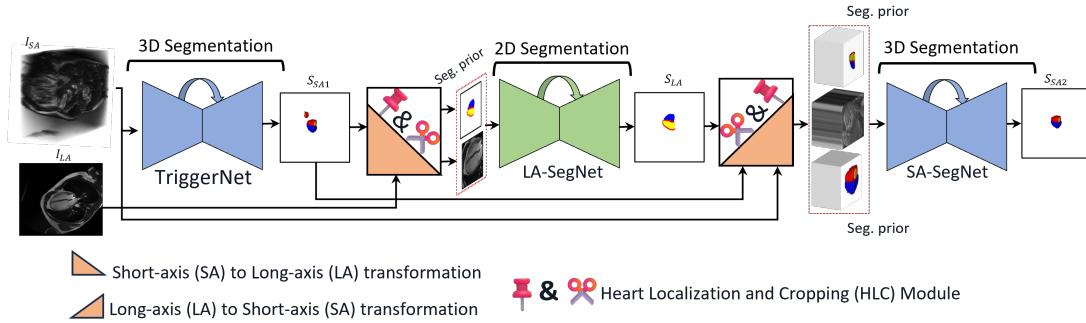# Graphical Abstract

## Multi-view Cardiac Image Segmentation via Trans-Dimensional Priors

Abbas Khan, Muhammad Asad, Martin Benning, Caroline Roney, Gregory Slabaugh

# Highlights

**Multi-view Cardiac Image Segmentation via Trans-Dimensional Priors**

Abbas Khan, Muhammad Asad, Martin Benning, Caroline Roney, Gregory Slabaugh

- We propose a sequential 3D-to-2D-to-3D approach for multi-view cardiac image segmentation by effectively utilizing the trans-dimensional segmentation priors (TDSP), which transform a segmentation from one view into another and serve as guidance.

- The TDSP provides a robust anatomical reference at the network's input and encourages the network to produce anatomically plausible segmentation maps.

- We also introduce a Heart Localization and Cropping (HLC) module to focus the segmentation on the heart region only. This strategy reduces the computation for the second and third-stage segmentation network and eliminates false positive predictions.

- Extensive experiments are conducted to showcase the efficacy of the proposed pipeline utilizing the HLC module and TDSP, where our proposed method outperforms the state-of-the-art as well as methods on the M&Ms-2 challenge leaderboard.

# Multi-view Cardiac Image Segmentation via Trans-Dimensional Priors

Abbas Khan[a,b], Muhammad Asad[b,c], Martin Benning[d], Caroline Roney[b,e], Gregory Slabaugh[a,b]

[a]*School of Electronic Engineering and Computer Science Queen Mary University of London UK*
[b]*Queen Mary's Digital Environment Research Institute (DERI) London UK*
[c]*School of Biomedical Engineering and Imaging Sciences Kings College London UK*
[d]*Department of Computer Science University College London UK*
[e]*School of Engineering and Materials Science Queen Mary University of London UK*

## Abstract

We propose a novel multi-stage trans-dimensional architecture for multi-view cardiac image segmentation. Our method exploits the relationship between long-axis (2D) and short-axis (3D) magnetic resonance (MR) images to perform a sequential 3D-to-2D-to-3D segmentation, segmenting the long-axis and short-axis images. In the first stage, 3D segmentation is performed using the short-axis image, and the prediction is transformed to the long-axis view and used as a segmentation prior in the next stage. In the second step, the heart region is localized and cropped around the segmentation prior using a Heart Localization and Cropping (HLC) module, focusing the subsequent model on the heart region of the image, where a 2D segmentation is performed. Similarly, we transform the long-axis prediction to the short-axis view, localize and crop the heart region and again perform a 3D segmentation to refine the initial short-axis segmentation. We evaluate our proposed method on the Multi-Disease, Multi-View & Multi-Center Right Ventricular Segmentation in Cardiac MRI (M&Ms-2) dataset, where our method outperforms state-of-the-art methods in segmenting cardiac regions of interest in both short-axis and long-axis images. The pre-trained models, source code, and implementation details will be publicly available.

