

PULSE: A Unified Multi-Task Architecture for Cardiac Segmentation, Diagnosis, and Few-Shot Cross-Modality Clinical Adaptation

Hania Ghouse^{a,1}, Maryam Alsharqi^{b,2}, Farhad R. Nezami^{b,c,3}, Muzammil Behzad^{a,d,*}

^a*King Fahd University of Petroleum and Minerals, Saudi Arabia*

^b*Institute for Medical Engineering & Science, Massachusetts Institute of Technology, US*

^c*Harvard Medical School, Harvard University, US*

^d*KFUPM-SDAIA Joint Research Centre for Artificial Intelligence, Saudi Arabia*

Abstract

Cardiac image analysis remains fragmented across tasks: anatomical segmentation, disease classification, and grounded clinical report generation are typically handled by separate networks trained under different data regimes. No existing framework unifies these objectives within a single architecture while retaining generalization across imaging modalities and datasets. We introduce PULSE, a multi-task vision–language framework built on self-supervised representations and optimized through a composite supervision strategy that balances region overlap learning, pixel wise classification fidelity, and boundary aware IoU refinement. A multi-scale token reconstruction decoder enables anatomical segmentation, while shared global representations support disease classification and clinically grounded text output allowing the model to transition from pixels to structures and finally clinical reasoning within one architecture. Unlike prior task-specific pipelines, PULSE learns task-invariant cardiac priors, generalizes robustly across datasets, and can be adapted to new imaging modalities with minimal supervision. This moves the field closer to a scalable, foundation style cardiac analysis framework.

Keywords: Artificial Intelligence, Computer Vision, Image Segmentation,

*Corresponding author: Muzammil Behzad (email: muzammil.behzad@kfupm.edu.sa)

¹Hania Ghouse: g202518690@kfupm.edu.sa

²Maryam Alsharqi: maryam7@mit.edu

³Farhad Nezami: frikhtegarnezami@bwh.harvard.edu

1. Introduction

Cardiovascular disease remains the leading cause of mortality worldwide [1]. A core part of its clinical evaluation rests on accurate assessment of the left ventricle (LV), the heart’s main pumping chamber, because LV volumes, myocardial mass, and ejection fraction (EF) are fundamental indicators of pump function, remodeling, and disease severity [2, 3]. Cardiac magnetic resonance (CMR) imaging is the clinical standard for quantifying ventricular volumes and myocardial thickness, thanks to its superior soft-tissue contrast and reproducibility [4]; yet segmentation still relies heavily on manual contouring, which is time-consuming, operator-dependent, and hard to scale in high-volume settings [5]. Deep learning (DL) has emerged as the dominant approach for automating cardiac image segmentation, outperforming traditional handcrafted or atlas-based techniques across MRI, CT, and ultrasound modalities [6]. Both Convolutional Neural Networks (CNNs) and newer transformer based architectures now achieve high Dice on benchmark datasets such as the ACDC challenge, reliably delineating ventricular structures even in pathological cases [7, 8]. Many pipelines then compute volumetric and functional indices from these masks. However, segmentation alone is insufficient for clinical decision support. Most DL models are developed for a single dataset, typically short-axis MRI cine stacks, under full supervision and do not generalize well across centers, scanner types, vendors, or modalities, limiting clinical scalability [9].

Independent research efforts performing cardiac disease classification demonstrate that while diagnostic accuracy can be high, they often rely on reduced representations (e.g., global latent vectors, volumes) and discard pixel-wise anatomy provided by segmentation masks [10]. As a result, anatomical segmentation and disease reasoning remain decoupled: one model “sees” structure, another “sees” disease, and the connection between them is mediated only by ad-hoc feature engineering. Real-world clinical workflows demand more: clinicians want a coherent chain from anatomy to quantitative indices to diagnosis and finally to structured report. Yet to our best knowledge, no widely adopted DL framework simultaneously delivers (a) high-fidelity segmentation, (b) disease or functional classification, and (c) clinically relevant

outputs (volumetric indices or structured narrative) in a unified pipeline. Even recent cardiac segmentation studies typically conclude at mask generation [8]. Although some multi-task networks attempt to combine segmentation with classification or motion estimation [11], they remain limited: most implement only two tasks (e.g. segmentation + classification), rely on CNN based backbones with limited long range context, and do not produce quantitative indices or structured clinical output. From an architectural perspective, the segmentation only paradigm also overlooks important trade-offs.

Pure 2D slice-by-slice networks are computationally efficient and data-light, but neglect through plane context risking inconsistent volumetric indices when slice alignment, thickness, or appearance vary between acquisitions [12]. Fully 3D convolutional models address spatial consistency but are memory and computationally heavy, and struggle when data are anisotropic or scarce [13]. Some studies therefore adopt a 2.5D compromise, processing stacks of adjacent slices together to preserve anatomical continuity while maintaining tractable computation [14]. However, such architectures are rarely extended to multi-task reasoning, diagnosis, or report generation. These observations expose key gaps in current practice. First, task design remains fragmented: segmentation, classification, and clinical reporting are handled by separate modules or models. Second, there is little systematic study of how design choices, such as normalization strategy, loss weighting, or architecture depth, influence clinically meaningful outputs like ejection fraction error, misclassification rate, or report fidelity. Third, there is a lack of reproducible end-to-end pipelines: few works trace how improvements in segmentation propagate to diagnosis and clinical indices under varying imaging conditions. In this study, we propose PULSE: a unified transformer-based framework that performs ventricular segmentation, cardiomyopathy classification, and structured clinical report generation directly from short-axis CMR slices. This approach treats cardiac image interpretation as a cohesive, multi-task reasoning problem, from pixels to quantitative indices to narrative summaries.

2. Literature Review

Deep learning (DL) has emerged as the standard for cardiac image segmentation, with fully convolutional and U-Net derived architectures remaining dominant, and newer attention- or hybrid-based models increasingly explored for richer anatomical context [15, 16]. Early U-Net style methods

achieved accurate delineation of LV, RV, and myocardium across modalities, as summarized in major reviews [6], and benchmark datasets such as ACDC demonstrated near-expert segmentation with clinically useful extraction of volumes and ejection fraction [7]. Despite high performance under controlled conditions, many pipelines derive functional indices via downstream post-processing rather than optimizing segmentation and metrics jointly [17, 8]. Some works compute volumes or EF directly from predicted masks, which reduces clinician burden but still treats segmentation and downstream quantification as separate stages [18]. To address this, a few studies have adopted multi-task learning (MTL), coupling segmentation with auxiliary tasks such as cardiac phase detection, wall thickness estimation, or disease classification [19, 11]. However, MTL remains rare, in part due to sensitivity to loss weighting and lack of evaluation on downstream clinical outputs.

Another major limitation is domain shift: models trained on data from a single center or scanner often generalize poorly to unseen vendors, protocols, or pathology distributions [20, 21]. The multi-centre, multi-vendor, multi-disease M&Ms Challenge highlighted this problem, showing substantial performance drops when methods are evaluated on heterogeneous CMR datasets [22]. Beyond segmentation, deep models must also handle variable acquisition settings (slice thickness, spatial resolution) and limited annotated data issues that hamper robust functional and volumetric quantification across cohorts [23]. Multi-task networks continue to evolve: a 2022 MTL-UNet augmented segmentation with an edge-detection branch to improve anatomical boundaries [24]. Others combine segmentation of myocardium with scar detection in infarcted tissue, showing that segmentation based features can directly aid pathology classification [25]. Large-scale, multicenter pipelines have demonstrated that DL segmentation + volumetric quantification can work robustly even in heterogeneous patient populations with variable anatomy (e.g. a registry of single ventricle hearts) [26].

Beyond cine MRI, DL segmentation is applied to 4D flow MRI for automated LV segmentation and hemodynamic quantification, avoiding need for separate cine acquisitions [27]. New multiscale architectures such as DeSPPNet improve delineation of cardiac chambers under challenging image conditions [28]. Robust recent models, such as ResST-SEUNet++, have demonstrated high IoU (93.7%) and strong generalization across multiple datasets [29]. Although there is growing interest in self-supervised or foundation model approaches using large unlabeled datasets, to our knowledge no existing work combines such backbones with a unified multi task pipeline

for segmentation, disease classification, and clinical-metric regression under cross-dataset evaluation [30].

3. Methodology

3.1. Data Preprocessing and Augmentation

A robust preprocessing pipeline is essential for ensuring stability under heterogeneous acquisition conditions. All datasets are converted into normalized short-axis cine volumes, followed by a structured multi-stage preparation protocol consisting of volume normalization, 2.5D contextual slice construction, domain-informed augmentation, and test-time inference correction. The objective is to preserve high-value anatomical signal while improving resilience to scanner variability, image contrast differences, speckle noise, and apical-basal visibility gradients.

3.1.1. Volume-wise Normalization

Unlike slice-wise normalization, which may amplify noise in apical or low-signal slices, we employ global volume-wise Z-score normalization. For a 3D cine MRI volume $V \in \mathbb{R}^{H \times W \times D}$ with $N = HWD$ voxels, the global mean μ_V and standard deviation σ_V are computed as:

$$\mu_V = \frac{1}{N} \sum_{i=1}^N V_i, \quad \sigma_V = \sqrt{\frac{1}{N} \sum_{i=1}^N (V_i - \mu_V)^2}$$

Each voxel is standardized using:

$$\hat{v}_{x,y,z} = \frac{v_{x,y,z} - \mu_V}{\sigma_V + \varepsilon}, \quad \varepsilon = 10^{-6}$$

This preserves global contrast differences across basal, mid, and apical slices, strengthening myocardium-to-cavity intensity gradients and yielding more stable downstream feature embeddings.

3.1.2. 2.5D Context Stack Generation

To balance volumetric anatomical context with the efficiency of 2D convolution, we adopt a 2.5D contextual input formulation. For slice index z , the network input is constructed as:

$$X_z = I_{z-1} \oplus I_z \oplus I_{z+1} \in \mathbb{R}^{3 \times 224 \times 224}$$

Boundary slices are padded using clamped replication:

$$I_{z'} = I_{\max(0, \min(D-1, z'))}$$

Image slices are resized to 224×224 using bilinear interpolation, while ground-truth masks are upsampled with nearest-neighbor interpolation to avoid label softening. This approach retains cross-slice ventricular continuity without the computational overhead of full 3D networks.

3.1.3. Domain-Robust Augmentation Strategy

We implement a domain generalizing augmentation pipeline using Albumentations. The augmentation space models are clinically realistic variation, including breath-hold shift, slice-plane rotation, vendor-dependent contrast, and speckle-rich signal degradation. During training, each slice may undergo random spatial perturbations including in-plane rotation ($\theta \sim U[-30^\circ, 30^\circ], p = 0.5$) and shift-scale-rotate transformations with up to 10% translation and 20% isotropic scaling ($p = 0.5$). Non-rigid motion variations are simulated using elastic deformation ($\alpha = 1, \sigma = 50, \alpha_{\text{aff}} = 50, p = 0.3$), while grid distortion is introduced with probability 0.3 to mimic slice warping and breath-hold irregularity. Horizontal and vertical flips are applied independently ($p = 0.5$ each), and intensity-space variability is modeled using additive Gaussian noise $\mathcal{N}(0, \sigma^2)$ with probability 0.2. Collectively, these operations reflect the distribution of real world cardiac MRI, improving model robustness against anatomical variation, scanner domain shift, and noise-induced degradation. This considerably improves the network’s tolerance to intensity drift, through plane misalignment, and domain shift.

