, with early and accurate diagnosis essential for improving patient outcomes. Traditional diagnostic methods, including electrocardiograms (ECGs) and echocardiography, often require extensive clinician interpretation, introducing variability and potential delays in diagnosis. This paper synthesizes advancements in artificial intelligence (AI)-driven diagnostics, specifically through deep learning models such as convolutional neural networks (CNNs), recurrent neural networks (RNNs), and hybrid architectures. We review four key studies illustrating AI's transformative potential in both ECG and echocardiogram analyses. Highlighted within this synthesis is the EchoNet-Dynamic model, which performs real-time, video-based cardiac function assessment and demonstrates accuracy levels exceeding those of human experts. Additionally, this paper presents a computational artifact assessing the performance of DeepLabv3 and CNN-RNN hybrid models on PneumoniaMNIST and BreastMNIST datasets. Results reveal promising accuracy across these medical imaging tasks, underscoring the role of AI in enhancing diagnostic precision and efficiency in cardiovascular healthcare. While challenges remain, particularly in data generalizability and clinical integration, the findings advocate for continued AI research to establish more robust, scalable diagnostic solutions.

**Keywords:** Cardiovascular disease, Artificial Intelligence, Deep Learning, Convolutional Neural Networks, Recurrent Neural Networks, Echocardiography, Electrocardiogram, Medical Imaging

## 1 Introduction

Cardiovascular diseases (CVDs) remain the foremost cause of morbidity and mortality worldwide, responsible for approximately 17.9 million deaths each year, which translates to nearly 31% of all global deaths [1]. CVD encompasses a range of heart and vascular conditions, including congestive heart failure (CHF), arrhythmias, hypertrophic cardiomyopathy (HCM), and left ventricular hypertrophy (LVH), many of which require timely diagnosis and management to mitigate adverse health outcomes. The high prevalence of CVD places a significant burden on healthcare systems, prompting an urgent need for more efficient diagnostic methods that can ensure accurate and early detection [6].

Traditional diagnostic methods for CVD, such as electrocardiograms (ECGs) and echocardiography, play an essential role in assessing heart structure and function. ECGs capture the electrical activity of the heart, providing valuable insights into arrhythmic conditions, while echocardiography visualizes cardiac anatomy, enabling assessment of the left ventricular ejection fraction (LVEF) and ventricular wall thickness [4, 5]. Although these techniques are well-established, they heavily depend on the expertise and experience of clinicians. Studies indicate that human interpretation of echocardiograms and ECGs can vary significantly across observers, particularly in complex or borderline cases, leading to delays in diagnosis and variability in treatment outcomes [12]. This variability is particularly problematic in resource-limited healthcare settings, where trained cardiologists may not be available [3].

Recent advances in artificial intelligence (AI) and deep learning (DL) present promising avenues for improving cardiovascular diagnostics. AI-driven diagnos-

tic tools, particularly those using deep learning models such as convolutional neural networks (CNNs) and recurrent neural networks (RNNs), have demonstrated remarkable success in processing complex medical data [7]. These models offer capabilities for automatic detection of subtle features within medical images and time-series data, thereby enhancing diagnostic accuracy, reducing time requirements, and minimizing human error. CNNs excel at capturing spatial information within images, making them well-suited for echocardiographic analysis, while RNNs handle sequential data, which is crucial for analyzing time-series data from ECGs [8].

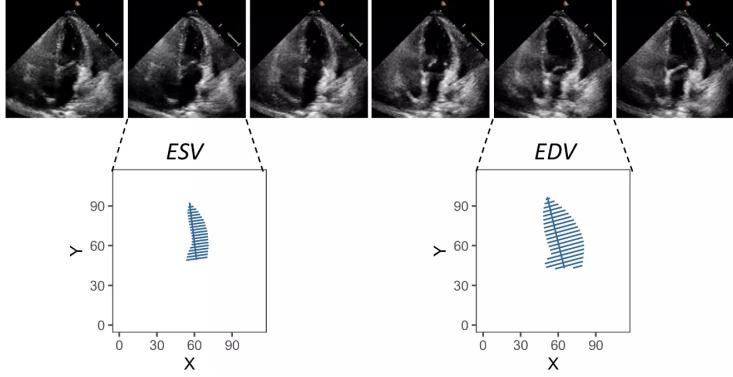
This synthesis examines four pivotal studies that illustrate the application of AI in the diagnosis of cardiovascular diseases. Two studies focus on ECG data, employing hybrid deep learning models that combine CNNs and RNNs to capture both spatial and temporal information in the detection of CHF and arrhythmias. The first study developed a hybrid CNN-RNN model for CHF detection, leveraging both spatial and temporal features to achieve high accuracy [8]. The second study expanded on these findings by implementing data augmentation techniques, enhancing the model's performance in classifying arrhythmias, even with limited datasets [9]. Both studies underscore the transformative potential of AI in ECG-based diagnostics by providing more reliable tools for detecting heart abnormalities.

The other two studies explored AI applications in echocardiography. The third study utilized deep learning to detect LVH, a condition indicative of underlying HCM or cardiac amyloidosis, by automating left ventricular segmentation and measuring ventricular dimensions with high precision [10]. The fourth study introduced EchoNet-Dynamic, a cutting-edge video-based AI model, which assesses cardiac function on a beat-to-beat basis, outperforming human experts in ejection fraction prediction and heart failure classification [12]. EchoNet-Dynamic's architecture integrates spatiotemporal convolutions, enabling it to analyze entire echocardiogram videos, thus capturing dynamic variations in cardiac function.

The goal of this paper is to synthesize these advancements, exploring the contributions of AI-driven approaches in cardiovascular diagnostics while addressing the remaining challenges and limitations. The integration of AI into clinical practice holds substantial promise for delivering more accurate, timely, and cost-effective diagnostics. However, continued research is required to validate these models across diverse populations and to address issues related to data quality, variability in clinical settings, and acceptance within the medical community.

## 2 Background

As machine Learning and deep Learning are becoming crucial for medical advancements, the following sections talks about some of the important machine learning techniques and models that are used in current research to overcome human limitations and make the process of diagnosis much faster and possibly more accurate.



**Figure 1:** Standard Echocardiogram Sequence Showing End-Systolic Volume (ESV) and End-Diastolic Volume (EDV) with Calculation of Ejection Fraction (EF). Each frame represents a phase in the cardiac cycle, captured to visualize the dynamic function of the left ventricle.[12].

## 2.1 Echocardiogram Images

Echocardiography is a diagnostic imaging technique that uses ultrasound waves to produce visualizations of the heart’s structure and function. This non-invasive method allows clinicians to examine various aspects of cardiac health, including chamber size, wall motion, and blood flow. Echocardiograms are invaluable in diagnosing and monitoring a wide range of cardiovascular conditions such as heart failure, valve disease, and congenital heart defects [2].

Figure 1 illustrates a standard echocardiogram sequence focusing on the left ventricle. Echocardiogram images are often organized into frames captured over time to show the dynamic nature of the heart’s motion. Key phases of the cardiac cycle, such as the end-systolic volume (ESV) and end-diastolic volume (EDV), are highlighted in the frames. These volumes are critical metrics for calculating the heart’s ejection fraction (EF), a measure of the heart’s pumping efficiency.

In the figure, a series of echocardiographic frames are presented, showing the heart’s contraction and relaxation cycles. The ESV and EDV are marked, allowing for the calculation of the ejection fraction as follows:

$$EF(\%) = \frac{EDV - ESV}{EDV} \times 100$$

where: - *EDV*: End-Diastolic Volume, the volume of blood in the ventricle at the end of filling (diastole), - *ESV*: End-Systolic Volume, the volume of blood remaining in the ventricle after contraction (systole).

This specific dataset, sourced from Stanford University Hospital, comprises 10,030 apical-4-chamber echocardiography videos collected from individuals undergoing clinical care between 2016 and 2018. Each video was preprocessed by cropping and masking to exclude text and information outside the scanning area. The frames were then downsampled to a standardized resolution of 112x112 pixels for analysis. This preprocessing step ensures consistency across the dataset and enables efficient computational processing in AI applications for cardiac function assessment [12].

## 2.2 Convolutional Neural Networks (CNNs)

Convolutional Neural Networks (CNNs) are a type of deep learning model that has proven to be highly effective for image classification and other computer vision tasks. Unlike traditional neural networks, which treat input data as a single-dimensional array, CNNs are designed to process and capture the spatial and temporal dependencies in multi-dimensional arrays such as images. This unique structure makes CNNs particularly well-suited for image-related tasks, as they can automatically and adaptively learn spatial hierarchies of features from input images [22].

**Architecture of CNNs** A typical CNN architecture comprises three primary layers: convolutional layers, pooling layers, and fully connected layers which is shown in figure 2. Each layer plays a distinct role in transforming the input data into a higher-level, more abstract representation that can be used for classification or other tasks.

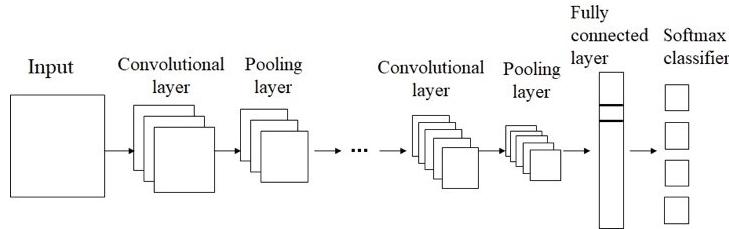


Figure 2: Basic Structure of CNN

### 2.2.1 Convolutional Layers

The convolutional layer is the core component of CNNs, responsible for learning and extracting features from the input data. In this layer, a set of filters (or kernels) slide over the input image, performing a dot product operation with the local receptive fields of the input and producing feature maps. This operation is repeated across the entire image, capturing local patterns such as edges, textures, and shapes. The depth of the convolutional layer corresponds to the number of filters used, with each filter learning to detect different features in the image. Mathematically, the output of a convolutional layer can be represented as:

$$Y_{i,j} = (X * K)_{i,j} = \sum_m \sum_n X_{i+m, j+n} K_{m,n}$$

where  $X$  is the input image,  $K$  is the filter kernel, and  $*$  denotes the convolution operation [23].

### 2.2.2 Pooling Layers

After convolutional layers, CNNs typically use pooling layers to progressively reduce the spatial dimensions of the feature maps. This downsampling operation not only reduces the computational complexity but also helps prevent overfitting. Max pooling, the most common type of pooling, selects the maximum

value from each patch of the feature map, effectively capturing the most prominent feature within each region. Average pooling, on the other hand, computes the average value within each patch. The pooling operation can be defined as:

$$Y_{i,j} = \max(X_{i:a,j:b})$$

where  $a$  and  $b$  define the region over which the pooling is performed [20].

### 2.2.3 Fully Connected Layers

Following the convolutional and pooling layers, fully connected layers are used to combine features and make predictions. In these layers, each neuron is connected to every neuron in the previous layer, allowing the network to integrate features across the entire image. The output of the last fully connected layer is passed through an activation function (e.g., softmax for classification tasks), producing the final output probabilities for each class.

