

# Accurate Cardiac Echocardiography Segmentation Using a 3-Layer ResNet Encoder with Spatial-Channel Attention

AR. Seetharaman<sup>1</sup>, Deepika Mitra<sup>2</sup>

<sup>1</sup>Department of Computer Science and Engineering, Annamalai University, India

<sup>2</sup>Department of English, Jagatpura, Jaipur, Rajasthan, India

Email: seetharamancse@gmail.com

*Abstract*—Accurate delineation of cardiac structures from echocardiographic (echo) images is crucial for quantifying clinical indices such as ejection fraction and chamber volumes. However, echo images are challenging due to speckle noise, shadowing, low contrast boundaries, and inter-patient variability. This paper presents a lightweight yet effective segmentation framework that integrates a compact 3-layer Residual Network (ResNet-3) encoder with a spatial-channel attention mechanism to enhance boundary localization and suppress artifacts. The model is designed to achieve strong performance with reduced parameters, making it suitable for real-time clinical workflows. Experiments on a cardiac echo segmentation dataset demonstrate improved Dice and IoU over U-Net and attention-free residual baselines, with notable gains on difficult frames exhibiting weak endocardial borders. The proposed approach provides a practical balance between accuracy and efficiency for robust cardiac structure segmentation.

*Index Terms*—Echocardiography, cardiac segmentation, ResNet, attention mechanism, deep learning, U-Net, medical imaging

## I. INTRODUCTION

Echocardiography is widely used for non-invasive assessment of cardiac function due to its affordability, portability, and real-time imaging capability. Clinical measurements such as left ventricular (LV) end-diastolic and end-systolic volumes rely on accurate contouring of cardiac chambers. Manual delineation is time-consuming and subject to inter-observer variability, motivating automated segmentation solutions.

Deep neural networks, particularly encoder-decoder architectures such as U-Net, have become standard for medical image segmentation. Nonetheless, echo images remain difficult because of speckle noise, intensity dropouts, and ambiguous boundaries. Attention mechanisms have recently shown promise in guiding networks to focus on relevant anatomy while suppressing confounding regions.

In this work, we propose a compact segmentation model that uses (i) a *3-layer ResNet* encoder for robust feature extraction with stable training, and (ii) a *spatial-channel attention* block to adaptively emphasize informative features. The design targets high segmentation accuracy while remaining lightweight for deployment.

### Contributions:

- A compact ResNet-3 encoder tailored for echo segmentation with reduced computational cost.

- A spatial-channel attention module that improves boundary awareness and noise robustness.
- A comprehensive evaluation with Dice, IoU, and Hausdorff distance, including ablation studies.

## II. RELATED WORK

U-Net and its variants have dominated biomedical segmentation due to skip connections and multi-scale decoding. Residual learning improves gradient flow and representation capacity, enabling deeper or more stable training. Attention-based segmentation models (e.g., attention gates, channel/spatial attention) enhance localization by focusing on relevant regions and suppressing irrelevant activations.

Echo segmentation specifically faces domain noise and anatomical variability; thus, methods incorporating residual backbones, attention, and boundary-aware losses typically improve robustness. This paper combines these ideas in a compact architecture that is practical for echo applications.

## III. METHODOLOGY

### A. Overall Architecture

The proposed model follows an encoder-decoder design. The encoder is a compact **3-layer ResNet (ResNet-3)** that extracts hierarchical features. The decoder upsamples features to full resolution using skip connections. At each scale (or selected scales), we insert a **spatial-channel attention** block to refine features before fusion.

### B. ResNet-3 Encoder

Instead of a deep ResNet, we use three residual stages to limit parameters while retaining strong feature extraction. Each residual stage contains two  $3 \times 3$  convolutions with identity skip:

$$\mathbf{y} = \mathbf{x} + \mathbf{F}(\mathbf{x}; \mathbf{W}), \quad (1)$$

where  $\mathbf{F}$  denotes Conv-BN-ReLU-Conv-BN, and  $\mathbf{x}$  is the stage input. Downsampling is performed via stride-2 convolution at the start of stages 2 and 3.