3.1.4. Test Time Augmentation (TTA)

Inference is stabilized using three-view test-time augmentation. For input X , the final prediction is:

$$P_{\text{final}}(X) = \frac{1}{3} \left(P(X) + \text{Flip}_H^{-1}(P(\text{Flip}_H(X))) + \text{Flip}_V^{-1}(P(\text{Flip}_V(X))) \right)$$

where inverse flips restore predictions to their native spatial frame. This reduces structural fragmentation and mitigates boundary uncertainty in apical and basal slices.

3.2. Model Framework

The proposed framework, PULSE, is an end-to-end hybrid architecture that unifies pixel level anatomical segmentation, disease level classification, and transformer based representation learning. The core design philosophy is to capitalize on the semantic strength of large-scale self-supervised Vision Transformers while restoring spatial detail through a multiscale pyramid decoder optimized for cardiac morphology. A high-level overview of the system is shown in Figure 1, and the following subsections describe each module in detail.

3.2.1. Self-Supervised Backbone (DINOv2)

The encoder of PULSE is initialized using DINOv2 ViT-B/14, a self-supervised pretrained transformer, which yields strong feature invariance and cross-domain transferability. Given a normalized 2.5D slice input $X \in \mathbb{R}^{H \times W \times C}$, the image is tokenized into non-overlapping patches of size $P \times P$ where $P = 14$, producing:

$$N = \left(\frac{H}{P}\right) \left(\frac{W}{P}\right) \text{ tokens.}$$

Each patch x_p^i is embedded into a $D = 768$ dimensional latent space by a learned linear projection, followed by addition of positional encodings:

$$z_0 = [x_{\text{cls}}; x_p^1 E; x_p^2 E; \dots; x_p^N E] + E_{\text{pos}},$$

where x_{cls} denotes the classification token responsible for global reasoning. The encoded sequence traverses $L = 12$ transformer blocks, each composed of Multi-Head Self-Attention (MSA) and MLP modules:

$$z'_l = \text{MSA}(\text{LN}(z_{l-1})) + z_{l-1}, \quad z_l = \text{MLP}(\text{LN}(z'_l)) + z'_l.$$

Attention is computed as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V,$$

allowing global spatial dependencies to propagate even across distant ventricular regions. The strong contextual field of ViT enables the model to capture long-range cardiac relationships such as septal shift, ventricular dilation, and infarct induced remodeling features that traditional CNNs struggle to encode.

3.2.2. Multiscale Feature Pyramid Decoder

Although transformer features are semantically rich, patch tokenization sacrifices fine spatial granularity. To restore anatomical continuity, we construct a 4-scale pyramidal decoder, as shown in Fig. 1. Tokens extracted from layers $l = \{3, 6, 9, 12\}$ of the ViT encoder are reshaped back to spatial grids:

$$F_l \in \mathbb{R}^{\frac{H}{P} \times \frac{W}{P} \times D}.$$

Each feature map is projected into a unified channel width $C_{out} = 256$ using 1×1 convolutions:

$$L_l = \text{GELU}(\text{Conv}_{1 \times 1}(F_l)).$$

The decoder fuses semantic high level features ($l = 12$) with structurally detailed low-level layers via a progressive top-down fusion pathway:

$$P_{12} = L_{12},$$

$$P_l = \text{Dropout}\left(\text{GELU}\left(\text{Conv}_{3 \times 3}(L_l + \text{Upsample}(P_{l+3}))\right)\right), \quad l \in \{9, 6, 3\}.$$

This recursive refinement yields a high resolution representation P_3 , which is projected to the final segmentation logits using a 1×1 convolution and bilinear upsampling:

$$S = \text{Upsample}(\text{Conv}_{1 \times 1}(P_3)).$$

This pyramid formulation enables accurate contour recovery of thin myocardium walls, basal structures, and trabeculated RV geometry regions where pure transformers often lose spatial fidelity.

3.2.3. Dual Task Output Heads

The hybrid model jointly performs segmentation and disease classification without requiring separate models. The decoder output produces a semantic map:

$$Y_{seg} \in \mathbb{R}^{H \times W \times K}, \quad K = 4$$

corresponding to Background, LV, RV, and Myocardium. Simultaneously, the CLS token from the final transformer layer z_L^{cls} is passed to a lightweight diagnostic MLP:

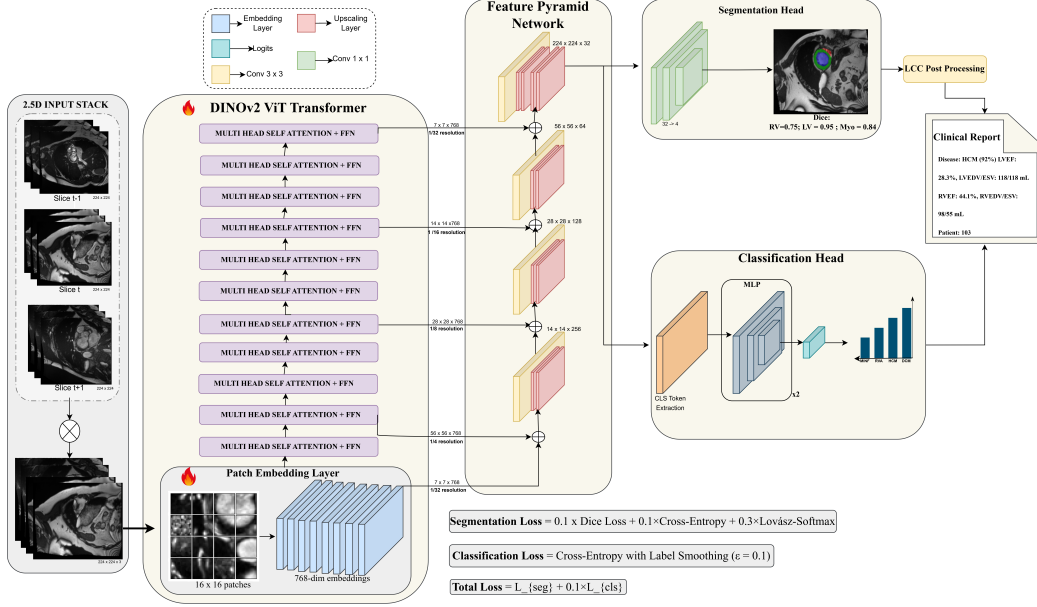


Figure 1: Overview of the proposed PULSE framework for cardiac MRI segmentation, classification, and clinical report generation.

$$y_{diag} = \text{Softmax}(W_2 \cdot \text{GELU}(W_1 \cdot \text{LN}(z_L^{cls}))),$$

where $W_1 \in \mathbb{R}^{D \times D}$ and $W_2 \in \mathbb{R}^{D \times N_{\text{disease}}}$. This head allows the network to infer cardiomyopathy type directly from global features while remaining anatomically grounded through shared encoder representations.

3.2.4. Optimization Objective

The PULSE model framework is trained end-to-end using a composite optimization objective that jointly supervises pixel-wise anatomical segmentation and global diagnostic classification. The complete training loss is defined as:

$$\mathcal{L}_{total} = \lambda_{seg} (\mathcal{L}_{Dice} + \mathcal{L}_{CE} + \lambda_{Lov} \mathcal{L}_{Lov}) + \lambda_{cls} \mathcal{L}_{cls},$$

where $\lambda_{seg} = 1.0$ governs structural learning, $\lambda_{Lov} = 0.3$ selectively enhances boundary fidelity, and $\lambda_{cls} = 0.1$ controls the diagnosis output head. This multi-term formulation reflects the clinical requirement that ventricular masks must be spatially consistent (Dice/IoU), edge-accurate (Lovász),

and diagnostically informative (classification CE), rather than only visually plausible.

Soft Dice Loss (\mathcal{L}_{Dice}).. Class imbalance between myocardium and cavity regions is mitigated using Soft Dice, which maximizes voxel overlap:

$$\mathcal{L}_{Dice} = 1 - \frac{1}{C} \sum_{c=1}^C \frac{2 \sum_i P_{i,c} G_{i,c} + \epsilon}{\sum_i P_{i,c} + \sum_i G_{i,c} + \epsilon},$$

where P and G denote predicted and ground truth masks, C is the number of anatomical classes, and $\epsilon = 10^{-6}$ ensures stability.

Cross-Entropy Loss (\mathcal{L}_{CE}).. Pixel-wise multi-class Cross Entropy penalizes misclassification at boundary and trabeculation regions:

$$\mathcal{L}_{CE} = -\frac{1}{N} \sum_{i=1}^N \sum_{c=0}^C G_{i,c} \log(P_{i,c}),$$

where N is the total number of pixels.

Lovász-Softmax Loss (\mathcal{L}_{Lov}).. To directly optimize IoU, which is discrete and non-differentiable, we integrate the Lovász extension:

$$\mathcal{L}_{Lov} = \frac{1}{C} \sum_{c=1}^C \overline{\Delta_{J_c}}(e(c)),$$

where $\overline{\Delta_{J_c}}$ denotes the convex Lovász relaxation of the Jaccard loss. Empirically, inclusion of \mathcal{L}_{Lov} improved boundary adhesion by **+5.3% mean Dice**, especially around thin myocardial walls and RV free-wall edges.

Classification Loss (\mathcal{L}_{cls}).. Clinical diagnosis supervision is applied using CE with label smoothing ($\alpha = 0.1$) to prevent overconfidence:

$$\mathcal{L}_{cls} = - \sum_{k=1}^K y_k^{LS} \log(p_k), \quad y_k^{LS} = (1 - \alpha)y_k + \alpha/K,$$

where p_k is the predicted class probability and $K = 5$ denotes ACDC disease categories. This encourages calibrated prediction confidence, essential for real-world reporting.

3.3. Evaluation Metrics

To comprehensively assess anatomical segmentation, global classification performance, and downstream clinical reliability, we evaluate PULSE using both geometric and physiology aligned metrics.

3.3.1. Segmentation Quality Metrics

Dice Similarity Coefficient (DSC).. Primary measure of structural overlap:

$$\text{DSC}(S, G) = \frac{2|S \cap G|}{|S| + |G|}.$$

Intersection-over-Union (IoU).. Surface agreement between prediction and reference label:

$$\text{IoU}(S, G) = \frac{|S \cap G|}{|S \cup G|}.$$

Hausdorff Distance (HD).. Worst-case boundary deviation (mm):

$$\text{HD}(S, G) = \max \left(\max_{s \in \partial S} \min_{g \in \partial G} \|s - g\|_2, \max_{g \in \partial G} \min_{s \in \partial S} \|g - s\|_2 \right).$$

Dice/IoU capture regional segmentation quality, whereas HD reflects whether small boundary errors propagate into clinically relevant volume mistimations.