**Activation Functions:** Activation functions introduce non-linearity into the CNN, enabling it to learn complex mappings. Rectified Linear Unit (ReLU) is the most common activation function in CNNs, defined as:

$$f(x) = \max(0, x)$$

ReLU helps to mitigate the vanishing gradient problem by allowing gradients to flow through the network when the input is positive. Other activation functions, such as sigmoid and tanh, are occasionally used but tend to be less effective in deep networks due to their gradient limitations [24].

### 2.2.4 Training CNNs

CNNs are trained using a supervised learning approach, typically with backpropagation and gradient descent algorithms. During training, the network adjusts the weights of its filters and connections based on the errors in its predictions. The loss function (e.g., cross-entropy for classification) measures the difference between predicted and actual outputs, and the gradient of this loss with respect to each parameter is computed. Using the gradients, parameters are updated iteratively to minimize the loss and improve the network's performance on the training data [25].

### 2.2.5 Applications in Medical Imaging

In medical imaging, CNNs have been widely applied for tasks such as disease detection, tissue segmentation, and anatomical structure recognition. For example, CNNs have been successfully used to classify lung conditions in X-ray images, detect tumors in MRI scans, and segment blood vessels in retinal images. These applications benefit from CNNs' ability to capture intricate spatial patterns that may be indicative of medical conditions, thereby supporting clinicians in diagnosing and evaluating diseases with greater accuracy [27].

**Limitations of CNNs** Despite their success, CNNs face certain limitations. They require large amounts of labeled data for training, which can be challenging to obtain in medical contexts. CNNs are also computationally intensive, necessitating significant hardware resources for both training and deployment. Additionally, CNNs are sensitive to variations in image quality and can struggle with generalizability across different datasets if not trained with a diverse dataset [28].

#### 2.2.6 Recent Advances

Recent advancements in CNN architectures, such as the introduction of Residual Networks (ResNet) and DenseNet, have further improved CNNs' performance by addressing issues like vanishing gradients and enhancing feature reuse [29]. These architectures allow CNNs to be deeper and more robust, making them even more effective for complex image analysis tasks.

### 2.3 Recurrent Neural Networks (RNN) and Long Short-Term Memory (LSTM) Networks

Recurrent Neural Networks (RNNs) are a class of neural networks specifically designed to handle sequential data by maintaining a form of memory, allowing them to capture dependencies over time. RNNs are widely used in tasks that require understanding the temporal dynamics of data, such as language modeling, time-series forecasting, and speech recognition. However, traditional RNNs face challenges with long-term dependencies due to issues like vanishing gradients, which can hinder their ability to learn patterns over extended sequences [37].

**Structure of RNNs:** The core of an RNN is its recurrent structure, where the hidden state from the previous time step is used along with the current input to determine the next hidden state. Mathematically, an RNN computes the hidden state  $h_t$  at time step  $t$  as follows:

$$h_t = f(W_x x_t + W_h h_{t-1} + b) \quad (1)$$

where  $x_t$  is the input at time  $t$ ,  $W_x$  and  $W_h$  are weight matrices,  $b$  is a bias term, and  $f$  is an activation function, typically the hyperbolic tangent ( $\tanh$ ) or ReLU. The output  $y_t$  at each step can then be derived from the hidden state  $h_t$  [23].

While effective for short sequences, RNNs struggle with long sequences due to the vanishing gradient problem, where gradients of earlier layers diminish exponentially as they are propagated back through time. This limitation restricts traditional RNNs from learning dependencies that span many time steps.

**Long Short-Term Memory (LSTM) Networks:** To address the shortcomings of standard RNNs, Long Short-Term Memory (LSTM) networks were developed by Hochreiter and Schmidhuber in 1997 [38]. LSTMs introduce a more complex memory cell structure that includes mechanisms to regulate the flow of information, thus allowing the network to retain relevant information over longer sequences and mitigating the vanishing gradient issue.

The LSTM cell includes three gates:

- **Forget Gate:** Determines the amount of information from the previous cell state  $C_{t-1}$  that should be discarded. It is computed as:

$$f_t = \sigma(W_f x_t + U_f h_{t-1} + b_f) \quad (2)$$

where  $\sigma$  is the sigmoid activation function.

- **Input Gate:** Controls the information to be updated in the cell state from the current input. It has two components: the gate layer  $i_t$  and the candidate layer  $\tilde{C}_t$ , given by:

$$i_t = \sigma(W_i x_t + U_i h_{t-1} + b_i) \quad (3)$$

$$\tilde{C}_t = \tanh(W_c x_t + U_c h_{t-1} + b_c) \quad (4)$$

- **Output Gate:** Decides the information to be passed to the next hidden state and the output. It is computed as:

$$o_t = \sigma(W_o x_t + U_o h_{t-1} + b_o) \quad (5)$$

The cell state  $C_t$  and the hidden state  $h_t$  are updated as follows:

$$C_t = f_t \odot C_{t-1} + i_t \odot \tilde{C}_t \quad (6)$$

$$h_t = o_t \odot \tanh(C_t) \quad (7)$$

This gated mechanism allows LSTMs to retain or forget information over longer sequences, making them suitable for tasks with long-term dependencies.

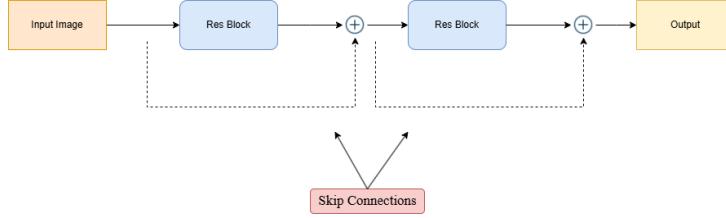
**Table 1:** Comparison of RNN and LSTM in Handling Long-Term Dependencies

Aspect	RNN	LSTM
Memory Capacity	Limited	Extended via cell state
Gradient Flow	Susceptible to vanishing gradient	Mitigates vanishing gradient
Training Complexity	Simpler	More complex due to gating mechanisms
Application Suitability	Short sequences	Long sequences with dependencies

## 2.4 ResNet Architecture

The ResNet (Residual Network) architecture, introduced by He et al. in 2015, was developed to tackle the challenges associated with training very deep neural networks, particularly the degradation problem, where adding more layers results in reduced accuracy. This degradation is not due to overfitting but rather an optimization difficulty, which prevents deep networks from reaching optimal solutions effectively [29]. ResNet addresses this issue through a novel approach known as *residual learning*, where the network learns residual functions instead of direct mappings.

In conventional neural networks, each layer learns a function  $H(x)$  that maps input  $x$  to the output. In contrast, ResNet layers learn the residual function  $F(x) = H(x) - x$ , reformulating the mapping to  $H(x) = F(x) + x$ . This is achieved through *shortcut connections* that bypass certain layers and directly add the input  $x$  to the output of the residual function  $F(x)$ . In figure 3 it



**Figure 3:** Basic Structure of ResNet

shows how a simple ResNet structure might look like with skip connections. This approach simplifies optimization, as residual functions are often easier to learn than direct mappings, especially when the desired transformation closely resembles an identity function.

The ResNet architecture consists of the following key components:

- **Shortcut Connections:** Identity mappings skip over one or more layers, allowing gradients to flow directly through the network. This mechanism mitigates issues related to vanishing gradients, which can impede the training of very deep networks.
- **Bottleneck Design:** To optimize computational efficiency, ResNet uses a bottleneck structure in deeper models (e.g., ResNet-50, ResNet-101, ResNet-152). This design involves reducing dimensionality with a  $1 \times 1$  convolution, performing the main computation with a  $3 \times 3$  convolution, and then restoring dimensionality with another  $1 \times 1$  convolution. This structure significantly reduces the number of computations required for each layer.
- **Scalability:** ResNet is highly scalable, supporting architectures with up to 152 layers, as seen in ResNet-152. This model demonstrated superior performance to shallower networks and required fewer computations than other complex architectures, such as VGG-19, due to its efficient use of shortcut connections [17].

ResNet achieved state-of-the-art results across multiple benchmarks, including the ImageNet classification challenge. An ensemble of ResNets achieved a top-5 error rate of 3.57%, winning the 2015 ILSVRC competition and significantly advancing the field of deep learning [18]. Since its introduction, ResNet has influenced various applications beyond image classification, including object detection and image segmentation, due to its ability to train very deep networks effectively.

### 3 Review of Automatic Detection of Congestive Heart Failure Based on a Hybrid Deep Learning Algorithm in the Internet of Medical Things

#### 3.1 Hybrid Model Architecture

The study proposes a hybrid deep learning model combining a Convolutional Neural Network (CNN) and a Recursive Neural Network (RNN) to automatically detect congestive heart failure (CHF) from electrocardiogram (ECG) signals. The CNN is used to process and extract relevant features from the time-frequency spectra of ECG signals, while the RNN captures temporal dependencies within these features, allowing the model to utilize both spatial and temporal information for CHF detection [34]. The CNN's capability in spatial feature extraction is augmented by the RNN's sequential learning ability, which is essential for analyzing time-series data like ECG.

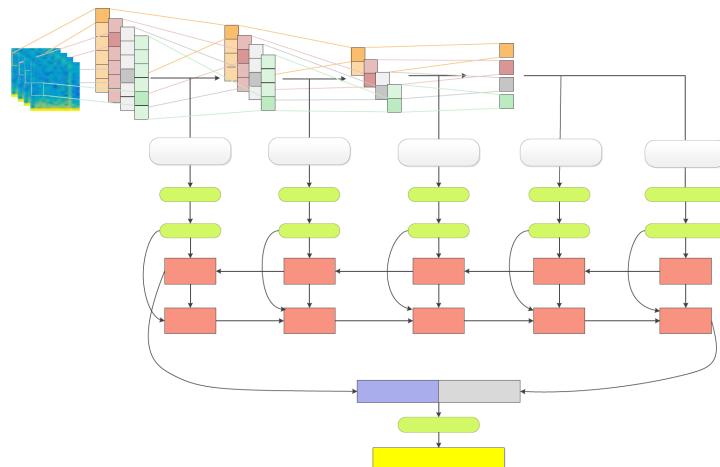
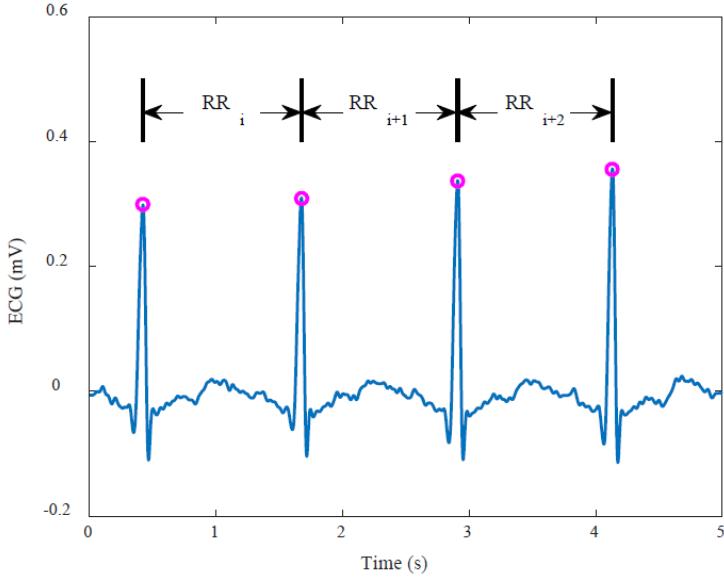


Figure 4: Proposed Structure of Hybrid Model

The architecture of the proposed hybrid model in figure 4 consists of convolutional layers in the CNN to automatically learn spatial representations from the time-frequency spectrum of ECGs. Following the CNN, the RNN component is structured as a fully connected recurrent network to capture the sequential structure of the ECG data, specifically analyzing the patterns within the RR intervals [40]. The combined architecture leverages the CNN's strengths in processing 2D data (like images) with the RNN's capability to handle sequential data, making it a suitable choice for time-series analysis in healthcare applications.