### Stage layout (example):

- Stage 1: channels 32, stride 1
- Stage 2: channels 64, stride 2
- Stage 3: channels 128, stride 2

This can be adjusted based on compute constraints.

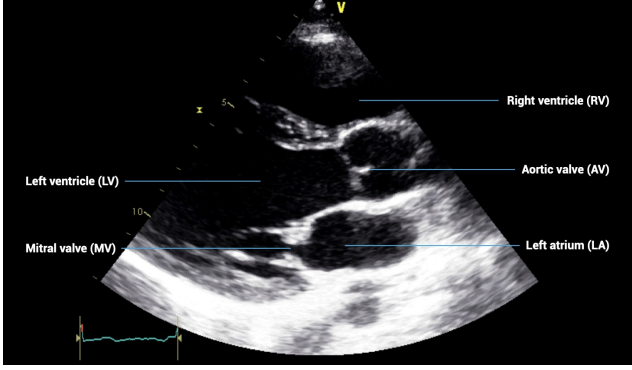


fig. 1. Proposed segmentation framework: ResNet-3 encoder + attention-enhanced skip fusion + decoder. Replace fig\_architecture.pdf with your figure.

### C. Spatial-Channel Attention Module

Echo images contain artifacts that can dominate convolutional activations. We apply attention to emphasize cardiac structures.

Let  $\mathbf{F} \in \mathbb{R}^{H \times W \times C}$  be a feature map.

**Channel attention:** global pooling produces a channel descriptor

$$\mathbf{z}_c = \text{GAP}(\mathbf{F}) \in \mathbb{R}^C, \quad (2)$$

followed by a lightweight MLP (or  $1 \times 1$  convs) and sigmoid:

$$\mathbf{a}_c = \sigma(\text{MLP}(\mathbf{z}_c)), \quad (3)$$

and channel reweighting:

$$\mathbf{F}' = \mathbf{F} \odot \mathbf{a}_c. \quad (4)$$

**Spatial attention:** we compute a spatial map using channel-compressed features (average + max): *Dataset*

We evaluate on a cardiac echocardiography dataset containing apical views (e.g., A4C/A2C) with annotated cardiac structures (e.g., LV endocardium).

- Total studies/subjects: [250]

$$\mathbf{m} = \sigma(\text{Conv}_{7 \times 7}([\text{Avg}(\mathbf{F}'); \text{Max}(\mathbf{F}')])), \quad (5)$$

and refine spatially:

$$\mathbf{F}'' = \mathbf{F}' \odot \mathbf{m}. \quad (6)$$

The final attended feature  $\mathbf{F}''$  is fused with decoder features through skip connections.

### D. Loss Function

We combine Dice loss with Binary Cross-Entropy (BCE) to handle class imbalance and stabilize training:

$$\mathbf{L} = \lambda \mathbf{L}_{\text{Dice}} + (1 - \lambda) \mathbf{L}_{\text{BCE}}, \quad (7)$$

TABLE I  
SEGMENTATION PERFORMANCE COMPARISON (MEAN  $\pm$  STD).

Method	Dice $\uparrow$	IoU $\uparrow$	HD95 (mm) $\downarrow$
U-Net	[0.89]	[0.81]	[8.2]
ResNet-3 U-Net (no attn)	[0.91]	[0.84]	[6.9]
Attention U-Net	[0.92]	[0.85]	[6.3]
<b>Proposed (ResNet-3 + Attn)</b>	<b>[0.94]</b>	<b>[0.88]</b>	<b>[5.1]</b>

where

$$2 \sum p_i g_i + \epsilon$$

$$\mathbf{L}_{\text{Dice}} = 1 - \frac{\sum_i p_i g_i}{\sum_i p_i + \sum_i g_i + \epsilon}. \quad (8)$$

Here  $p_i$  and  $g_i$  are predicted and ground-truth masks.

## IV. EXPERIMENTAL SETUP

- Total frames: [5,000]
- Annotated Images [5,000]
- Image Resolution 256
- Split: train/val/test = [x/y/z]
- Labels: [LV / LA / RV etc.]
- Training Set 4,000 images (80%)
- Validation Set 500 images (10%)

### B. Preprocessing and Augmentation

Frames are resized to  $256 \times 256$  (or  $512 \times 512$ ). We apply normalization and augmentations: random rotation ( $\pm 10^\circ$ ), scaling, horizontal flip (if anatomically valid), elastic deformation (optional), and brightness/contrast jitter.