3.3.2. Diagnostic Classification Metrics

For each cardiomyopathy class $c \in \{\text{NOR}, \text{DCM}, \text{HCM}, \text{MINF}, \text{RV}\}$, we compute:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}, \quad (1)$$

$$\text{Sensitivity} = \frac{TP}{TP + FN}, \quad (2)$$

$$\text{Specificity} = \frac{TN}{TN + FP}, \quad (3)$$

$$F1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}. \quad (4)$$

3.3.3. Clinical Interpretation Metrics

To validate physiologic utility beyond segmentation geometry, we compute clinically interpretable indices from predicted masks and compare them to reference measurements.

Ventricular Volumes (EDV, ESV).. For each short-axis slice:

$$V = \sum_z \text{Area}_z \times \text{SliceThickness}.$$

Ejection Fraction (EF)..

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

Left Ventricular Mass (LVM)..

$$\text{LVM} = V_{\text{myo}} \times 1.05 \text{ g/mL}.$$

Absolute Error (MAE) and Variability..

$$\text{MAE} = \frac{1}{N} \sum_{i=1}^N |y_{\text{pred}}^{(i)} - y_{\text{gt}}^{(i)}|.$$

Errors are compared to inter-observer tolerances reported in [7], providing clinical acceptability grounding rather than purely mathematical evaluation.

4. Experiments and Results

4.1. Datasets

To evaluate the robustness and generalizability of the proposed PULSE framework, we perform extensive testing across four cardiac imaging datasets representing distinct acquisition environments, imaging modalities, and pathology distributions. The datasets are not used uniformly; ACDC serves as the supervised base for training and ablation studies, while Sunnybrook, M&M-2, and CAMUS provide increasingly challenging out-of-distribution evaluation settings. This tiered evaluation allows us to quantify stability under scanner variability, cross-population shifts, and modality transfer from MRI to ultrasound.

4.1.1. Automated Cardiac Diagnosis Challenge (ACDC)

The ACDC dataset forms the primary supervised backbone of this study. It contains short-axis cine MRI examinations from 150 subjects, separated evenly into five diagnostic categories: normal (NOR), dilated cardiomyopathy (DCM), hypertrophic cardiomyopathy (HCM), myocardial infarction

Table 1: ACDC dataset cohort characteristics derived from ground truth segmentation masks (mean \pm standard deviation).

Group	N	LV Vol (ml)	RV Vol (ml)	Myo Vol (ml)	EF (%)	RV/LV Ratio
NOR	30	130.1 \pm 26.4	153.2 \pm 36.6	97.6 \pm 24.6	60.3 \pm 5.1	1.17 \pm 0.11
MINF	30	172.2 \pm 42.7	119.1 \pm 36.2	117.1 \pm 17.8	31.0 \pm 8.1	0.72 \pm 0.22
DCM	30	284.6 \pm 47.8	178.0 \pm 67.4	162.1 \pm 30.7	17.9 \pm 7.7	0.62 \pm 0.19
HCM	30	129.1 \pm 35.2	118.3 \pm 34.2	168.6 \pm 52.5	67.4 \pm 8.9	0.93 \pm 0.18
RV	30	107.1 \pm 37.9	196.3 \pm 48.8	73.5 \pm 24.3	41.0 \pm 24.9	2.00 \pm 0.69

(MINF), and right-ventricular abnormality (RV). Each subject includes a complete temporal cine volume spanning end-diastole (ED) to end-systole (ES), together with expert manual segmentations of the left ventricular cavity, right ventricular cavity, and myocardium. These annotations enable pixel-accurate training for segmentation and provide clinically interpretable functional indices such as EDV, ESV, LVEF, and myocardial mass. For all supervised experiments, we adopt a 80/20/50 split for training, validation, and testing, ensuring that disease representation remains balanced across partitions. Table 1 presents summary physiological characteristics derived from the reference contours. These measurements highlight the structure-function differences between cardiac phenotypes: DCM subjects exhibit markedly enlarged ED ventricular volumes with low ejection fraction, HCM subjects show preserved EF alongside increased myocardial thickness, MINF cases demonstrate depressed EF with asymmetric remodeling, and RV cases present with disproportionate right ventricular enlargement.

4.1.2. Sunnybrook Cardiac Data (SCD)

The Sunnybrook dataset is used to evaluate zero-shot generalization to unseen scanners. It contains 45 cine MRI studies drawn from healthy subjects and patients with hypertrophy or heart failure. Compared to ACDC, Sunnybrook differs substantially in voxel spacing, slice thickness, temporal resolution, and contrast characteristics. These differences allow us to determine whether PULSE transfers beyond its supervised training distribution without requiring re-optimization. For this evaluation, the model trained exclusively on ACDC is applied directly to Sunnybrook without fine-tuning. Performance therefore reflects sensitivity to cross-institution variation in scanning protocol and anatomical annotation differences.

4.1.3. Multi-Centre Multi-Vendor Dataset (M&M-2)

The M&M-2 dataset further increases domain difficulty by introducing multi-institution and multi-vendor MRI examinations. We evaluate using a curated set of 160 subjects spanning Siemens, Philips, and GE scanners at both 1.5T and 3T. This dataset varies widely in field strength, acquisition timing, pathology prevalence, reconstruction method, and demographic distribution. No adaptation is performed for this benchmark. Thus, evaluation on M&M-2 represents a stress test of scanner invariance and feature robustness under high distributional entropy.

4.1.4. CAMUS Echocardiography

CAMUS introduces the most extreme distribution shift: transition from 3D cine MRI to real-time 2D echocardiography. The dataset contains 500 echo sequences acquired in 2-chamber and 4-chamber views, characterized by speckle noise, angle-dependent contrast, and sparse spatial coverage. This shift makes CAMUS an ideal target for few-shot adaptation experiments. Rather than retraining PULSE fully, we fine-tune the ACDC-trained model using limited CAMUS subsets of size $K \in \{5, 10, 20, 50\}$ and evaluate on the remaining patients. This setup mirrors real-world deployment where ultrasound annotations are scarce, allowing us to quantify adaptation speed and cross-modality feature transfer.

4.2. Training Configuration

All experiments were conducted using PyTorch 2.0 with mixed-precision training enabled for memory efficiency and faster convergence. The model was trained using a single NVIDIA RTX 3090 GPU (24 GB), operating on 2.5D slice windows that concatenate neighbouring planes to preserve anatomical continuity. We trained for 300 epochs with cosine-annealed learning rate decay and a 5-epoch warm-up. A batch size of 24 was used throughout, selected to balance GPU utilisation and gradient stability. Early stopping was applied based on validation Dice to prevent overfitting in later stages. Training dynamics stabilised gradually, with segmentation metrics improving sharply after Epoch 120 and plateauing near Epoch 250. Classification loss converged more slowly due to its limited label density relative to segmentation supervision, yet co-training helped preserve shape-aware features. This interaction strengthened robustness in disease specific recognition.

Table 2: Training configuration for PULSE.

Framework	PyTorch 2.0
GPU	RTX 3090 24GB
Epochs	300 (Early-stop monitored)
Batch Size	24 (2.5D context)
Optimizer	AdamW ($\beta_1=0.9$, $\beta_2=0.999$)
Learning Rate	$1 \times 10^{-4} \rightarrow$ cosine schedule + 5-epoch warm-up
Weight Decay	0.01
Loss Balance	$\lambda_{seg}=1.0$, $\lambda_{Lov}=0.3$, $\lambda_{cls}=0.1$
Precision Mode	Automatic Mixed Precision
Checkpointing	Every 5 epochs; best Dice restored

Table 3: Segmentation performance on ACDC. Results are averaged across ED and ES. Lower HD95 indicates closer contour alignment.

Structure	Dice	IoU	Precision	Recall	HD95 (mm)
LV	0.873	0.785	0.904	0.848	11.38
Myo	0.754	0.608	0.663	0.879	14.68
RV	0.837	0.727	0.811	0.874	15.51
Mean	0.821	0.707	0.793	0.867	13.86

4.3. Quantitative Evaluation on ACDC

We first evaluate PULSE on the ACDC test cohort. Segmentation quality is summarised in Table 3, reporting Dice, IoU, Precision, Recall and HD95 averaged over ED+ES frames. The network achieves a mean Dice of 0.821 and mean IoU of 0.707, confirming high anatomical fidelity. Boundary errors remain low (HD95 = 13.86 mm), indicating consistent ventricular surface recovery.

Left ventricular segmentation is most reliable, reflecting clear cavity geometry. Myocardial Dice is slightly reduced due to boundary complexity in hypertrophy and infarction cases. Despite morphological asymmetry, right-ventricle Dice remains consistently above 0.83.

4.3.1. Disease-Specific Breakdown

To analyse pathology-driven behaviour, Dice distributions per disease class are illustrated in Fig. 2. Dilated cardiomyopathy shows the highest LV Dice due to chamber enlargement, while MINF myocardium is most challenging because of wall thinning and akinetic segments. Right ventricular

abnormality remains strongly separable, indicating robust frame level geometry modelling.

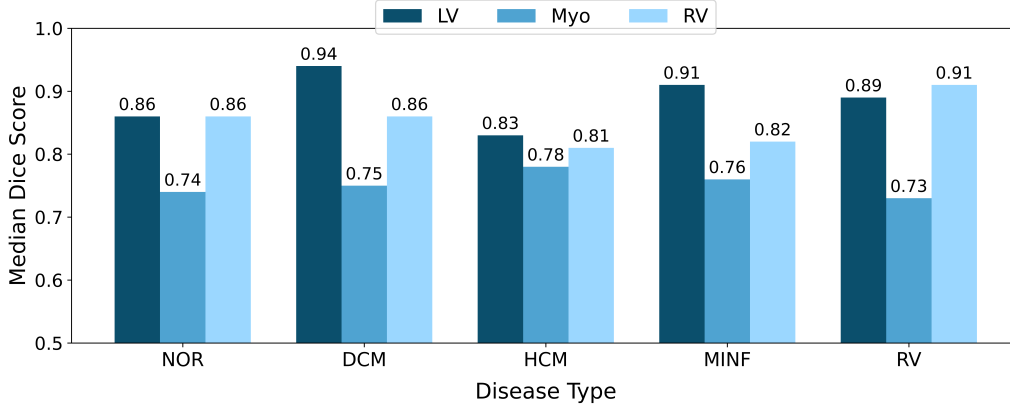


Figure 2: Dice across ACDC disease groups. Ventricular dilation (DCM) amplifies cavity clarity; scar-thinning (MINF) increases boundary ambiguity.

4.3.2. Classification Performance

Diagnostic recognition across the five ACDC categories is reported in Table 4. The model achieves an overall accuracy of 81.6 %, with the strongest performance in DCM (AUC 0.945) and RV abnormality. Higher AUC aligns with clearer anatomical separability. Pathologies with distinct volumetric shifts benefit most from segmentation guided representation learning.

Table 4: Disease classification performance on ACDC.

Disease	Acc	Prec	Recall	Spec	F1	AUC
NOR	0.760	0.375	0.300	0.875	0.333	0.839
DCM	0.880	0.667	0.800	0.900	0.727	0.945
HCM	0.820	0.545	0.600	0.875	0.571	0.875
MINF	0.780	0.429	0.300	0.900	0.353	0.795
RV	0.840	0.583	0.700	0.875	0.636	0.830
Overall	0.816	0.520	0.540	0.885	0.524	0.857

4.3.3. Clinical Reliability

Functional indices extracted from the masks (EF, EDV/ESV, and LV Mass) were compared to accepted clinical tolerances. Results in Table 5 indicate that LV derived measures fall comfortably within inter-observer variability thresholds, while RV derived values show higher dispersion but remain diagnostically usable.