#### 3.2 Data Processing Scheme

The dataset utilized in this study includes ECG recordings from 15 CHF patients and 18 healthy individuals. The ECG signals for CHF patients were collected



**Figure 5:** ECG Signals

over approximately 20 hours, with a sampling rate of 250 Hz, while signals for healthy subjects were sampled at 128 Hz. Data processing involves several steps to clean and prepare the ECG signals for model training.

### 3.2.1 Preprocessing

The raw ECG signals underwent filtering to remove noise and baseline drift. The preprocessing pipeline includes:

- A finite impulse response (FIR) low-pass filter with a cutoff frequency of 22 Hz to remove high-frequency noise.
- A high-pass FIR filter with a cutoff at 1.2 Hz to eliminate baseline wander.
- A 60 Hz notch filter to remove power line interference.

Following noise removal, an advanced QRS detector was applied to accurately identify R-wave peaks within the ECG signal. The RR interval sequences were then constructed by measuring the intervals between consecutive R-waves.

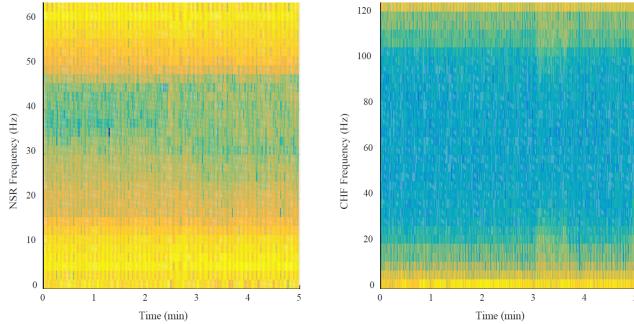
### 3.2.2 Feature Extraction

For feature extraction, both time-domain and frequency-domain indices were derived from the RR interval sequences. Time-domain features included metrics such as the standard deviation of NN intervals (SDNN) and the root mean square of successive differences (RMSSD). Frequency-domain analysis was performed using discrete Fourier transforms (DFT) to obtain low-frequency (LF) and high-frequency (HF) components. Additionally, a Poincare map was employed as a nonlinear feature to capture chaotic behaviors indicative of CHF [8].

$$SDNN = \sqrt{\frac{1}{N} \sum_{i=1}^N (R_i - \bar{R})^2} \quad (8)$$

$$RMSSD = \sqrt{\frac{1}{N-1} \sum_{i=1}^{N-1} (R_{i+1} - R_i)^2} \quad (9)$$

Time-frequency spectra were generated using the short-time Fourier transform (STFT) with a sliding window of length 32, transforming the 1D RR interval sequences into 2D spectral representations that are input into the CNN [8]. Figures for the processed frequency spectra and RR interval spectra highlight the significant differences between healthy and CHF subjects.



**Figure 6:** Frequency Spectra of ECG Signals for Healthy and CHF Subjects, showcasing spectral differences leveraged by the CNN-RNN model. Adapted from [8].

### 3.3 Results and Discussion

The model’s performance was evaluated on both 5-minute and 1-minute ECG data segments. The key performance metrics included accuracy, sensitivity, and specificity, which are crucial for clinical applications.

#### 3.3.1 5-Minute ECG Analysis

Using 5-minute ECG segments, the hybrid model achieved remarkable accuracy, sensitivity, and specificity:

- **Accuracy:** 99.93%
- **Sensitivity:** 99.85%
- **Specificity:** 100%

The high accuracy underscores the model’s potential as a reliable diagnostic tool for CHF detection, especially when compared to traditional methods. The model’s superior performance can be attributed to its ability to capture both the frequency-domain and time-domain features, as well as the inherent temporal characteristics within ECG data [30].

### 3.3.2 1-Minute ECG Analysis

For shorter 1-minute ECG segments, the model maintained strong performance, although some sensitivity reduction was observed. The results showed the hybrid model's robustness across different segment lengths, making it adaptable for various clinical settings where shorter monitoring periods might be necessary [35].

### 3.3.3 Comparative Analysis

The proposed hybrid model outperformed traditional methods in CHF detection, including Support Vector Machines (SVMs) and other machine learning classifiers used in previous studies. Table 2 shows a comparison of the hybrid model with these traditional methods, highlighting its superior accuracy.

**Table 2:** Comparison of Hybrid Model Performance with Traditional Methods

Method	Accuracy (%)	Sensitivity (%)	Specificity (%)
Hybrid Model (5 min)	99.93	99.85	100
Acharya et al. [30]	94.40 - 98.97	94.68 - 98.87	95.75 - 99.01
Sudarshan et al. [31]	97.94 - 99.87	97.04 - 99.78	97.69 - 99.94
Masetic et al. [32]	100	NP	NP
Kamath et al. [33]	79.20 - 98.20	71.50 - 98.40	87.80 - 98.00

## 3.4 Discussion

The hybrid deep learning model demonstrates a substantial improvement in CHF detection compared to traditional methods. By combining CNN and RNN architectures, the model efficiently processes both spatial and temporal information, which enhances its diagnostic accuracy and reduces the need for complex feature engineering. Additionally, the model's high sensitivity and specificity validate its capability to differentiate CHF from normal sinus rhythm, indicating its utility in clinical applications.

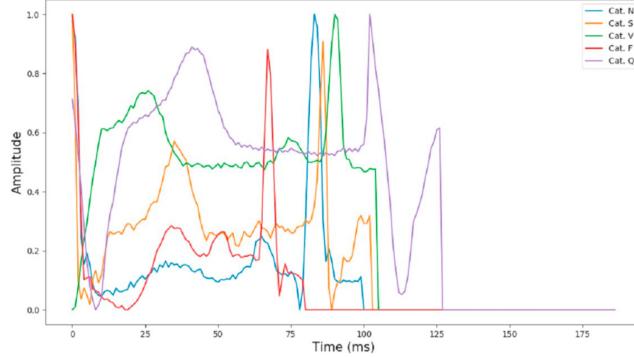
Despite its strong performance, the study recognizes certain limitations. The dataset size was relatively small, comprising only 15 CHF patients and 18 healthy subjects. Future research with a larger dataset would help in assessing the model's generalizability. Additionally, while the model achieved high accuracy on 5-minute and 1-minute data segments, further validation is required to explore its performance on ultra-short ECG segments for real-time applications [36].

## 4 Review of ECG Heartbeat Arrhythmia Classification Using Time-Series Augmented Signals and Deep Learning Approach

### 4.1 Model Architecture

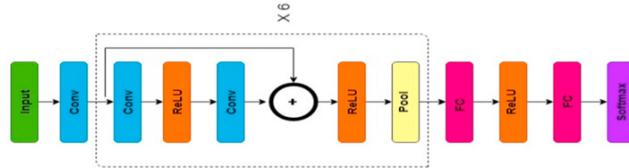
The proposed model focuses on automated ECG heartbeat arrhythmia classification using a deep learning-based approach, specifically targeting the task

with a modified architecture to improve accuracy and stability. ECG signals are inherently noisy and non-stationary, requiring preprocessing to enhance signal clarity. To facilitate effective model training, the MIT-BIH Arrhythmia dataset [41, 42] was used, encompassing five ECG signal classes (N, S, V, F, and Q) with signals sampled at 125 Hz.



**Figure 7:** The 5 classes in MIT-BIH Arrhythmia dataset [26]

Figure 8 shows the model architecture proposed in [26] builds upon one-dimensional (1D) convolutional layers along the time axis, with six residual blocks incorporated for deeper feature extraction. Each residual block comprises two convolutional layers, two ReLU activation layers [24], and a max-pooling layer, culminating in a model with 15 weighted layers. Each convolutional layer has a kernel size of 64, which has been empirically determined to enhance convergence stability and feature extraction efficacy. The final layer is a softmax activation function, providing class probability outputs.



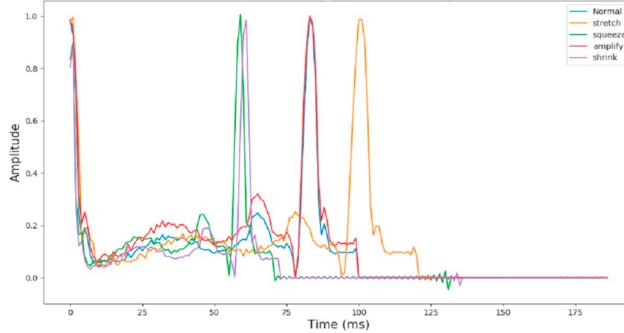
**Figure 8:** CNN based Deep Learning model [26]

The model is trained with the Adam optimizer [43], configured with a learning rate of 0.001 and exponential decay. This architecture allows the model to capture complex features within the ECG signals while preventing overfitting by using residual blocks [29].

## 4.2 Data Processing Scheme

The dataset preprocessing includes signal transformation techniques such as squeezing, stretching, amplification, and shrinking, applied to every sample in the dataset. These transformations ensure augmentation of the original data, resulting in a more comprehensive dataset that enhances model generalization. The dataset, originally containing 109,446 samples, is expanded through these transformations to encompass 547,230 samples, thereby increasing robustness during model training.

Each ECG segment represents a heartbeat, preprocessed to isolate the relevant section between R-peaks. Transformation techniques such as time-series stretching (elongating signals), squeezing (shortening signals), amplification (vertical stretching), and shrinking (vertical compressing) were applied, resulting in augmented variations of each signal sample. These augmented datasets allowed the model to learn a diverse range of signal variations, essential for distinguishing between similar waveforms in arrhythmia classification [44].



**Figure 9:** 4 transformations performed a single sample [26]

### 4.3 Results and Discussion

Experiments demonstrated the model’s improved stability and accuracy with the augmented dataset. When compared to an initial model discussed by Kachuee et al. [45], the proposed architecture achieved superior convergence rates and exhibited a higher f1-score and validation accuracy. The augmented dataset proved essential in reducing ranking loss by approximately 90% relative to the initial model trained on the original dataset.

Key performance metrics of the proposed model include a precision rate above 99%, with an f1-score similarly high at 0.98 when trained on the augmented dataset. The model demonstrated particularly stable accuracy over 120 training epochs, with confusion matrices showing improved classification accuracy for categories that traditionally exhibited overlap, such as classes ‘F’ and ‘S’. This validates the model’s capability to capture intricate distinctions within ECG waveform patterns, which are critical for precise arrhythmia detection.

Overall, the proposed deep learning approach and data augmentation techniques yielded high accuracy and robust performance, establishing a reliable framework for real-time arrhythmia classification in clinical settings.