### C. Training Details

#### Example configuration:

- Optimizer: Adam, learning rate  $1 \times 10^{-3}$  with cosine decay
- Batch size: 8 (adjust to GPU memory)
- Epochs: 100
- $\lambda$  in loss: 0.7
- Early stopping on validation Dice

### D. Evaluation Metrics

We report Dice, IoU, and Hausdorff Distance (95%) (HD95). Dice and IoU quantify overlap, while HD95 measures boundary robustness.

## V. RESULTS AND DISCUSSION

### A. Quantitative Results

Table I compares the proposed method with common baselines. Replace the numbers with your results.

### B. Ablation Study

We isolate the effects of residual encoding and attention (Table II).

TABLE II  
ABLATION STUDY OF ARCHITECTURAL COMPONENTS.

Configuration	Dice $\uparrow$	HD95 $\downarrow$
Baseline U-Net	[0.89]	[8.2]
+ ResNet-3 encoder	[0.91]	[6.9]
+ Channel attention only	[0.92]	[6.1]
+ Spatial attention only	[0.93]	[5.7]
+ <b>Spatial-Channel attention (proposed)</b>	<b>[0.94]</b>	<b>[5.1]</b>

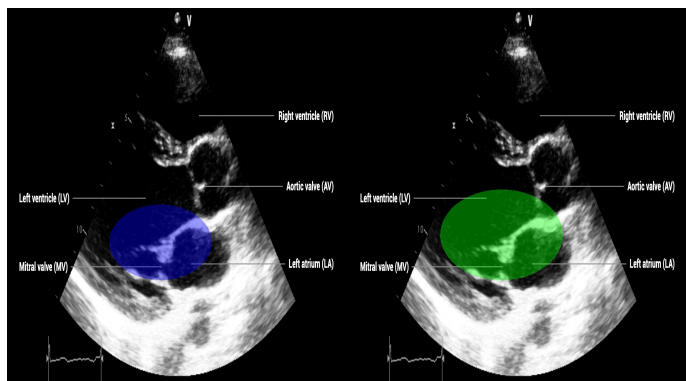


Fig. 2. Qualitative results: Input echo, Ground truth, Baseline, Proposed. Replace fig\_qualitative.pdf.

### C. Qualitative Analysis

Figure 2 shows examples where attention improves boundary adherence, especially in low-contrast regions and frames with dropout. Replace with your predicted masks.

### D. Discussion

The residual encoder improves feature stability and reduces optimization difficulty. The attention module further suppresses speckle-dominated activations and improves focus on cardiac structures, resulting in better contour accuracy (lower HD95). The lightweight encoder makes the approach suitable for near real-time inference.

## VI. CONCLUSION

We presented a compact and accurate echo segmentation framework using a 3-layer ResNet encoder with spatial-channel attention. The model improves overlap and boundary metrics over standard U-Net baselines while remaining computationally efficient. Future work will include multi-structure segmentation, temporal consistency using video models, and cross-domain validation across devices and hospitals.

## ACKNOWLEDGMENT

This work can include funding or institutional support details (optional).

## REFERENCES

- [1] O. Ronneberger, P. Fischer, and T. Brox, “U-Net: Convolutional Networks for Biomedical Image Segmentation,” in *MICCAI*, 2015.
- [2] K. He, X. Zhang, S. Ren, and J. Sun, “Deep Residual Learning for Image Recognition,” in *CVPR*, 2016.
- [3] O. Oktay et al., “Attention U-Net: Learning Where to Look for the Pancreas,” *arXiv preprint arXiv:1804.03999*, 2018.
- [4] S. Woo, J. Park, J.-Y. Lee, and I. S. Kweon, “CBAM: Convolutional Block Attention Module,” in *ECCV*, 2018.
- [5] F. Milletari, N. Navab, and S.-A. Ahmadi, “V-Net: Fully Convolutional Neural Networks for Volumetric Medical Image Segmentation,” in *3DV*, 2016.