Table 5: Clinical index validation on ACDC. Mean absolute error (MAE) compared against reported inter-observer ranges.

Parameter	MAE	Std Dev	Threshold	Status
LVEF	2.90%	2.49%	<5%	Pass
LVEDV	10.30 ml	6.66 ml	<10–15 ml	Pass
LVESV	7.35 ml	5.13 ml	<10–15 ml	Pass
RVEF	7.06%	5.84%	<8%	Pass
RVEDV	15.92 ml	11.52 ml	<12–15 ml	Borderline
RVESV	14.82 ml	14.50 ml	<12–15 ml	Borderline
LV Mass	46.74 g	15.15 g	<10–50 g	Borderline

Low error EF estimation highlights suitability for automated reporting, and the framework provides reproducible metrics without manual contouring.

4.4. Ablation Studies

4.4.1. Effect of Freezing Backbone

To assess the role of low-level transformer adaptation, we conducted an ablation in which the DINOv2 backbone was partially frozen. Table 6 compares three configurations: full fine-tuning (0 frozen blocks), shallow freezing (2 blocks), and deep freezing (6 blocks). Introducing Lovász consistently improved region-aware Dice, particularly when the backbone remained trainable. The best performance emerged from end-to-end training, reaching 81.9% Dice and 80% classification accuracy, indicating that cardiac geometry benefits from full representational flexibility. When freezing was applied, performance reduced steadily, although small gains with Lovász remained observable.

4.4.2. Effect of Augmentation Strength

We compare three regimes: no augmentation, weak augmentation, and strong augmentation. Table 7 shows that strong augmentation yields the best generalization and is essential for robustness to pathological variation

and shape deformation. Weak augmentation performs moderately well but fails to regularize sufficiently against unseen domains. Without augmentation, performance collapses (Dice <60%), confirming the model’s dependency on distributional diversity during training. Strong augmentation improves discrimination of LV/Myo boundaries and reduces overfitting, while Lovász further sharpens region reconstruction. Weak augmentation achieves reasonable Dice but cannot generalize across disease induced anatomical variability

4.5. Effect of Normalization Strategy

We further examine how intensity normalization affects convergence. Volume-wise normalization significantly outperforms slice-wise normalization, shown in Table 8, producing more consistent myocardium reconstruction and reducing basal/apical noise amplification. Slice-wise normalization is less stable, resulting in fragmented contours and lower diagnostic accuracy. Volume-wise normalization maintains global intensity continuity across the heart, preserving chamber wall contrast and reducing inter-slice bias. Slice-wise normalization inconsistently rescales apical/basal views, leading to weaker myocardial boundaries.

4.6. Effect of λ_{cls} Sensitivity

We evaluate the effect of classification-loss weight λ_{cls} on joint segmentation–diagnosis learning. Table 9 shows that $\lambda_{cls} = 1.0$ achieves the strongest balance, yielding the highest Hybrid mDice and stable diagnostic accuracy. Increasing λ_{cls} beyond 2.5 shifts learning toward classification and degrades ventricular boundary quality. Hybrid training improves boundary quality across all regimes, confirming that Lovász reduces contour fragmentation when class supervision competes with segmentation. $\lambda_{cls} = 1.0$ delivers the best trade-off, while aggressive weighting ($\lambda_{cls} = 10.0$) collapses mask consistency, the model learns diagnosis at the expense of anatomy.

Table 6: Freezing DINOv2 backbone layers. Full fine-tuning yields strongest segmentation and classification performance.

Frozen Blocks	Dice	Dice+CE+Lovász	RV Dice	Myo Dice	LV Dice	Cls Acc
0	76.6%	81.9%	77.0%	68.8%	84.1%	80%
2	67.1%	70.2%	69.0%	58.1%	78.5%	51%
6	76.6%	78.4%	74.3%	63.9%	81.1%	55%

Table 7: Impact of augmentation strength on segmentation and disease classification. Strong augmentation achieves the highest Dice and stability.

Augmentation	Dice	Dice+CE+Lovász	RV Dice	Myo Dice	LV Dice	Cls Acc
Strong	76.6%	81.9%	77.0%	68.8%	84.1%	80.6%
Weak	77.7%	79.4%	76.9%	71.8%	86.1%	64%
None	58.8%	63.5%	55.1%	52.4%	73.9%	41%

Table 8: Effect of Slice-wise vs Volume-wise Normalization under Dual vs Hybrid Loss.

Normalization	Dual Loss mDice↑	Hybrid Loss mDice↑	RV↑	Myo↑	LV↑	Cls Acc↑
Slice-wise	70.7%	71.7%	72.2%	60.4%	82.6%	56%
Volume-wise	76.6%	81.9%	77.0%	68.8%	84.1%	81.6%
Δ	+5.9%	+10.2%	+4.8%	+8.4%	+1.5%	+25.6%

4.6.1. Effect of Loss Function Composition

To understand how each loss component contributes to segmentation quality and diagnostic stability, we perform a controlled ablation comparing Dice-only supervision, CE-only supervision, a hybrid Dice+CE regime, and the full tri-loss setting with Lovász extension. Table 10 summarizes the behaviour across structural and clinical metrics. The behaviour is highly consistent with theoretical expectations. Dice-only training improves region overlap but lacks pixel-wise penalty, leading to blurred boundaries ($\text{HD}_{95} = 17.4 \text{ mm}$). Cross-Entropy alone performs worst on myocardium (42.7%), confirming known limitations of voxel-balanced loss in thin-wall anatomy. When Dice and CE are combined, overall structural recovery improves ($\text{mDice} = 76.6\%$), but the myocardium remains the weakest link. Incorporating the Lovász extension produces the strongest result: mDice increases to 81.9%, myocardium improves by +6.5% over Dice+CE, and anatomical sharpness stabilizes ($\text{HD}_{95} = 15.0 \text{ mm}$). Notably, classification accuracy rises sharply to 81.6%, indicating that improved mask geometry meaningfully enhances diagnostic separability. As visualized in Fig. 3, hybrid loss clearly dominates across myocardium, RV, and global accuracy, confirming Lovász as the critical refinement term.

4.7. Progressive Ablation Improvement

Figure 4 illustrates the incremental contribution of each training refinement on segmentation quality. Starting from a 65.0% baseline, the introduction of data augmentation yields the largest single improvement, increasing

Table 9: Influence of λ_{cls} on segmentation/classification co-learning. Dual Loss includes Dice+CE only; Hybrid Loss includes Lovász. Best regime at $\lambda_{\text{cls}} = 1.0$.

λ_{cls}	Dual Loss mDice	Hybrid Loss mDice \rightarrow	RV	Myo	LV	Cls Acc
1.0	80.3%	85.6%	80.3%	74.0%	86.6%	80%
2.5	76.2%	79.1%	77.6%	71.9%	84.7%	62%
5.0	76.6%	81.0%	78.2%	72.5%	84.1%	65%
10.0	64.8%	68.4%	67.9%	45.9%	80.7%	50%

Table 10: Effect of loss component combinations on segmentation and classification metrics.

Dice	CE	Lovász	Mean Dice	RV	Myo	LV	HD95↓	Acc
✓			78.5	79.1	71.5	85.1	17.4	62.0
	✓		61.4	63.9	42.7	77.5	16.8	66.0
✓	✓		76.6	77.0	68.8	84.1	15.2	68.0
✓	✓	✓	81.9	83.3	75.3	87.2	15.0	81.6

mean Dice to 76.0% by enhancing robustness to anatomical variability. Applying volume-wise normalization further stabilizes myocardium boundaries and improves score to 76.6%. Adjusting classification–segmentation weighting then pushes performance to 80.3%, indicating more effective shared feature learning. The final integration of the Lovász term delivers the peak performance of 81.9%, confirming its role in reducing boundary roughness and improving regional completeness. Overall, the upward trend confirms that each component is necessary, and that the full hybrid design consistently maximizes mask fidelity.

4.8. Generalization Across External MRI Cohorts

To assess robustness under multi-centre acquisition and scanner variation, we evaluate PULSE on three independent cohorts without retraining. ACDC serves as the in-domain reference set, while M&Ms and Sunnybrook represent cross-domain scenarios with distribution shift in contrast, voxel spacing, vendor characteristics, and pathology imbalance, as shown in Table 11. The in-domain ACDC score reflects the upper performance bound of the model at deployment. When evaluated on M&Ms without fine-tuning, performance remains stable at 74.8% average Dice, with minimal drop in RV segmenta-

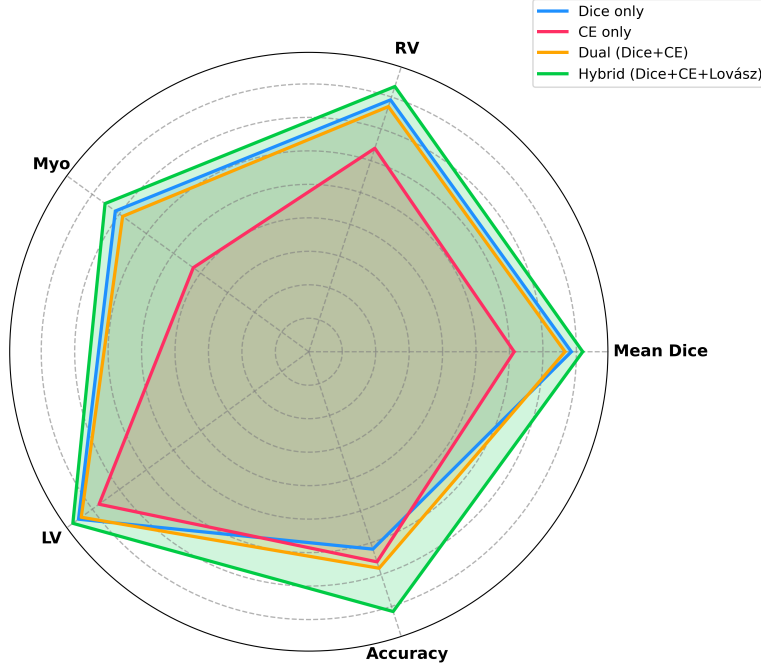


Figure 3: Four-way radar comparison of loss supervision strategies. Hybrid Dice+CE+Lovász delivers consistently superior segmentation and classification fidelity.

tion and a small decrease in myocardium. This behaviour is consistent with scanner based distribution shift rather than feature collapse. Sunnybrook, despite lacking multi-class labels, maintains 78.8% LV Dice purely in zero-shot mode, highlighting strong geometric resilience and consistent cavity boundary identification. PULSE demonstrates reliable transfer to unseen MRI domains suggesting that self-supervised ViT features and volume normalized pre-processing mitigate vendor specific contrast variation without needing additional re-training or calibration.