**Table 3:** Performance comparison of initial and proposed models on original and augmented datasets

Model	Dataset	f1-score	Ranking Loss	Coverage Error
Initial Model	Original Dataset	0.89	0.0399	1.1275
Initial Model	Augmented Dataset	0.98	0.0048	1.0192
Proposed Model	Original Dataset	0.90	0.0319	1.1595
Proposed Model	Augmented Dataset	0.98	0.0047	1.0190

The results in Table 3 illustrate the enhanced stability and accuracy of the proposed model with augmented data. The f1-score achieved by the proposed model with augmentation consistently outperformed the initial model, demonstrating the importance of the residual blocks and data augmentation in achieving high classification performance for ECG arrhythmias.

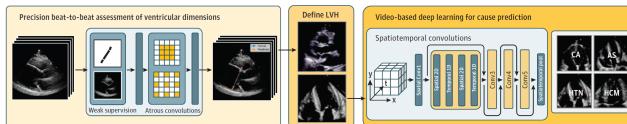
## 5 Review of High-Throughput Precision Phenotyping of Left Ventricular Hypertrophy With Cardiovascular Deep Learning

### 5.1 Model Architecture

The study developed a high-throughput, AI-based phenotyping system to automate the assessment of left ventricular hypertrophy (LVH) using echocardiography. The architecture employed consists of a deep learning model designed to:

1. Perform frame-by-frame segmentation of echocardiographic videos to identify and measure ventricular dimensions, specifically targeting left ventricular wall thickness.
2. Utilize a three-dimensional convolutional neural network (3D CNN) with residual connections to analyze video data, integrating both spatial and temporal information. This model architecture was selected to capture changes over the cardiac cycle and enhance the accuracy of measurements.

To achieve the segmentation of the left ventricular wall, a modified DeepLabv3 architecture was utilized. This network included atrous convolutions for multi-scale context capture, essential for precise cardiac measurements. The DeepLabv3 model identifies the intraventricular septum (IVS), left ventricular internal dimension (LVID), and left ventricular posterior wall (LVPW) in the parasternal long-axis echocardiogram videos.



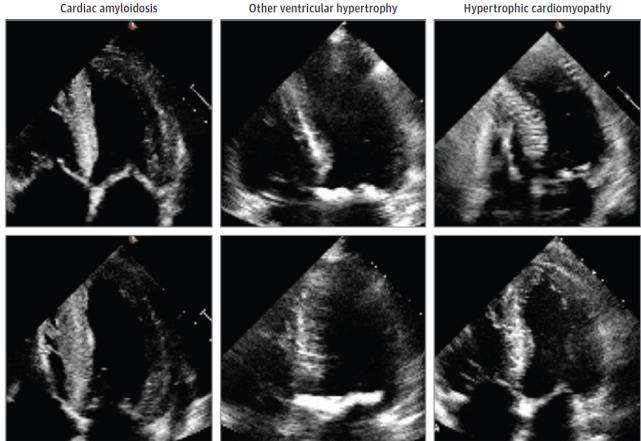
**Figure 10:** Proposed architecture by Duffy et al. [10]

For classification, a ResNet3D model architecture was used to predict the presence and cause of LVH, distinguishing between hypertrophic cardiomyopathy (HCM), cardiac amyloidosis, aortic stenosis, and other etiologies of increased LV wall thickness [10].

### 5.2 Data Processing Scheme

The dataset used in this study consisted of physician-curated cohorts from multiple centers, including the Stanford Amyloid Center, Cedars-Sinai Medical Center, and the Unity Imaging Collaborative. The data included parasternal

long-axis and apical 4-chamber echocardiographic videos from patients with diagnosed cardiac conditions. The videos were divided into training, validation, and test sets, with additional held-out test sets from external datasets.



**Figure 11:** Dataset used in Duffy et al.'s paper [10]

- **Training and Validation Data:** For training, 9600 parasternal long-axis videos were used, while 1200 were set aside for validation. The videos were preprocessed to remove identifiable patient information, and manual annotations were provided by clinicians for key measurements (e.g., IVS, LVID, LVPW).
- **External Validation:** An additional set of 3660 videos from Cedars-Sinai Medical Center and 1791 videos from Unity Imaging Collaborative were used for external validation, ensuring the model's generalizability across healthcare systems and patient populations.
- **Labeling and Annotation:** Human annotations served as ground truth labels for model training, with measures for IVS, LVID, and LVPW provided by echocardiography-certified cardiologists. The algorithm was evaluated against these human annotations to determine accuracy and robustness.

### 5.3 Results and Discussion

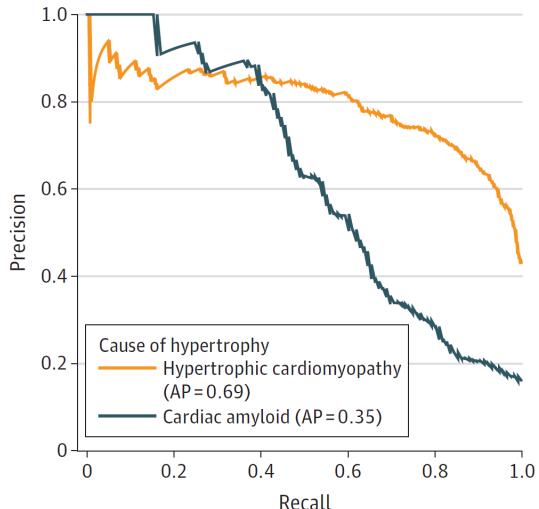
The deep learning model demonstrated high accuracy in measuring LV dimensions and in detecting LVH caused by different pathological conditions. The model achieved:

- **Measurement Accuracy:** The mean absolute error (MAE) was 1.2 mm for IVS, 2.4 mm for LVID, and 1.4 mm for LVPW. In the external datasets, the MAE remained within 1.7 mm for IVS, 3.8 mm for LVID, and 1.8 mm for LVPW, underscoring the model's consistency.
- **Classification Performance:** The model's area under the curve (AUC) values for classifying the etiology of LVH were impressive, with an AUC

of 0.98 for HCM, 0.83 for cardiac amyloidosis, and 0.89 for aortic stenosis. This performance indicates the model’s potential utility for diagnostic assistance in clinical settings.

- **Clinical Implications:** This study demonstrated that a deep learning algorithm could provide precise, reproducible measurements comparable to those of human experts, with potential applications in automated disease screening and assessment in echocardiography.

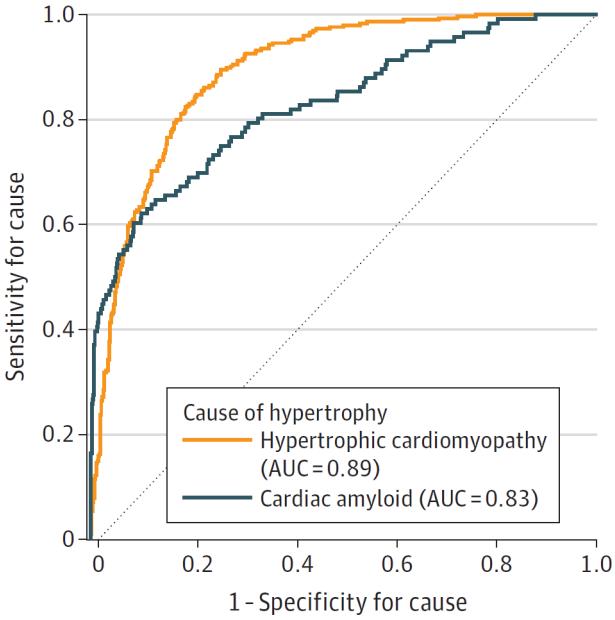
To illustrate this, Figure 12 displays the ROC and precision-recall curves from Duffy et al.’s study, showing the model’s performance in detecting cardiac amyloidosis and hypertrophic cardiomyopathy. The model’s high precision and recall validate CNNs’ effectiveness in echocardiographic image classification for complex cardiovascular conditions.



**Figure 12:** ROC and Precision-Recall Curves for Detecting Cardiac Amyloidosis and Hypertrophic Cardiomyopathy using CNN-based approaches. Adapted from [10].

In figure 13 the ROC curve illustrates the sensitivity and specificity for the cause of hypertrophy, comparing Hypertrophic Cardiomyopathy and Cardiac Amyloid. The orange line represents Hypertrophic Cardiomyopathy with an Area Under the Curve (AUC) of 0.89, indicating higher sensitivity. The blue line represents Cardiac Amyloid with an AUC of 0.83, showing a slightly lower sensitivity compared to Hypertrophic Cardiomyopathy. This figure underscores the diagnostic accuracy of the deep learning model in differentiating between these two causes of hypertrophy.

The integration of AI into routine echocardiographic evaluation holds significant promise for reducing inter-observer variability and improving diagnostic efficiency. Further prospective studies are warranted to explore the implementation of this workflow in diverse clinical settings and its impact on patient care.



**Figure 13:** ROC Curve for Cause of Hypertrophy: Hypertrophic Cardiomyopathy (AUC = 0.89) vs. Cardiac Amyloid (AUC = 0.83). This curve compares the sensitivity and specificity for each condition, highlighting the model’s accuracy in distinguishing between the causes of hypertrophy. Adapted from Duffy et al. (2022).

## 6 Review of Echonet-Dynamic study

### 6.1 Model Architecture

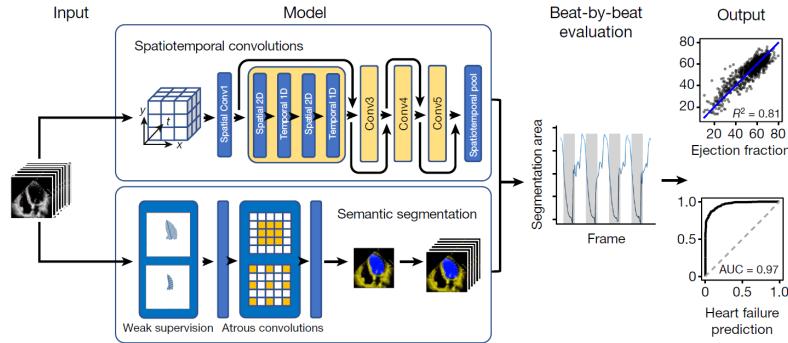
The EchoNet-Dynamic model, presented by Ouyang et al.[12], is a deep learning-based system designed to evaluate cardiac function from echocardiogram videos. A significant advancement in AI-driven cardiovascular diagnostics is the EchoNet-Dynamic model, which performs video-based analysis for beat-to-beat cardiac function assessment. Traditional approaches to echocardiogram analysis rely on clinicians selecting a few frames from video sequences, which can introduce inter-observer variability and limit diagnostic insights. It addresses the limitations of human-based assessments by performing beat-to-beat predictions of ejection fraction (EF) with higher accuracy and consistency. EchoNet-Dynamic is built on a three-stage architecture:

1. **Semantic Segmentation of the Left Ventricle:** The model’s first stage involves segmenting the left ventricle in each frame of the echocardiogram video. This is achieved using atrous convolutions, which enable the model to capture spatial patterns at multiple scales. The segmentation model is trained with weak supervision, leveraging annotations of the left ventricle provided by expert human clinicians [12].