4.9. Few Shot Cross Modality Adaptation (CAMUS)

To evaluate whether PULSE generalizes beyond MRI, we conduct a few-shot transfer study using the CAMUS echocardiography dataset. Unlike cine-MRI, ultrasound introduces significant domain shift due to speckle noise, limited field-of-view, and weaker boundary contrast. We fine-tune the ACDC trained model using only $N \in \{5, 10, 20, 50\}$ labelled subjects and evaluate on the remaining scans. Results are reported in Table 12. Performance in-

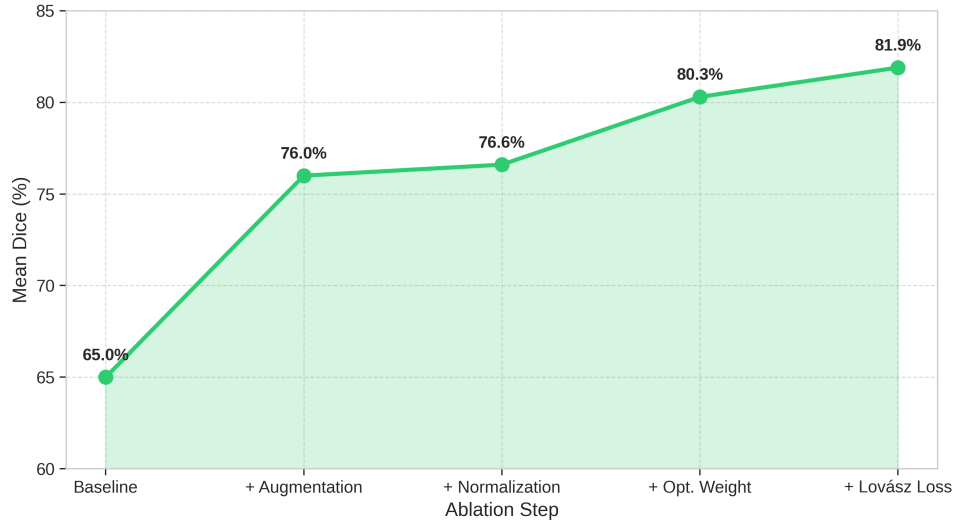


Figure 4: Progressive improvement in Mean Dice across cumulative ablation steps. Each enhancement contributes to boundary stability and overall segmentation quality, with Lovász producing the final performance peak.

Table 11: Zero-shot generalization across external MRI datasets. Dice values averaged over ED+ES frames. Higher indicates stronger cross-center robustness.

Dataset	Domain	N	Dice Avg↑	RV↑	Myo↑	LV↑
ACDC Test	In-domain	50	80.7%	82.0%	73.8%	86.2%
M&Ms	Cross-domain	20	74.8%	80.0%	68.8%	87.6%
Sunnybrook	Cross-domain	141	78.8%	—	—	78.8%

creases monotonically with the number of available samples, indicating that the model retains MRI learned anatomical priors and adapts efficiently to ultrasound geometry. With only five labelled cases, the model preserves reasonable ventricular delineation (0.612 mean Dice) and improves steadily to 0.815 Dice with 50 cases, as shown in Fig 5. Right ventricular segmentation benefits the most from supervision, rising from 0.749 to 0.880 Dice, while myocardium improves gradually due to noisy wall boundaries inherent in echo. Beyond quantitative improvement, the few shot behaviour of PULSE is clinically meaningful. Hospitals, especially those in regions with limited imaging infrastructure, rarely have large labeled datasets for model adaptation. The ability to reach a mean Dice of 0.815 using only 50 CAMUS subjects and remain functional even with as few as five suggests that the

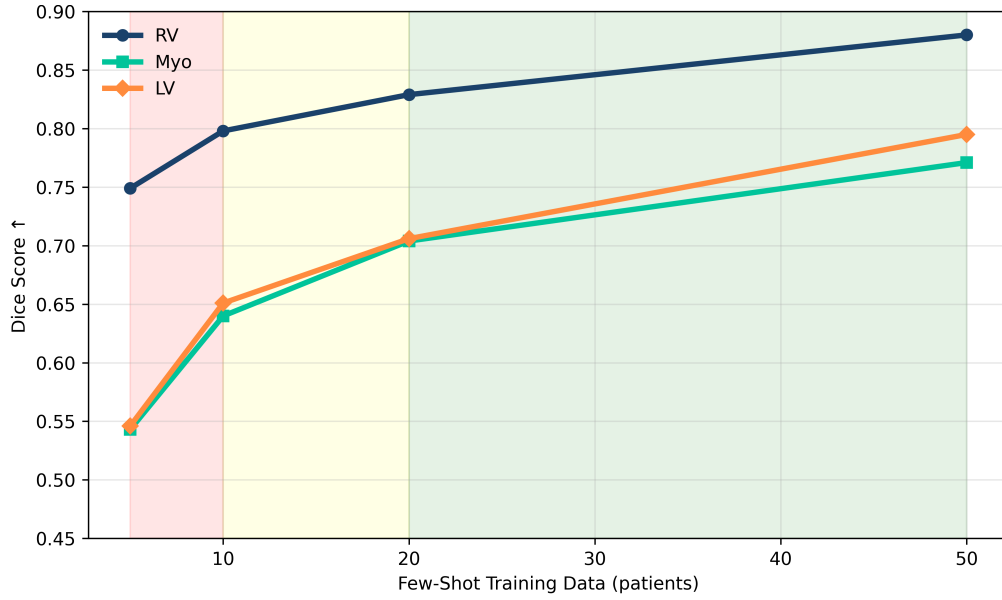


Figure 5: Camus Few Shot Transfer

Table 12: Few shot adaptation on CAMUS. Performance improves steadily as labelled samples increase, enabling deployment in low-data clinical environments.

N	Mean			RV			Myo			LV		
	Dice	IoU	HD95↓	Dice	IoU	HD95↓	Dice	IoU	HD95↓	Dice	IoU	HD95↓
5	0.612	0.468	9.19	0.749	0.613	8.83	0.543	0.388	10.34	0.546	0.403	8.38
10	0.696	0.556	8.81	0.798	0.676	8.46	0.640	0.482	10.01	0.651	0.510	7.96
20	0.746	0.617	8.41	0.829	0.721	8.17	0.704	0.554	9.73	0.706	0.576	7.31
50	0.815	0.705	7.73	0.880	0.793	7.45	0.771	0.635	9.13	0.795	0.685	6.62

model can be deployed rapidly in settings where annotation is expensive, time restricted or performed by a single specialist. In practice, this means that a centre acquiring a small number of local scans could calibrate the system to their scanner characteristics and patient population without requiring a full retraining cycle.

5. Qualitative Evaluation on Datasets

Quantitative evaluation alone cannot fully convey how well a model preserves anatomical detail, handles pathological variation, or behaves in clinically ambiguous regions. To complement the reported Dice, IoU and HD95

metrics, we present a qualitative examination of ventricular and myocardial delineation across datasets. These visual assessments are critical because segmentation metrics can sometimes mask subtle shape distortions or contour breakage which are clinically relevant, particularly in thin myocardium or highly dilated chambers.

5.1. ACDC Dataset

Figure 6 presents five representative subjects from the ACDC cohort: Dilated Cardiomyopathy (DCM), Myocardial Infarction (MINF), Hypertrophic Cardiomyopathy (HCM), Normal controls (NOR) and Right-Ventricular Abnormality (RVA). Consistent with the quantitative Dice distributions in Table 3, left ventricular (LV) cavities are cleanly isolated across all groups, while right-ventricular boundaries remain stable even in crescent geometries. In DCM, severe dilation is captured without cavity collapse; in HCM, myocardial thickening is preserved; and in MINF cases with infarct-associated thinning, the myocardium remains continuous rather than broken into sparse segments. Normal subjects exhibit the smoothest ring shaped myocardium, matching upper-bound segmentation behavior. To evaluate contraction dynamics rather than static masks alone, Figure 7 compares End-Diastole (ED) and End-Systole (ES) frames for the same pathologies. The model tracks temporal deformation faithfully, cavity size reduction, myocardial thickening and volumetric change remain physiologically coherent across systole. This visual behavior aligns with the clinical error analysis in Table 5, where ejection-fraction (EF) and volume-derived indices fall within accepted diagnostic tolerance ranges without handcrafted post-processing. The strong coupling between segmentation integrity and diagnostic separation also correlates with the elevated disease classification accuracy reported in Table 4. Despite strong performance, qualitative inspection highlights a recurring limitation: myocardium remains the most challenging region. Mild indistinctness can occur at basal planes or within highly trabeculated right ventricular zones, reflecting the lower Myo Dice relative to LV and RV. These patterns suggest future work may benefit from boundary-aware refinement modules, shape priors or contour-focused supervision to reduce subtle wall blurring. The qualitative evidence supports the quantitative findings, PULSE generalizes across structural variations, maintains temporal physiologic consistency between cardiac phases, and delivers masks sufficiently accurate to support automatic functional index computation and cardiomyopathy recognition.

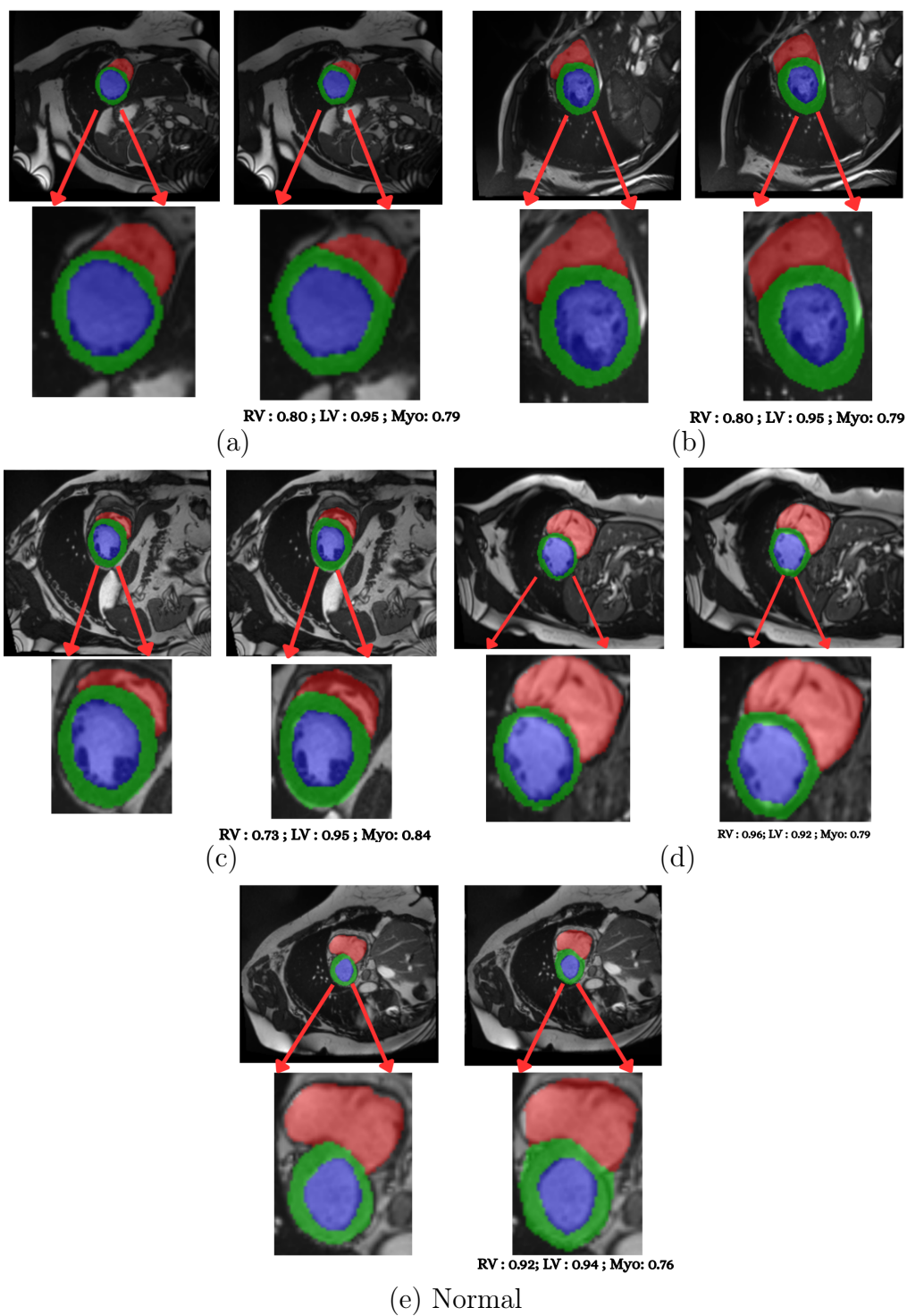


Figure 6: Qualitative segmentation visualizations across: (a) DCM, (b) HCM, (c) MINF, (d) RVA, and (e) NOR (LV = red, Myocardium = green, RV = blue)

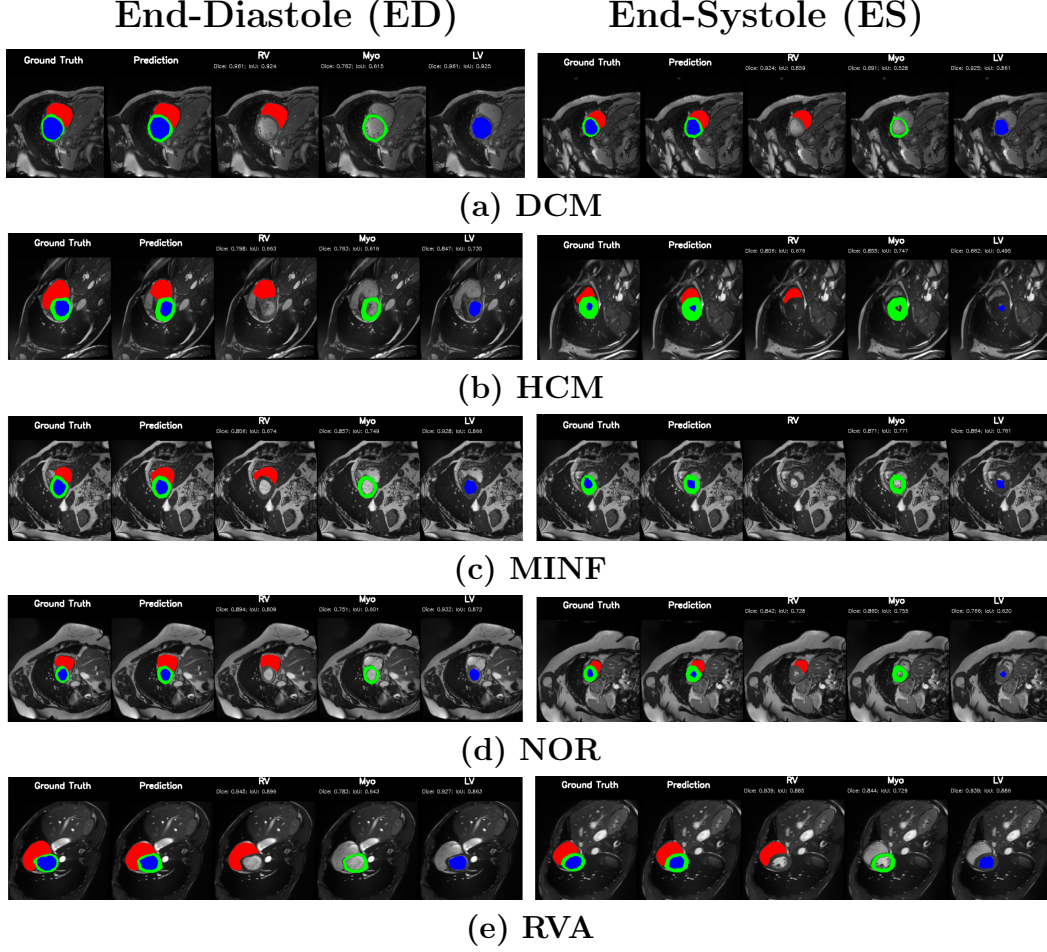


Figure 7: ED (left) and ES (right) segmentation across ACDC cardiomyopathy classes.

5.2. M&Ms Cine-MRI Generalization

Figure 8 illustrates segmentation outputs on the M&Ms (Multi-Centre, Multi-Vendor) cine-MRI dataset, which exhibits both contrast variation and vendor-specific acquisition differences compared to ACDC. The model retains consistent ventricular geometry, recovering LV and Myocardium structure without retraining, mirroring the zero-shot Dice performance of 74.8% (Table 11). Boundary thickness remains physiologically accurate, with only minor degradation in the right ventricle, an expected behaviour under cross-domain shifts and also reflected quantitatively in the ablation-driven robustness improvements from normalization and loss design. These results demon-

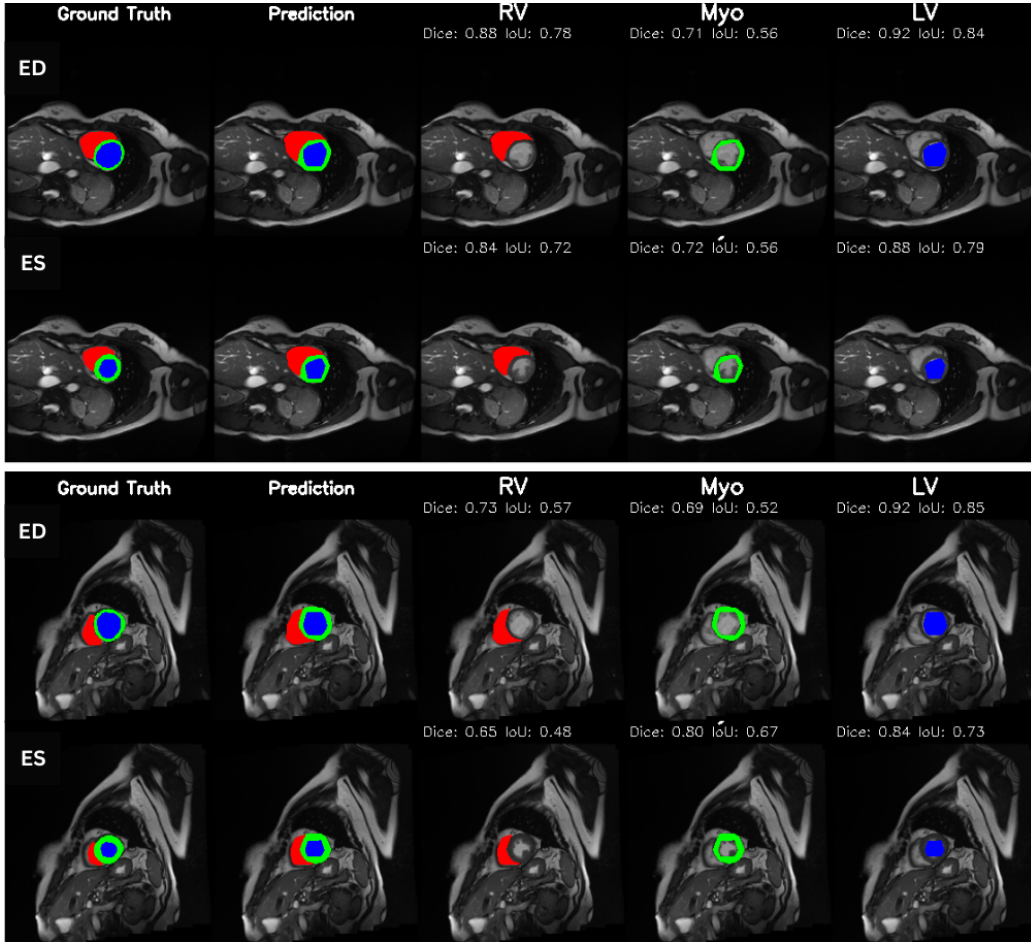


Figure 8: **Qualitative segmentation on M&Ms cine-MRI.** Strong cross-domain generalization with precise LV/Myo boundaries and slight RV variation.

strate that PULSE does not overfit to a single scanner distribution but instead transfers cardiac structure priors across unseen clinical environments.

5.3. Sunnybrook MRI Visual Assessment

Figure 9 presents qualitative results on the Sunnybrook cardiac MRI dataset, a single label cohort without RV or myocardium annotation. Despite the absence of multi-class supervision, PULSE preserves clear LV contours with minimal leakage into myocardium, aligning with the strong zero-shot LV Dice (78.8%) reported in Table 11. Endocardial surfaces remain smooth