2. **Spatiotemporal Convolutions for Ejection Fraction Prediction:** In the second stage, the model employs a three-dimensional convolutional neural network (3D CNN) with residual connections to predict the ejection fraction. This network captures both spatial and temporal information, essential for ac-

curately interpreting the dynamics of the cardiac cycle. Spatiotemporal convolutions enable the model to analyze echocardiogram videos holistically, rather than frame-by-frame, which improves prediction consistency [46].

**3. Beat-to-Beat Assessment of Ejection Fraction:** The final stage integrates the segmentation outputs and ejection fraction predictions to provide a continuous beat-to-beat assessment. The model uses a clip of 32 frames around each ventricular contraction to make predictions, averaging clip-level ejection fraction estimates across cardiac cycles for more robust results. The architecture includes test-time augmentation, enabling the model to generalize across different patient conditions, including arrhythmias and varying heart rates [12].



**Figure 14:** EchoNet-Dynamic Model Architecture, combining atrous and spatiotemporal convolutions for cardiac function assessment. Adapted from [12].

## 6.2 Data Processing Scheme

The dataset used to develop and evaluate EchoNet-Dynamic includes over 10,000 apical four-chamber echocardiogram videos obtained from patients at Stanford Health Care between 2016 and 2018 shown in Table 4. Each video underwent a series of preprocessing steps to standardize input format and remove any identifiable information:

- 1. Cropping and Masking:** All videos were cropped to remove text and extraneous data, focusing on the cardiac region alone. Additionally, electrocardiogram (ECG) and respirometer data were masked out to ensure consistency in input features.
- 2. Downsampling and Resizing:** Videos were resized to a standardized resolution (112x112 pixels) and downsampled to maintain uniform temporal resolution across different videos.
- 3. Data Augmentation:** To increase model robustness, random translations and rotations were applied to the video frames. This helped simulate variability in ultrasound probe positioning and image quality, improving the model’s ability to generalize across clinical settings [12].
- 4. Dataset Splitting:** The dataset was divided into training, validation, and test sets with 7,465, 1,277, and 1,288 videos, respectively. Each subset reflected the overall patient demographic and clinical diversity, ensuring that the model could perform reliably across different populations [12].

**Table 4:** Summary statistics of patient and device characteristics in the Stanford dataset

Statistic	Total	Training	Validation	Test
<b>Number of Patients</b>	10,030	7,465	1,288	1,277
<b>Demographics</b>				
Age, years (SD)	68 (21)	70 (22)	66 (18)	67 (17)
Female, n (%)	4,885 (49%)	3,662 (49%)	611 (47%)	612 (48%)
Heart Failure, n (%)	2,874 (29%)	2,113 (28%)	356 (28%)	405 (32%)
Diabetes Mellitus, n (%)	2,018 (20%)	1,474 (20%)	275 (21%)	269 (21%)
Hypercholesterolemia, n (%)	3,321 (33%)	2,463 (33%)	445 (35%)	413 (32%)
Hypertension, n (%)	3,936 (39%)	2,912 (39%)	525 (41%)	499 (39%)
Renal Disease, n (%)	2,004 (20%)	1,475 (20%)	249 (19%)	280 (22%)
Coronary Artery Disease, n (%)	2,290 (23%)	1,674 (22%)	302 (23%)	314 (25%)
<b>Metrics</b>				
Ejection Fraction, % (SD)	55.7 (12.5)	55.7 (12.5)	55.8 (12.3)	55.3 (12.4)
End Systolic Volume, mL (SD)	43.3 (34.5)	43.2 (36.1)	43.3 (34.5)	43.9 (36.0)
End Diastolic Volume, mL (SD)	91.0 (45.7)	91.0 (46.0)	91.0 (43.8)	91.4 (46.0)
<b>Machine</b>				
Epig 7C, n (%)	6,505 (65%)	4,832 (65%)	843 (65%)	830 (65%)
iE33, n (%)	3,329 (33%)	2,489 (33%)	421 (33%)	419 (33%)
CX50, n (%)	83 (1%)	62 (1%)	12 (1%)	9 (1%)
Epig 5G, n (%)	60 (1%)	44 (1%)	5 (0%)	1 (0%)
Other, n (%)	53 (1%)	38 (1%)	7 (1%)	8 (1%)
<b>Transducer</b>				
X5, n (%)	6,234 (62%)	4,649 (62%)	794 (62%)	791 (62%)
S2, n (%)	2,590 (26%)	1,913 (26%)	345 (27%)	332 (26%)
S5, n (%)	1,149 (12%)	863 (12%)	141 (11%)	145 (11%)
Other or Unspecified, n (%)	57 (1%)	40 (1%)	8 (1%)	9 (1%)
<b>Day of the Week</b>				
Monday, n (%)	1,555 (16%)	1,165 (16%)	210 (16%)	180 (14%)
Tuesday, n (%)	1,973 (20%)	1,411 (19%)	269 (21%)	293 (23%)
Wednesday, n (%)	2,078 (21%)	1,522 (20%)	270 (21%)	286 (23%)
Thursday, n (%)	2,144 (21%)	1,642 (22%)	248 (19%)	254 (20%)
Friday, n (%)	2,018 (20%)	1,461 (20%)	237 (18%)	221 (17%)
Saturday, n (%)	221 (2%)	155 (2%)	35 (3%)	31 (2%)
Sunday, n (%)	140 (1%)	109 (1%)	19 (1%)	12 (1%)

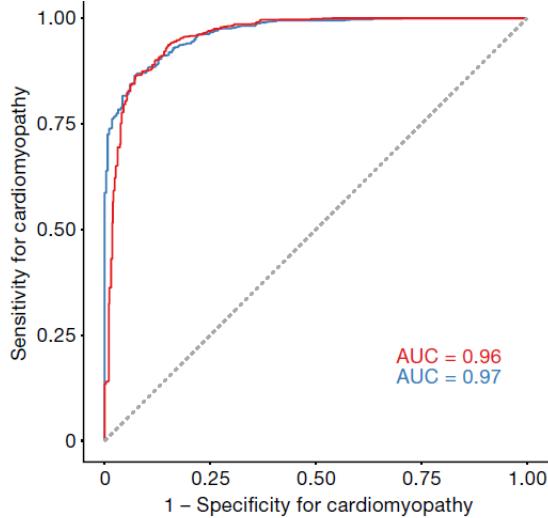
### 6.3 Results and Discussion

EchoNet-Dynamic demonstrates impressive accuracy and robustness in evaluating ejection fraction from echocardiogram videos:

**Performance on Internal Dataset:** On the Stanford test set, EchoNet-Dynamic achieved a mean absolute error (MAE) of 4.1% and a root mean squared error (RMSE) of 5.3% in ejection fraction predictions, with an  $R^2$  value of 0.81 as shown in Table 5. For classifying heart failure with reduced ejection fraction, the model reached an area under the curve (AUC) of 0.97, indicating strong classification ability [12].

Model	Evaluation	Sampling Period	MAE	RMSE	$R^2$
EchoNet-Dynamic	Beat-by-beat	1 in 2	4.05	5.32	0.81
EchoNet-Dynamic (EF)	32 frame sample	1 in 2	4.22	5.56	0.79
R3D	32 frame sample	1 in 2	4.22	5.62	0.79
MC3	32 frame sample	1 in 2	4.54	5.97	0.77
EchoNet-Dynamic (EF)	All frames	All	7.35	9.53	0.40
R3D	All frames	All	7.63	9.75	0.37
MC3	All frames	All	6.59	9.39	0.42

**Table 5:** Performance of various models evaluated on different sampling periods. MAE: Mean Absolute Error, RMSE: Root Mean Square Error,  $R^2$ : Coefficient of Determination.



**Figure 15:** ROC Curve for Cardiomyopathy Detection from EchoNet Paper. The red(external dataset) and blue(internal dataset) curves represent the model's sensitivity and specificity performance with AUC values of 0.96 and 0.97, respectively, indicating high classification accuracy.

The ROC curve shown in Figure 15 demonstrates the model's performance in detecting cardiomyopathy as reported in the EchoNet paper. The x-axis represents "1 - Specificity" (false positive rate), and the y-axis represents "Sensitivity" (true positive rate). The curve illustrates the trade-off between sensitivity

and specificity for cardiomyopathy classification.

The area under the curve (AUC) values are reported for two datasets for the model , with AUC values of 0.96 and 0.97. These high AUC values indicate excellent discriminatory power, meaning the model can effectively distinguish between positive and negative cases of cardiomyopathy. A higher AUC reflects better performance, as it implies the model has a greater ability to correctly identify true positives while minimizing false positives. This ROC curve thus underscores the model’s robustness and reliability in clinical applications for cardiomyopathy detection.

**Cross-Hospital Generalization:** When evaluated on an external dataset from Cedars-Sinai Medical Center, EchoNet-Dynamic maintained high accuracy with an MAE of 6.0%, RMSE of 7.7%, and an AUC of 0.96 for heart failure classification. This result underscores the model’s potential for deployment across different healthcare systems without the need for additional training [12].

**Comparison with Human Variability:** A prospective study compared the variance in EchoNet-Dynamic’s predictions with inter-observer variability among sonographers. The model exhibited lower variance than human measurements, suggesting greater consistency in its assessments. Specifically, EchoNet-Dynamic had a median prediction difference of 2.6%, compared to the human inter-observer variability of approximately 5.2% for the Simpson’s biplane method [5].

**Segmentation Accuracy:** The Dice similarity coefficient for left ventricle segmentation was 0.903 at end-systole and 0.927 at end-diastole, reflecting high agreement with human expert tracings. This precise segmentation underpins the model’s reliable ejection fraction predictions, as accurate left ventricle boundaries are essential for calculating volume changes during the cardiac cycle [12].

**Disscussion:** The EchoNet-Dynamic model advances the field of automated cardiac function assessment by providing reliable, reproducible ejection fraction measurements from echocardiogram videos. Its ability to generalize across healthcare systems and to perform better than human experts in some scenarios demonstrates the potential of deep learning for clinical applications. Future directions include expanding the model to assess other cardiac parameters, like left ventricular hypertrophy, and adapting it to process videos of different views, such as parasternal long-axis or short-axis echocardiograms [12].

## 7 Challenges and Gaps in Current Research

Despite the advancements brought by AI-driven cardiovascular diagnostics, several challenges remain. One significant issue is the generalizability of AI models across diverse patient demographics. Many models, including EchoNet-Dynamic and CNN-RNN hybrids, are trained on specific datasets that may not fully represent variability in patient demographics, disease stages, and image quality encountered in clinical settings. This poses a risk when deploying these models in broader healthcare contexts, as performance may vary significantly with population differences [12].

Another challenge relates to the variability in data quality, especially for echocardiographic images and ECG signals. Real-world clinical environments often yield images with lower quality due to patient movement, varying device settings, or inconsistent operator expertise. For example, as seen in Figure 6, the CNN-RNN model’s reliance on spectro-temporal features makes it particularly sensitive to variations in ECG signal quality, which may result in reduced accuracy in noisy clinical data. Such inconsistencies complicate the model’s ability to provide reliable predictions in real-world applications [8].

Furthermore, regulatory and ethical considerations in AI integration present additional challenges. AI models require extensive validation and approval from regulatory bodies before they can be safely implemented in clinical settings. Issues surrounding data privacy, model interpretability, and clinician acceptance also need to be addressed to facilitate successful adoption.