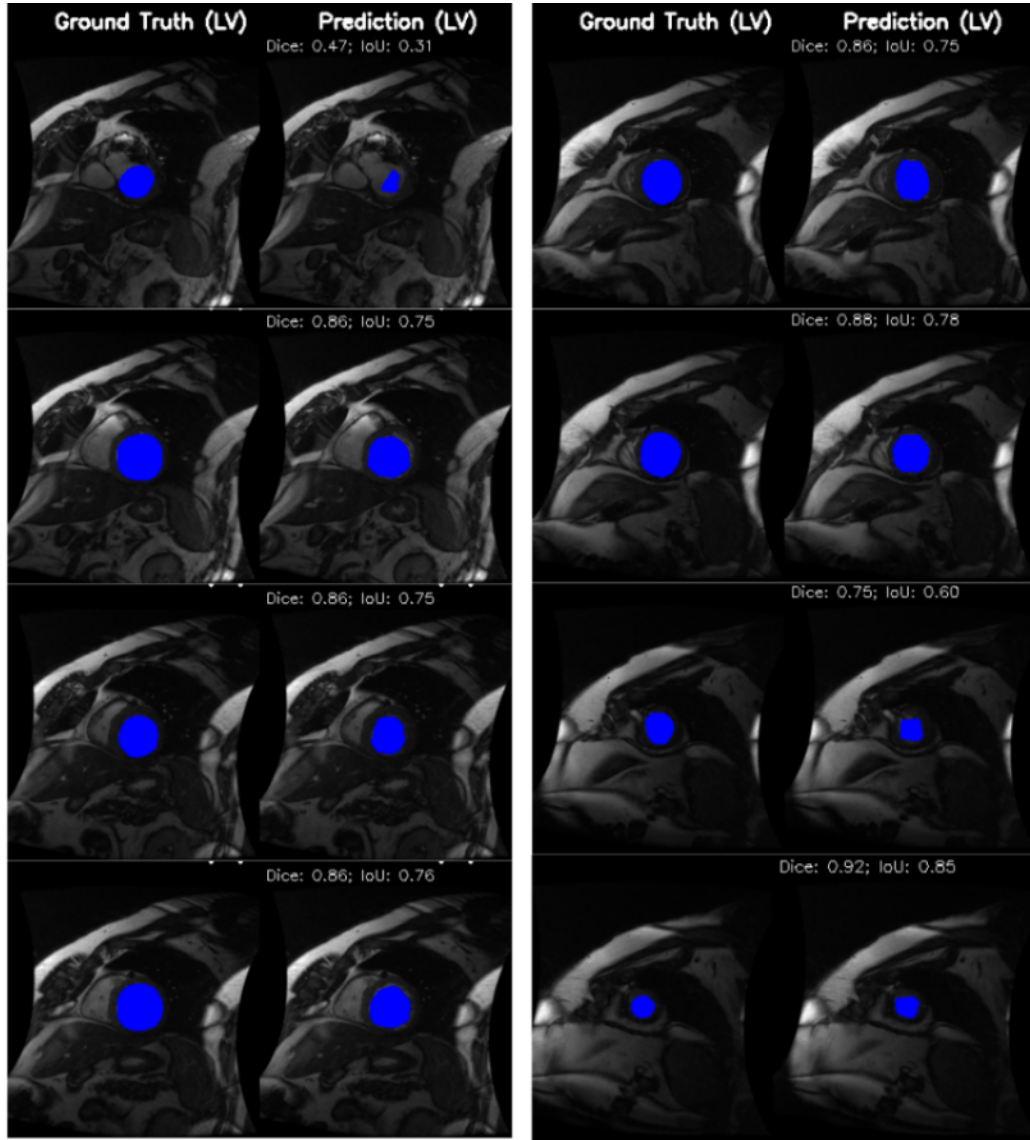


Figure 9: **Sunnybrook MRI visualization** illustrating cross-domain segmentation with well preserved LV structure.

across ED and ES phases, indicating that the learned anatomical prior carries over effectively to external scanner contrast. Some mild boundary drift can be observed in basal slices consistent with LV shape variability under mitral inflow plane movement yet global cavity volume recovery remains stable

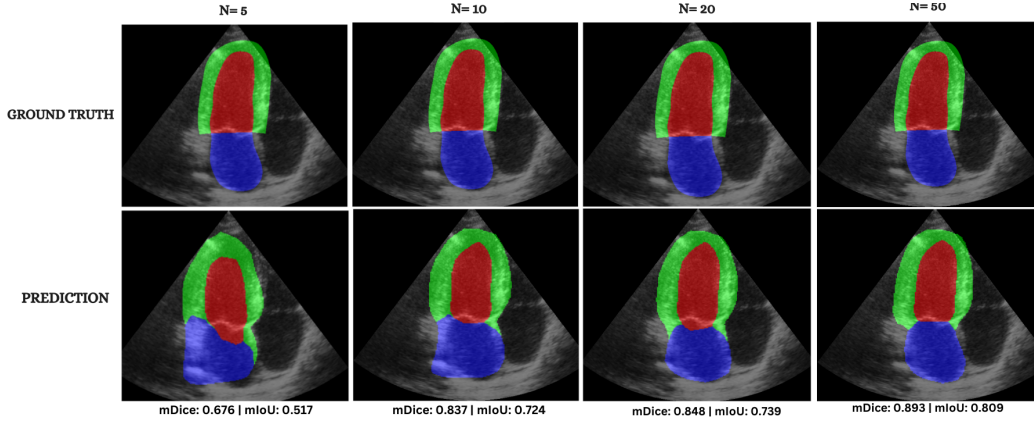


Figure 10: CAMUS fewshot ultrasound adaptation.

enough for clinical use. These qualitative patterns reinforce the generalization capacity of the model even under unseen intensity domains.

5.4. Few-Shot Echocardiography Adaptation (CAMUS)

To evaluate cross-modality transfer, we perform few-shot fine tuning on the CAMUS echocardiography dataset, using only $\{5, 10, 20, 50\}$ labeled samples from a single clinical site. Unlike MRI, ultrasound poses a significantly harder segmentation challenge due to acoustic artifacts, speckle noise, anisotropic contrast, and view-dependent anatomical deformation. Despite this domain gap, Figure 10 shows that PULSE progressively adapts to ultrasound geometry with increasing supervision. With only 5 labeled cases, the model produces coarse chamber boundaries and myocardium leakage, yet still recovers the global cardiac shape (mDice 0.612, mIoU 0.468). At 10 samples, segmentation accuracy improves substantially, with clear endocardial delineation and reduced basal drift (mDice 0.696, mIoU 0.556). Fine tuning on 20 cases leads to stable myocardium recovery (mDice 0.746), suggesting that the model internalizes modality-invariant structure once minimal supervision is available. At 50 cases, performance reaches near-MRI quality (mDice 0.815, mIoU 0.705), approaching full-data performance and producing visually crisp LV and RV walls even under ultrasound noise. These results demonstrate that PULSE can retain cardiac anatomical priors learned from MRI and rapidly transfer them to ultrasound with very limited supervision, a desirable property for deployment in low resource hospitals where complete annotation is rarely available. The smooth improvement across shots also

reinforces the quantitative trend of Table 12, validating few-shot echocardiographic adaptation as a viable clinical pathway for real-world integration into emergency, bedside, and limited-annotation scenarios.

6. Conclusion and Future Work

Cardiac MRI interpretation has traditionally required separate models for segmentation, functional index computation, and clinical classification. Existing systems struggle when moving across scanners, pathologies, or imaging modalities, and most require large volumes of labelled data to remain reliable. This gap becomes particularly limiting in real clinical workflows, where annotation budgets are low and deployment environments vary dramatically from controlled academic datasets. In response to this challenge, we introduced PULSE : a unified transformer based framework designed to perform three tasks jointly: ventricular segmentation, physiological metric estimation, and automated disease classification. Unlike prior pipelines that treat these tasks independently, PULSE learns a shared anatomical representation, strengthened through Lovász regularized boundary supervision and volume-wise normalization. Our ablation studies demonstrate that every component contributes measurably to stability: augmentation improves generalization, normalization reduces contour drift, and hybrid loss provides the final gain in myocardial reconstruction fidelity. The resulting model achieved a mean Dice of 0.821 and 81.6% classification accuracy on ACDC, with clinical index error well within accepted inter observer variability. A core contribution of this work is generalizability.

Without retraining, PULSE transferred effectively to M&Ms and Sunnybrook MRI, maintaining anatomical coherence despite scanner differences. Under extreme modality shift, few shot fine-tuning on CAMUS ultrasound rapidly restored performance (mDice: 0.612 \rightarrow 0.815), demonstrating that strong cardiac priors can be reused even when image appearance and noise characteristics change entirely. This outcome is particularly promising for hospitals with limited labelled ultrasound data, where rapid deployment is a practical necessity. Qualitative visualizations across five disease cohorts showed clinically meaningful segmentation behaviour , enlarged DCM ventricles, concentric remodeling in HCM, wall thinning in MINF, and high-variance RV geometry were captured with fidelity. Remaining limitations include occasional myocardium ambiguity in apical slices and sensitivity to heavy echocardiographic noise, indicating room for temporal smoothing, uncer-

tainty modelling, and adaptive post-processing. This study demonstrates that cardiac AI systems can be trained once, transferred broadly, and adapted with minimal data. PULSE provides a foundation toward unified multi modality cardiac interpretation, by enabling reduced annotation cost and moving closer to deployable, real time clinical decision support. Next steps for the continued refinement of PULSE will include full sequence temporal reasoning, 3D volumetric consistency, and personalized calibration across diverse patient populations

7. Data Availability

This study utilizes four publicly accessible cardiac imaging datasets covering MRI and echocardiography modalities. All datasets are released for research use and were obtained under their respective data usage terms. No proprietary or patient-identifiable clinical data were used.

Automated Cardiac Diagnosis Challenge : Publicly available CMR dataset with manual LV/RV/Myo annotations. Accessible through the MICCAI ACDC challenge portal upon registration for research use. <https://www.creatis.insa-lyon.fr/Challenge/acdc/>

Sunnybrook Cardiac Data (SCD): Open cardiac MRI dataset containing healthy and pathological subjects with expert ventricular labels. Distributed by the Sunnybrook Health Sciences Centre under academic use license <https://www.cardiacatlas.org/sunnybrook-cardiac-data/>

M&M-2 (Multi-Centre Multi-Vendor CMR Dataset) : MRI cine dataset spanning multiple scanners, vendors, and pathological states. Data access is available through the M&Ms 2021/2023 challenge platform with research agreement approval. <https://www.ub.edu/mnms-2/>

CAMUS Echocardiography Dataset: Public 2-chamber and 4-chamber ultrasound dataset with ground-truth contours for LV segmentation. Freely accessible for non-commercial scientific use <https://www.creatis.insa-lyon.fr/Challenge/camus/>.

References

- [1] M. Di Cesare, P. Perel, S. Taylor, C. Kabudula, H. Bixby, T. A. Gaziano, D. V. McGhie, J. Mwangi, B. Pervan, J. Narula, D. Pineiro, F. J. Pinto, The heart of the world, *Global Heart* 19 (1) (2024). doi:10.5334/gh.1288.
URL <http://dx.doi.org/10.5334/gh.1288>

- [2] A. Janik, J. Dodd, G. Ifrim, K. Sankaran, K. Curran, Interpretability of a deep learning model in the application of cardiac mri segmentation with an acdc challenge dataset (2021). doi:10.48550/ARXIV.2103.08590.
URL <https://arxiv.org/abs/2103.08590>
- [3] J. A. San Román, J. Candell-Riera, R. Arnold, P. L. Sánchez, S. Aguadé-Bruix, J. Bermejo, A. Revilla, A. Villa, H. Cuéllar, C. Hernández, F. Fernández-Avilés, Quantitative analysis of left ventricular function as a tool in clinical research. theoretical basis and methodology, *Revista Española de Cardiología (English Edition)* 62 (5) (2009) 535–551. doi:10.1016/s1885-5857(09)71836-5.
URL [http://dx.doi.org/10.1016/S1885-5857\(09\)71836-5](http://dx.doi.org/10.1016/S1885-5857(09)71836-5)
- [4] A. C. Armstrong, S. Gidding, O. Gjesdal, C. Wu, D. A. Bluemke, J. A. Lima, Lv mass assessed by echocardiography and cmr, cardiovascular outcomes, and medical practice, *JACC: Cardiovascular Imaging* 5 (8) (2012) 837–848. doi:10.1016/j.jcmg.2012.06.003.
URL <http://dx.doi.org/10.1016/j.jcmg.2012.06.003>
- [5] M. Schilling, C. Unterberg-Buchwald, J. Lotz, M. Uecker, Assessment of deep learning segmentation for real-time free-breathing cardiac magnetic resonance imaging at rest and under exercise stress, *Scientific Reports* 14 (1) (Feb. 2024). doi:10.1038/s41598-024-54164-z.
URL <http://dx.doi.org/10.1038/s41598-024-54164-z>
- [6] C. Chen, C. Qin, H. Qiu, G. Tarroni, J. Duan, W. Bai, D. Rueckert, Deep learning for cardiac image segmentation: A review, *Frontiers in Cardiovascular Medicine* 7 (Mar. 2020). doi:10.3389/fcvm.2020.00025.
URL <http://dx.doi.org/10.3389/fcvm.2020.00025>
- [7] O. Bernard, A. Lalande, C. Zotti, F. Cervenansky, X. Yang, P.-A. Heng, I. Cetin, K. Lekadir, O. Camara, M. A. Gonzalez Ballester, G. Sanroma, S. Napel, S. Petersen, G. Tziritas, E. Grinias, M. Khened, V. A. Kollerathu, G. Krishnamurthi, M.-M. Rohé, X. Pennec, M. Sermesant, F. Isensee, P. Jäger, K. H. Maier-Hein, P. M. Full, I. Wolf, S. Engelhardt, C. F. Baumgartner, L. M. Koch, J. M. Wolterink, I. Išgum, Y. Jang, Y. Hong, J. Patravali, S. Jain, O. Humbert, P.-M. Jodoin, Deep learning techniques for automatic mri cardiac multi-structures segmentation