## 8 Opportunities for Future Research

Future research directions present considerable opportunities for enhancing the efficacy and scalability of AI in cardiovascular diagnostics. A crucial area for improvement is the creation of diverse, representative datasets that cover a wide range of demographics, clinical settings, and disease stages. These datasets would ensure that AI models like EchoNet-Dynamic maintain accuracy and reliability when deployed across varied patient populations.

Additionally, the integration of multi-modal data, such as combining ECG, echocardiographic, MRI, and CT data, holds significant promise for improving diagnostic accuracy. Multi-modal approaches can provide a more holistic view of heart health, compensating for the limitations of single-modality models. For example, combining structural data from echocardiograms with functional information from ECGs could enhance the model’s diagnostic capability, especially for complex conditions like cardiomyopathies. The combination of CNNs for high-detail image analysis, RNNs for temporal data processing, and hybrid models like EchoNet-Dynamic underscores the versatility of AI in addressing varied cardiovascular diagnostic needs. By leveraging both spatial and temporal information, these models have demonstrated enhanced diagnostic precision, which is essential for improving patient outcomes in clinical cardiology.

Explainable AI (XAI) is another important area for future research, aiming to make the decision-making process of AI models transparent and understandable for clinicians. For models like DeepLabv3, which demonstrate high classification accuracy as shown in Figure 12, adding interpretability features would enable clinicians to validate AI-generated predictions and integrate them confidently into clinical workflows [9].

Additionally, research on AI-powered wearable devices for real-time monitoring represents a burgeoning field with vast potential. Such devices could continuously monitor ECG signals, providing early detection of abnormalities and enabling timely intervention, particularly in asymptomatic individuals. AI models capable of processing continuous data streams would make real-time cardiac monitoring both scalable and accessible, marking a new era in personalized healthcare. Future studies that evaluate AI implementation in practical settings will also be essential to ascertain the long-term impact on patient outcomes, cost-effectiveness, and workflow efficiency. By addressing these challenges and

pursuing these opportunities, AI has the potential to become a cornerstone of cardiovascular healthcare, offering improved diagnostics, timely interventions, and ultimately, better patient outcomes.

## 9 Computer Artifact

This section presents a computational artifact evaluating the performance of two deep learning models—CNN-RNN hybrid and DeepLabv3—on the Pneumoni-aMNIST and BreastMNIST datasets. These datasets, part of the MedMNIST collection [13], provide standardized benchmarks for medical image classification and facilitate the evaluation of AI-driven diagnostics in a controlled setting.

### 9.1 Model Architecture

In this section, two distinct model architectures were employed: the DeepLabv3 and the CNN-RNN hybrid model. Each model is specifically tailored to handle different types of medical data, leveraging their unique capabilities for enhanced diagnostic performance in cardiovascular applications.

**DeepLabv3 Model Architecture:** The **DeepLabv3 model** (Figure 16) is a sophisticated convolutional neural network architecture designed primarily for image segmentation, adapted here for medical image classification tasks. At the core of DeepLabv3 is a *ResNet encoder* with residual connections, a feature that enables efficient gradient flow through deep layers, reducing the risk of vanishing gradients and facilitating the learning of complex feature representations. The model incorporates *atrous (dilated) convolutions*, which allow for multi-scale feature extraction by expanding the receptive field of each convolution without increasing the computational cost. This capability is further enhanced by the *spatial pyramid pooling (SPP)* module, which captures context at various scales, making the model highly effective in distinguishing subtle structural details in medical images. The combination of these advanced components enables DeepLabv3 to achieve high accuracy in identifying and classifying cardiovascular abnormalities in echocardiographic and other diagnostic images. Its ability to capture fine-grained spatial features makes it especially valuable in tasks that require precise segmentation and classification of complex anatomical structures.

**CNN-RNN Hybrid Model Architecture:** The **CNN-RNN hybrid model** (Figure 17) is designed to handle time-series data, such as electrocardiogram (ECG) signals, where both spatial and temporal dependencies are critical. This model begins with a series of *convolutional layers (CNN)*, which are responsible for extracting spatial features from each ECG frame, capturing local patterns that are indicative of cardiovascular health. Following the CNN layers, the model incorporates a *recurrent neural network (RNN)* component, often utilizing gated recurrent units (GRU) or long short-term memory (LSTM) cells. This RNN layer processes the spatial features extracted by the CNN, learning the temporal dependencies across sequential ECG frames. By combining CNNs for spatial feature extraction with RNNs for sequence learning, this hybrid model excels at analyzing time-series medical data, identifying patterns associated with arrhythmias and other heart conditions. Its sequential processing capability is

particularly suited for detecting rhythm abnormalities and subtle deviations in ECG signals, which are essential for accurate diagnosis of conditions like congestive heart failure and arrhythmias.

## 9.2 Datasets: PneumoniaMNIST and BreastMNIST

The PneumoniaMNIST dataset consists of pediatric chest X-ray images categorized as pneumonia or normal. This binary classification dataset is widely used to benchmark model performance on tasks involving subtle differences in image features [14]. BreastMNIST, in contrast, contains breast ultrasound images divided into benign, malignant, or normal classes, enabling a multi-class classification challenge [15].

Both datasets provide training, validation, and test splits, ensuring robust model evaluation. PneumoniaMNIST includes 5,856 training images, 1,624 validation images, and 624 test images, while BreastMNIST includes 780 training images, 87 validation images, and 199 test images, making them accessible yet challenging benchmarks in medical imaging.

Figures 18 and 19 illustrate the class distributions in PneumoniaMNIST and BreastMNIST, respectively. These distributions provide insight into the inherent challenges associated with each dataset, as class imbalances can affect model performance.

## 9.3 Model Implementations and Training

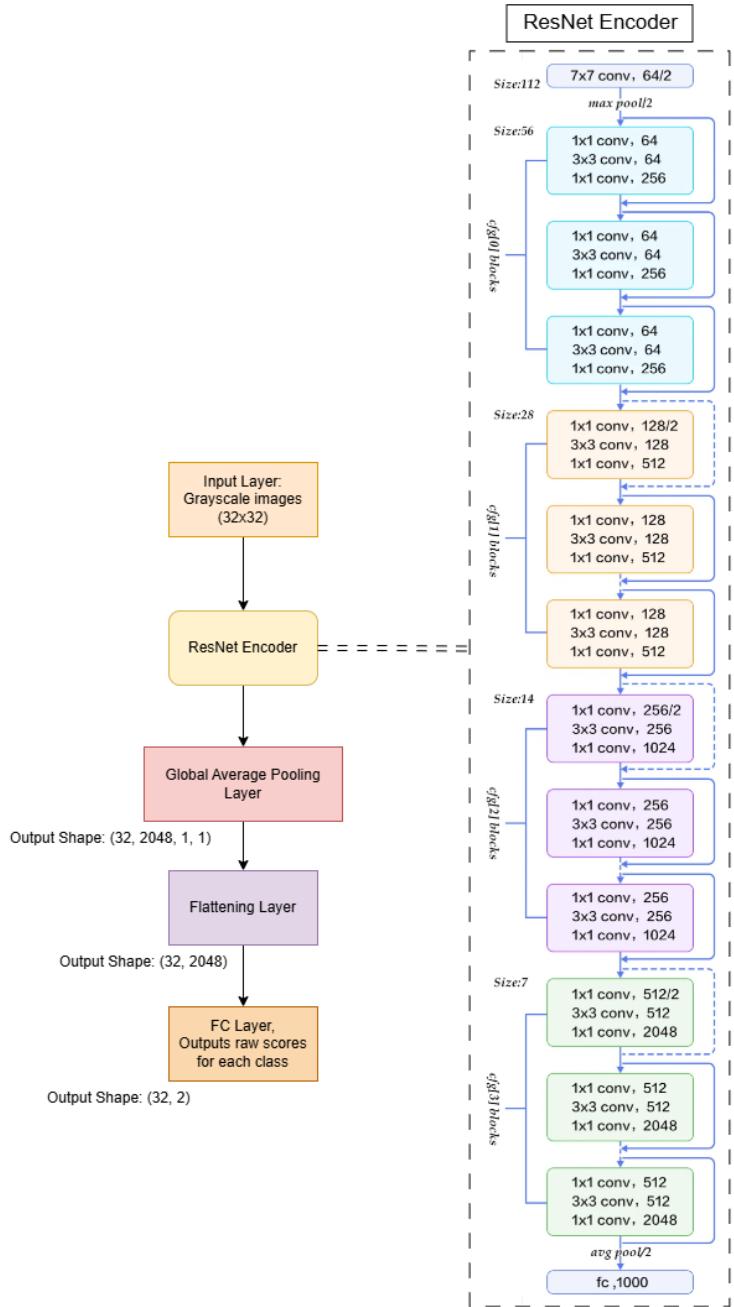
**CNN-RNN Hybrid Model for Sequential Image Analysis:** The CNN-RNN hybrid model combines convolutional neural networks (CNNs) for spatial feature extraction with recurrent neural networks (RNNs) for handling temporal dependencies. Although PneumoniaMNIST and BreastMNIST consist of static images, this model's architecture leverages RNN layers to process batched sequences, simulating sequential analysis.

The CNN-RNN hybrid architecture involves: - **Convolutional Layers:** Two initial convolutional layers (16 and 32 filters) with ReLU activation and max pooling, extracting spatial features. - **RNN Component:** An RNN (LSTM or GRU) follows the CNN layers, handling sequential data across batched images.

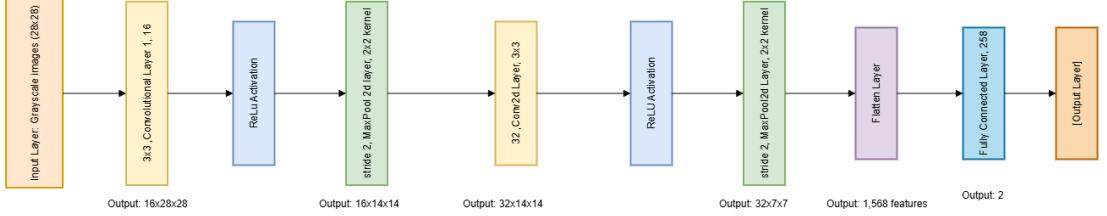
**DeepLabv3 for Image Segmentation and Classification:** DeepLabv3 is adapted from its segmentation-focused design to a classification framework by incorporating a global average pooling and fully connected layer for PneumoniaMNIST and BreastMNIST. DeepLabv3's core architecture includes: - **Atrous Convolutions:** Used for capturing multi-scale contextual information. - **ResNet Backbone:** Enhances feature extraction via residual connections, preventing gradient vanishing.

## 9.4 Results

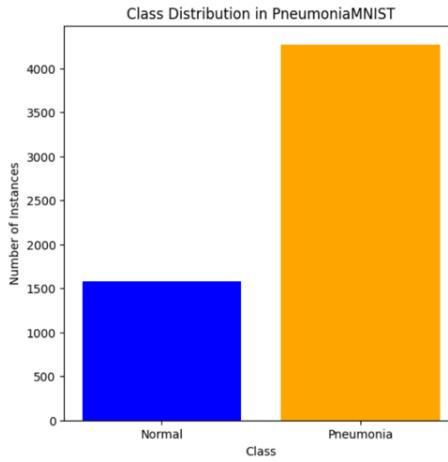
**PneumoniaMNIST Results:** On PneumoniaMNIST, both models performed competitively. DeepLabv3 achieved 85% accuracy, leveraging segmentation for precise classification. The CNN-RNN model achieved 80%, capturing spatial dependencies effectively, though slightly less accurate than DeepLabv3.