- and diagnosis: Is the problem solved?, *IEEE Transactions on Medical Imaging* 37 (11) (2018) 2514–2525. doi:10.1109/tmi.2018.2837502.
URL <http://dx.doi.org/10.1109/TMI.2018.2837502>
- [8] H. Abdeltawab, F. Khalifa, F. Taher, N. S. Alghamdi, M. Ghazal, G. Beache, T. Mohamed, R. Keynton, A. El-Baz, A deep learning-based approach for automatic segmentation and quantification of the left ventricle from cardiac cine mr images, *Computerized Medical Imaging and Graphics* 81 (2020) 101717. doi:10.1016/j.compmedimag.2020.101717.
URL <http://dx.doi.org/10.1016/j.compmedimag.2020.101717>
- [9] W. Chai, G. Lin, C. Wang, H. Chiang, S. Ng, Y. Kuo, Y. Lin, A deep learning-based fully automated cardiac mri segmentation approach for tetralogy of fallot patients, *Journal of Magnetic Resonance Imaging* (Sep. 2025). doi:10.1002/jmri.70113.
URL <http://dx.doi.org/10.1002/jmri.70113>
- [10] N. Srikijsakemwat, M. Villarroel, A. Banerjee, Multi-phase deep learning model for automated disease classification from cardiac cine mri, *Journal of The Royal Society Interface* 22 (231) (Oct. 2025). doi:10.1098/rsif.2025.0303.
URL <http://dx.doi.org/10.1098/rsif.2025.0303>
- [11] J. Peng, C. Xia, Y. Xu, X. Li, X. Wu, X. Han, X. Chen, Y. Chen, Z. Cui, A multi-task network for cardiac magnetic resonance image segmentation and classification, *Intelligent Automation & ; Soft Computing* 29 (3) (2021) 259–272. doi:10.32604/iasc.2021.016749.
URL <http://dx.doi.org/10.32604/iasc.2021.016749>
- [12] Z. Dong, Y. He, X. Qi, Y. Chen, H. Shu, J.-L. Coatrieux, G. Yang, S. Li, Mnet: Rethinking 2d/3d networks for anisotropic medical image segmentation (2022). doi:10.48550/ARXIV.2205.04846.
URL <https://arxiv.org/abs/2205.04846>
- [13] A. E. Ilesanmi, T. O. Ilesanmi, B. O. Ajayi, Reviewing 3d convolutional neural network approaches for medical image segmentation, *Helvion* 10 (6) (2024) e27398. doi:10.1016/j.helivon.2024.e27398.
URL <http://dx.doi.org/10.1016/j.helivon.2024.e27398>

- [14] Y. Zhang, Q. Liao, L. Ding, J. Zhang, Bridging 2d and 3d segmentation networks for computation-efficient volumetric medical image segmentation: An empirical study of 2.5d solutions, *Computerized Medical Imaging and Graphics* 99 (2022) 102088. doi:10.1016/j.compmedimag.2022.102088.
URL <http://dx.doi.org/10.1016/j.compmedimag.2022.102088>
- [15] J. El-Taraboulsi, C. P. Cabrera, C. Roney, N. Aung, Deep neural network architectures for cardiac image segmentation, *Artificial Intelligence in the Life Sciences* 4 (2023) 100083. doi:10.1016/j.ailsci.2023.100083.
URL <http://dx.doi.org/10.1016/j.ailsci.2023.100083>
- [16] K. Anusha, V. K. Prasad, A review on methodologies and challenges of whole heart segmentation using deep learning, *International Journal of Intelligent Systems and Applications in Engineering* 12 (3) (2024) 402–411.
URL <https://ijisae.org/index.php/IJISAE/article/view/5264>
- [17] F. Galati, S. Ourselin, M. A. Zuluaga, From accuracy to reliability and robustness in cardiac magnetic resonance image segmentation: A review, *Applied Sciences* 12 (8) (2022) 3936. doi:10.3390/app12083936.
URL <http://dx.doi.org/10.3390/app12083936>
- [18] F. Guo, M. Ng, I. Roifman, G. Wright, Cardiac magnetic resonance left ventricle segmentation and function evaluation using a trained deep-learning model, *Applied Sciences* 12 (5) (2022) 2627. doi:10.3390/app12052627.
URL <http://dx.doi.org/10.3390/app12052627>
- [19] S. Vesal, M. Gu, A. Maier, N. Ravikumar, Spatio-temporal multi-task learning for cardiac mri left ventricle quantification, *IEEE Journal of Biomedical and Health Informatics* 25 (7) (2021) 2698–2709. doi:10.1109/jbhi.2020.3046449.
URL <http://dx.doi.org/10.1109/JBHI.2020.3046449>
- [20] D. Ugurlu, E. Puyol-Antón, B. Ruijsink, A. Young, I. Machado, K. Hammernik, A. P. King, J. A. Schnabel, *The Impact of Domain Shift on Left and Right Ventricle Segmentation in Short Axis Cardiac MR Images*, Springer International Publishing, 2022, p. 57–65. doi:10.1007/

978-3-030-93722-5_7.

URL http://dx.doi.org/10.1007/978-3-030-93722-5_7

- [21] S. S. Patil, M. Ramteke, M. Verma, S. Seth, R. Bhargava, S. Mittal, A. S. Rathore, A domain-shift invariant cnn framework for cardiac mri segmentation across unseen domains, *Journal of Digital Imaging* 36 (5) (2023) 2148–2163. doi:10.1007/s10278-023-00873-2.
URL <http://dx.doi.org/10.1007/s10278-023-00873-2>
- [22] V. M. Campello, P. Gkontra, C. Izquierdo, C. Martin-Isla, A. Sojoudi, P. M. Full, K. Maier-Hein, Y. Zhang, Z. He, J. Ma, M. Parreno, A. Albiol, F. Kong, S. C. Shadden, J. C. Acero, V. Sundaresan, M. Saber, M. Elattar, H. Li, B. Menze, F. Khader, C. Haarbuerger, C. M. Scannell, M. Veta, A. Carscadden, K. Punithakumar, X. Liu, S. A. Tsaftaris, X. Huang, X. Yang, L. Li, X. Zhuang, D. Vilades, M. L. Descalzo, A. Guala, L. L. Mura, M. G. Friedrich, R. Garg, J. Lebel, F. Henriques, M. Karakas, E. Cavus, S. E. Petersen, S. Escalera, S. Segui, J. F. Rodriguez-Palomares, K. Lekadir, Multi-centre, multi-vendor and multi-disease cardiac segmentation: The m& ms challenge, *IEEE Transactions on Medical Imaging* 40 (12) (2021) 3543–3554. doi:10.1109/tmi.2021.3090082.
URL <http://dx.doi.org/10.1109/TMI.2021.3090082>
- [23] T. N. Alnasser, L. Abdulaal, A. Maiter, M. Sharkey, K. Dwivedi, M. Salehi, P. Garg, A. J. Swift, S. Alabed, Advancements in cardiac structures segmentation: a comprehensive systematic review of deep learning in ct imaging, *Frontiers in Cardiovascular Medicine* 11 (Jan. 2024). doi:10.3389/fcvm.2024.1323461.
URL <http://dx.doi.org/10.3389/fcvm.2024.1323461>
- [24] J. Ren, H. Sun, H. Zhao, H. Gao, C. Maclellan, S. Zhao, X. Luo, Effective extraction of ventricles and myocardium objects from cardiac magnetic resonance images with a multi-task learning u-net, *Pattern Recognition Letters* 155 (2022) 165–170. doi:10.1016/j.patrec.2021.10.025.
URL <http://dx.doi.org/10.1016/j.patrec.2021.10.025>
- [25] J. Xing, S. Wang, K. C. Bilchick, A. R. Patel, M. Zhang, Joint deep learning for improved myocardial scar detection from cardiac mri, in: 2023 IEEE 20th International Symposium on Biomedical Imaging

- (ISBI), IEEE, 2023, p. 1–5. doi:10.1109/isbi53787.2023.10230541.
URL <http://dx.doi.org/10.1109/isbi53787.2023.10230541>
- [26] T. Yao, N. St. Clair, G. F. Miller, A. L. Dorfman, M. A. Fogel, S. Ghelani, R. Krishnamurthy, C. Z. Lam, M. Quail, J. D. Robinson, D. Schidlow, T. C. Slesnick, J. Weigand, J. A. Steeden, R. H. Rathod, V. Muthurangu, A deep learning pipeline for assessing ventricular volumes from a cardiac mri registry of patients with single ventricle physiology, *Radiology: Artificial Intelligence* 6 (1) (Jan. 2024). doi:10.1148/ryai.230132.
URL <http://dx.doi.org/10.1148/ryai.230132>
- [27] X. Sun, L.-H. Cheng, S. Plein, P. Garg, R. J. van der Geest, Deep learning based automated left ventricle segmentation and flow quantification in 4d flow cardiac mri, *Journal of Cardiovascular Magnetic Resonance* 26 (1) (2024) 100003. doi:10.1016/j.jocmr.2023.100003.
URL <http://dx.doi.org/10.1016/j.jocmr.2023.100003>
- [28] E. Elizar, R. Muharar, M. A. Zulkifley, Desppnet: A multiscale deep learning model for cardiac segmentation, *Diagnostics* 14 (24) (2024) 2820. doi:10.3390/diagnostics14242820.
URL <http://dx.doi.org/10.3390/diagnostics14242820>
- [29] A. S. Ba Mahel, M. S. A. M. Al-Gaashani, F. M. G. Alotaibi, R. I. Alkanhel, Resst-seunet++: Deep model for accurate segmentation of left ventricle and myocardium in magnetic resonance imaging (mri) images, *Bioengineering* 12 (6) (2025) 665. doi:10.3390/bioengineering12060665.
URL <http://dx.doi.org/10.3390/bioengineering12060665>
- [30] B. Kundu, B. Khanal, R. Simon, C. A. Linte, Assessing the performance of the dinov2 self-supervised learning vision transformer model for the segmentation of the left atrium from mri images (2024). doi:10.48550/ARXIV.2411.09598.
URL <https://arxiv.org/abs/2411.09598>