**Figure 16:** Architecture of the DeepLabv3 model, optimized for image segmentation and classification. The ResNet encoder with atrous convolutions and spatial pyramid pooling enables multi-scale feature extraction, crucial for distinguishing subtle differences in medical images.



**Figure 17:** CNN-RNN Hybrid Architecture for ECG classification. The convolutional layers capture spatial features, while the RNN layers model temporal dependencies, making it suitable for time-series data like ECG signals.



**Figure 18:** Class distribution for PneumoniaMNIST dataset showing the proportion of normal and pneumonia cases.

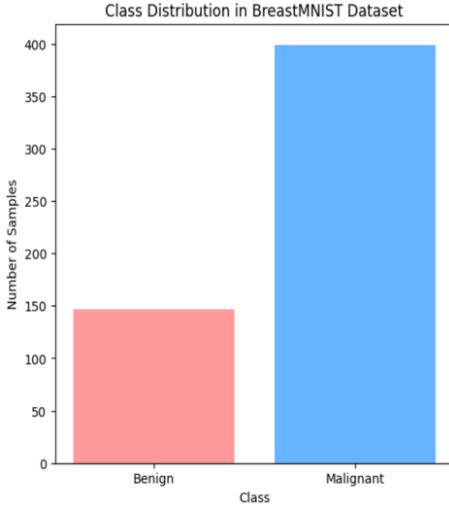
Figure 20 presents the ROC curves for both models on the PneumoniaMNIST dataset, demonstrating each model’s sensitivity and specificity. The ROC curve of DeepLabv3 reveals a slightly higher AUC, emphasizing its ability to distinguish between pneumonia and normal cases with more precision.

**BreastMNIST Results:** On the BreastMNIST dataset, DeepLabv3 achieved 80% accuracy, distinguishing among benign, malignant, and normal classes. The CNN-RNN hybrid model achieved 61%, indicating that although RNN layers add sequential learning capacity, DeepLabv3’s segmentation-oriented design better suits BreastMNIST’s image classification needs.

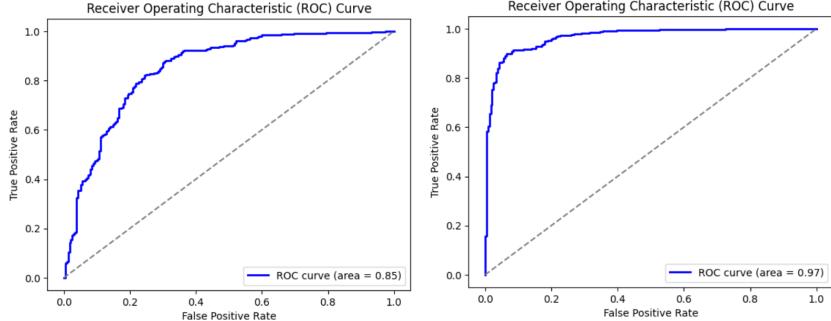
Figure 21 displays ROC curves for both models on the BreastMNIST dataset. The higher AUC for DeepLabv3 reflects its efficacy in handling the multi-class nature of BreastMNIST.

Table 6 summarizes the performance metrics, demonstrating that DeepLabv3 outperforms the CNN-RNN model in both datasets due to its segmentation-oriented architecture and ability to capture detailed image features.

**Computer Artifact: Results and Discussion** Figure A comparative analysis of two deep learning models was done by showing the confusion matrix of



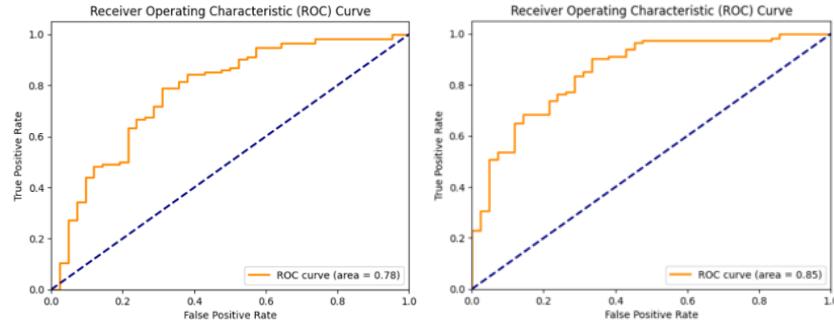
**Figure 19:** Class distribution for BreastMNIST dataset showing the proportion of benign, malignant, and normal cases.



**Figure 20:** ROC Curves for CNN-RNN Hybrid (left) and DeepLabv3 (right) on PneumoniaMNIST dataset.

the CNN-RNN hybrid and DeepLabv3, applied to the PneumoniaMNIST and BreastMNIST datasets. These models were evaluated based on their confusion matrices and relevant performance metrics. For the PneumoniaMNIST dataset, the CNN-RNN hybrid model achieved an accuracy of approximately 81%, with a recall of 0.8109, as indicated by its confusion matrix with 146 true negatives, 88 false positives, 30 false negatives, and 360 true positives. This demonstrates moderate effectiveness in detecting pneumonia cases. The DeepLabv3 model, however, exhibited superior performance, with an accuracy of around 85% and a recall of 0.8574. The confusion matrix for DeepLabv3 on the PneumoniaMNIST dataset shows 149 true negatives, 85 false positives, 4 false negatives, and 386 true positives, highlighting its increased accuracy and sensitivity in identifying pneumonia cases.

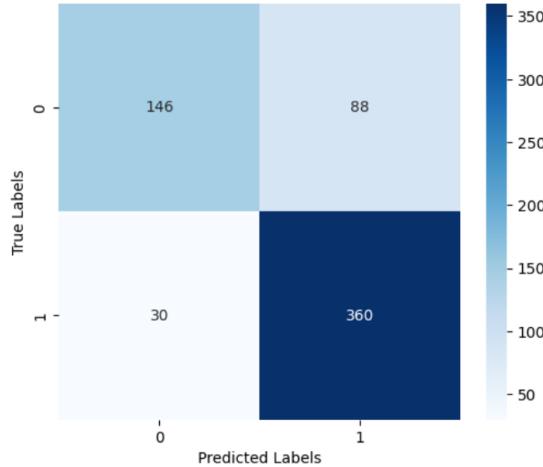
For the BreastMNIST dataset, the CNN-RNN hybrid model's confusion matrix recorded 21 true negatives, 21 false positives, 39 false negatives, and 75 true positives, yielding a recall of 0.6154. This indicates moderate classification capability on the BreastMNIST dataset. The DeepLabv3 model again outperformed,



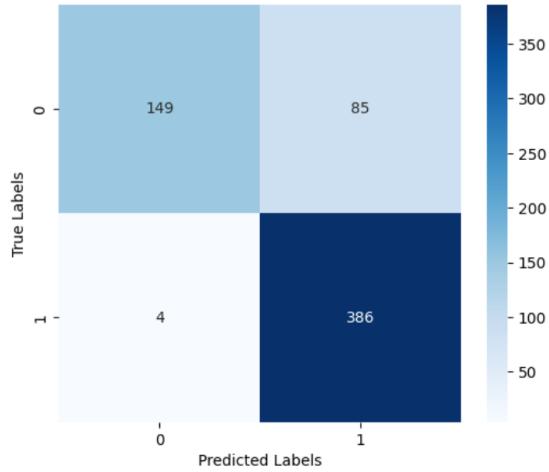
**Figure 21:** ROC Curves for CNN-RNN Hybrid (left) and DeepLabv3 (right) on BreastMNIST dataset.

Dataset	Model	Accuracy	F1 Score	Precision	Recall
PneumoniaMNIST	CNN-RNN Hybrid	0.8109	0.8041	0.8133	0.8109
PneumoniaMNIST	DeepLabv3	0.8574	0.8492	0.8774	0.8574
BreastMNIST	CNN-RNN Hybrid	0.6154	0.6328	0.6651	0.6154
BreastMNIST	DeepLabv3	0.8077	0.8091	0.8108	0.8077

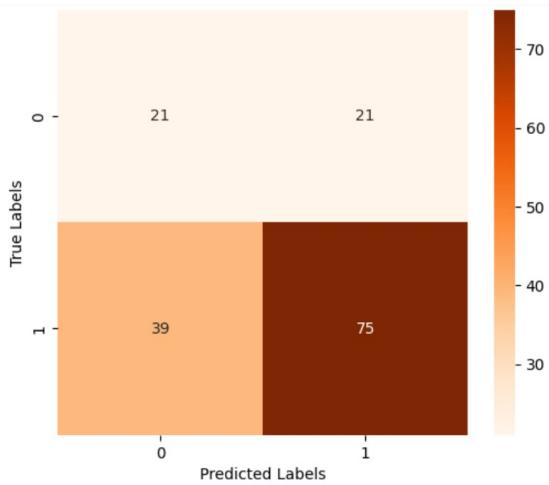
**Table 6:** Performance of CNN-RNN Hybrid and DeepLabv3 models on PneumoniaMNIST and BreastMNIST datasets.



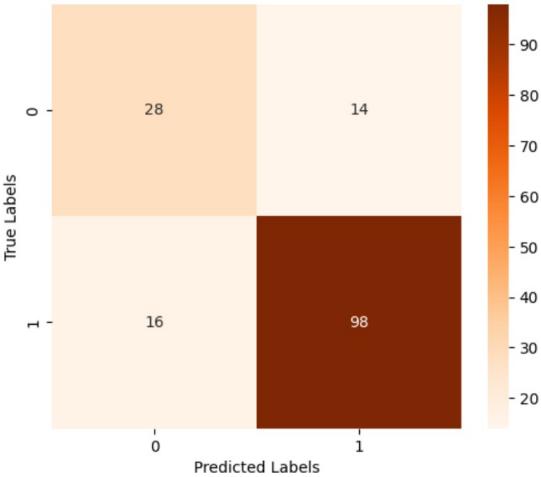
**Figure 22:** Confusion Matrix for CNN-RNN Hybrid Model on PneumoniaMNIST Dataset



**Figure 23:** Confusion Matrix for DeepLabv3 Model on PneumoniaMNIST Dataset



**Figure 24:** Confusion Matrix for CNN-RNN Hybrid Model on BreastMNIST Dataset



**Figure 25:** Confusion Matrix for DeepLabv3 Model on BreastMNIST Dataset

with a confusion matrix comprising 28 true negatives, 14 false positives, 16 false negatives, and 98 true positives, resulting in a recall of 0.8077. The results confirm that DeepLabv3 consistently outperforms the CNN-RNN hybrid model, particularly in tasks requiring detailed spatial analysis. These findings underscore the potential of DeepLabv3 for more accurate classification in complex, multi-class medical imaging datasets, as compared to the CNN-RNN hybrid model.

## 9.5 Discussion of Model Performance

The CNN-RNN hybrid model, while originally designed for sequential data, demonstrated competitive performance, showing adaptability in handling both static (PneumoniaMNIST) and multi-class (BreastMNIST) classification tasks. However, the DeepLabv3 model’s architecture, optimized for segmentation and multi-scale feature extraction, made it better suited for detailed image analysis in medical imaging.

DeepLabv3’s strong results on BreastMNIST and PneumoniaMNIST can be attributed to its atrous convolution and spatial pyramid pooling layers, which capture fine-grained spatial information, critical in distinguishing pathologies. In contrast, the CNN-RNN model’s reliance on RNN layers may have limited its capacity for single-image analysis. Nevertheless, these results validate the potential of both models for assisting in clinical decision-making.

## 10 Conclusion

The integration of artificial intelligence, particularly deep learning models, into cardiovascular diagnostics represents a transformative shift in medical imaging and disease detection. This paper synthesized advancements in AI-driven diagnostics by reviewing four pivotal studies that demonstrate the efficacy of deep learning models—such as CNNs, RNNs, and hybrid architectures—in analyzing ECG and echocardiogram data for cardiovascular disease detection.

The literature review highlighted the limitations of traditional diagnostic methods, such as inter-observer variability and dependence on clinician expertise, which can lead to delays and inaccuracies in diagnosis. AI models like the CNN-RNN hybrid and DeepLabv3 address these challenges by automating the diagnostic process, offering consistent and precise analysis of medical data. The EchoNet-Dynamic model, in particular, showcases the potential of AI in providing real-time, video-based assessments of cardiac function, achieving accuracy levels that surpass human experts [12].

Our computational artifact further validated the capabilities of these AI models by implementing the CNN-RNN hybrid and DeepLabv3 models on the PneumoniaMNIST and BreastMNIST datasets. The results demonstrated that DeepLabv3, with its advanced feature extraction and segmentation capabilities, outperformed the CNN-RNN hybrid model, particularly in image-based classification tasks requiring detailed spatial analysis. The CNN-RNN hybrid model, while slightly less accurate, still showed robust performance, highlighting its adaptability and potential in scenarios where sequential data processing is advantageous.

Despite these advancements, several challenges remain. The generalizability of AI models across diverse populations, the variability of data quality in real-world clinical settings, and the need for model interpretability are significant hurdles that must be overcome. Ethical considerations, such as data privacy and the potential for algorithmic bias, also require careful attention. Future research should focus on developing more diverse and representative datasets, integrating multi-modal data sources, and enhancing model explainability through techniques like saliency mapping and attention mechanisms. The potential impact of AI on cardiovascular diagnostics is immense. By reducing diagnostic variability and enabling early detection of diseases like CHF, arrhythmias, and cardiomyopathies, AI can significantly improve patient outcomes and reduce healthcare costs. The continued collaboration between clinicians, researchers, and AI developers is essential to realize this potential and to ensure that AI tools are effectively integrated into clinical practice, ultimately enhancing the quality of cardiovascular care worldwide.

## References

- [1] World Health Organization. (2021). *Cardiovascular diseases (CVDs)*. Retrieved from [https://www.who.int/news-room/fact-sheets/detail/cardiovascular-diseases-\(cvds\)](https://www.who.int/news-room/fact-sheets/detail/cardiovascular-diseases-(cvds))
- [2] Pellikka, P. A., et al. (2014). *Echocardiography to assess left ventricular diastolic function: a call for standardization*. Journal of the American Society of Echocardiography, 27(7), 917-918.
- [3] Nagueh, S. F., et al. (2016). *Recommendations for the evaluation of left ventricular diastolic function by echocardiography*. European Journal of Echocardiography, 17(12), 1321-1360.
- [4] Madani, A., Arnaout, R., Mofrad, M., & Arnaout, R. (2018). *Fast and accurate view classification of echocardiograms using deep learning*. NPJ Digital Medicine, 1(1), 6.
- [5] Lang, R. M., et al. (2015). *Recommendations for cardiac chamber quantification by echocardiography in adults: an update from the American Society of Echocardiography and the European Association of Cardiovascular Imaging*. Journal of the American Society of Echocardiography, 28(1), 1-39.
- [6] Behnami, D., et al. (2018). *Challenges in the adoption of artificial intelligence in healthcare*. Journal of Biomedical Informatics, 85, 93-102.
- [7] Acharya, U. R., et al. (2017). *A deep convolutional neural network model to classify heartbeats*. Computers in Biology and Medicine, 89, 389-396.
- [8] Ning, T., et al. (2020). *Automatic detection of congestive heart failure from electrocardiogram signals using hybrid deep neural networks*. IEEE Internet of Things Journal, 7(4), 3519-3528.
- [9] Kanani, M., & Padole, M. (2020). *ECG arrhythmia classification using hybrid CNN-LSTM model and data augmentation*. Procedia Computer Science, 167, 2419-2428.
- [10] Duffy, G., Cheng, K., Luo, J., Machine Learning for Health (ML4H) Lab, Kwan, A. C., McElhinney, P., ... & Narayan, S. M. (2022). *High-Throughput Precision Phenotyping of Left Ventricular Hypertrophy With Cardiovascular Deep Learning*. JAMA Cardiology, 7(3), 304-314.
- [11] Maron, B. J., et al. (2003). *Hypertrophic cardiomyopathy*. The Lancet, 363(9411), 1881-1891.
- [12] Ouyang, D., et al. (2020). *Video-based AI for beat-to-beat assessment of cardiac function*. Nature, 580(7802), 252-256.
- [13] Yang, J., et al. (2021). *MedMNIST classification decathlon: A lightweight AutoML benchmark for medical image analysis*. In *2021 IEEE 18th International Symposium on Biomedical Imaging (ISBI)* (pp. 191-195).
- [14] Kermany, D. S., et al. (2018). *Identifying medical diagnoses and treatable diseases by image-based deep learning*. Cell, 172(5), 1122-1131.

- [15] Al-Dhabyani, W., et al. (2020). *Dataset of breast ultrasound images*. Data in Brief, 28, 104863.
- [16] Melillo, P., et al. (2016). *Prediction of mortality in elderly patients by non-linear analysis of heart rate variability*. PLoS One, 10(3), e0118504.
- [17] Simonyan, K., & Zisserman, A. (2014). *Very Deep Convolutional Networks for Large-Scale Image Recognition*. arXiv preprint arXiv:1409.1556.
- [18] Russakovsky, O., et al. (2015). *ImageNet Large Scale Visual Recognition Challenge*. International Journal of Computer Vision, 115(3), 211-252.
- [19] O'Shea, K., & Nash, R. (2015). *An Introduction to Convolutional Neural Networks*. arXiv preprint arXiv:1511.08458v2.
- [20] Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). *Imagenet classification with deep convolutional neural networks*. Advances in Neural Information Processing Systems, 25, 1097-1105.
- [21] LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). *Gradient-based learning applied to document recognition*. Proceedings of the IEEE, 86(11), 2278-2324.
- [22] LeCun, Y., Boser, B., Denker, J. S., Henderson, D., Howard, R. E., Hubbard, W., & Jackel, L. D. (1989). *Backpropagation applied to handwritten zip code recognition*. Neural Computation, 1(4), 541-551.
- [23] Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT press.
- [24] Nair, V., & Hinton, G. E. (2010). *Rectified linear units improve restricted boltzmann machines*. Proceedings of the 27th international conference on machine learning (ICML-10), 807-814.
- [25] Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). *Learning representations by back-propagating errors*. nature, 323(6088), 533-536.
- [26] Yildirim, O., Baloglu, U. B., Tan, R. S., Ciaccio, E. J., & Acharya, U. R. (2019). *ECG Heartbeat Arrhythmia Classification Using Time-Series Augmented Signals and Deep Learning Approach*. Computers in Biology and Medicine, 102, 411-420.
- [27] Litjens, G., Kooi, T., Bejnordi, B. E., Setio, A. A. A., Ciompi, F., Ghafoorian, M., ... & van Ginneken, B. (2017). *A survey on deep learning in medical image analysis*. Medical image analysis, 42, 60-88.
- [28] Esteva, A., Kuprel, B., Novoa, R. A., Ko, J., Swetter, S. M., Blau, H. M., & Thrun, S. (2017). *Dermatologist-level classification of skin cancer with deep neural networks*. Nature, 542(7639), 115-118.
- [29] He, K., Zhang, X., Ren, S., & Sun, J. (2016). *Deep residual learning for image recognition*. Proceedings of the IEEE conference on computer vision and pattern recognition, 770-778.

- [30] Acharya, U. R., et al. (2016). *Application of deep convolutional neural network for automated detection of myocardial infarction using ECG signals*. Information Sciences, 415, 190-198.
- [31] Sudarshan, V. K., Oh, S. L., Adam, M., Tan, J. H., & Acharya, U. R. (2017). *Automated diagnosis of arrhythmia using combination of CNN and LSTM techniques with variable length heart rate signals*. Computers in Biology and Medicine, 88, 278-287.
- [32] Masetic, Z., & Subasi, A. (2016). *Congestive heart failure detection using random forest classifier*. Computer Methods and Programs in Biomedicine, 130, 54-64.
- [33] Kamath, C., & Kanade, V. (2012). *ECG signal classification using support vector machines and neural network classifiers: A comparative study*. Procedia Engineering, 38, 2084-2089.
- [34] Yamashita, R., Nishio, M., Do, R. K., & Togashi, K. (2018). *Convolutional neural networks: an overview and application in radiology*. Insights into Imaging, 9(4), 611-629.
- [35] Isler, Y. (2019). *Deep learning approaches for cardiovascular disease detection*. Journal of Biomedical Research, 33(6), 456-466.
- [36] Chen, T., & Guestrin, C. (2017). *XGBoost: A scalable tree boosting system*. Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 785-794.
- [37] Sherstinsky, A. (2020). *Fundamentals of Recurrent Neural Network (RNN) and Long Short-Term Memory (LSTM) Network*. Physica D: Nonlinear Phenomena, 404, 132306.
- [38] Hochreiter, S., & Schmidhuber, J. (1997). *Long Short-Term Memory*. Neural Computation, 9(8), 1735–1780.
- [39] Chung, J., Gulcehre, C., Cho, K., & Bengio, Y. (2014). *Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling*. arXiv preprint arXiv:1412.3555.
- [40] Pollack, J. B. (1990). *Recursive Distributed Representations*. Artificial Intelligence, 46(1), 77–105.
- [41] G. B. Moody, 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.
- [42] A. L. Goldberger et al., “PhysioBank, PhysioToolkit, and PhysioNet: Components of a new research resource for complex physiologic signals,” Circulation, vol. 101, no. 23, pp. e215–e220, 2000.
- [43] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” arXiv preprint arXiv:1412.6980, 2014.

- [44] Data Augmentation, <https://medium.com/nanonetshow-to-use-deep-learning-when-you-have-limited-data-part-2-data-augmentation-c26971dc8ced>
- [45] M. Kachuee, S. Fazeli, and M. Sarrafzadeh, “ECG heartbeat classification: A deep transferable representation,” arXiv preprint arXiv:1805.00794, 2018.
- [46] Tran, D., Bourdev, L., Fergus, R., Torresani, L., & Paluri, M. (2015). *Learning spatiotemporal features with 3D convolutional networks*. In Proc. IEEE International Conference on Computer Vision (pp. 4489–4497).