*Keywords:* Cardiac MRI, Image Segmentation, Short-Axis, Long-Axis, Transformation Priors, Sequential Segmentation

## 1. Introduction

Cardiovascular disease is the leading cause of death, with a yearly toll of 23.6 million lives due to heart disease and stroke globally [1]. This underscores the need to identify and treat cardiac disorders. Cardiologists have focused on early diagnosis as part of the clinical workflow [2]. Deep learning architectures have achieved a wide range of competencies for computational cardiac imaging [3],[4],[5] including segmentation [6],[7],[8].

Modern non-invasive medical imaging techniques, such as ultrasound, magnetic resonance imaging (MRI), and computed tomography, are widely used to capture detailed images of the structure and function of the heart and its associated vessels [9]. However, detection of disease and quantification often requires a laborious process of manual segmentation to identify the areas of the anatomy in these scans. Recent advances in artificial intelligence [10] are improving automation to segment a medical image into meaningful areas of interest. In the context of cardiac imaging, areas of interest include the left atrium, right atrium, left ventricle, right ventricle, and myocardium to diagnose different cardiac pathologies. Many image segmentation methods have been proposed, including active shape models [11], active appearance models [12], atlas-based methods [13], convolutional neural network (CNN)-based approaches [14],[15] including those with self-attention-based architectures [16],[17].

Among the successful cardiac image segmentation methods, most rely on a single view, i.e., short-axis (SA) or long-axis (LA), where the segmentation is performed. However, capturing both SA and LA MR images is considered standard practice [18], [19], and segmentation of one view can be utilized to improve the segmentation of the other. Here, we propose a novel framework that performs accurate cardiac image segmentation by transferring the segmentation of one view to guide the segmentation of the other. Our proposed method sequentially utilizes the multi-view images. Despite being based on single encoder-decoder segmentation networks, the proposed pipeline still benefits from multi-view data.

Fig. 1(d) depicts the overall architecture of the proposed pipeline. TriggerNet, which functions as a 3D segmentation model, generates the initial segmentation for the short-axis denoted as $S_{SA1}$. Subsequently, utilizing the transformation parameters for the given volume, $S_{SA1}$ undergoes transformation to the LA view to produce a SA-to-LA map (SA2LAmap), a trans-dimensional segmentation prior. The SA2LAmap and LA image ($I_{LA}$) are

fed as input to the Heart Localization and Cropping (HLC) module, resulting in a cropped $I_{LA}$ and SA2LAmap containing only the heart region. The SA2LAmap is input along with the cropped $I_{LA}$ to the LA-SegNet model that generates a segmentation for the long-axis named $S_{LA}$.

We next refine the short axis segmentation $S_{SA1}$. $S_{LA}$ is transformed to the SA view, resulting in the LA-to-SA map (LA2SAmap), another trans-dimensional segmentation prior. Here, we again use the HLC module to obtain cropped LA2SAmap, SA image ($I_{SA}$), and the TriggerNet output ($S_{SA1}$). Finally, the SA-SegNet utilizes these cropped outputs of the HLC module and generates the final segmentation for the short-axis named $S_{SA2}$.

In our proposed framework, integrating the segmentation from alternate views (SA to LA and LA to SA) acts as a segmentation prior, and provides HLC module guidance to remove the surrounding background regions and improve overall segmentation accuracy for the respective views. This framework enables LA-SegNet and SA-SegNet segmentation to outperform the existing state-of-the-art methods on the Multi-Disease, Multi-View & Multi-Center Right Ventricular Segmentation in Cardiac MRI (M&Ms-2) dataset's challenge [20]. The proposed framework efficiently utilizes the multi-view aspect of the M&Ms-2 dataset, as the challenge provides the images and labels of both views (LA and SA) for each instance, compared to the previous datasets M&Ms [21] and ACDC [22].

## 2. Related Work

This section lists some of the most well-known deep learning-based segmentation architectures, including those unifying the power of CNN and self-attention-based mechanisms. We also detail existing cardiac image segmentation approaches, specifically from the M&Ms-2 dataset's challenge [20] leaderboard and subsequent publications leveraging the dataset. We note that our proposed pipeline can be implemented with any segmentation backbone, provided that the architecture can segment both 2D and 3D views.

UNet [23] revolutionized deep learning-based medical image segmentation by proposing a symmetric encoder-decoder architecture. The encoder part extracts the features from the image, and the decoder reconstructs the segmentation map, while the skip-connections help to propagate information across different stages. The no-new-Net (nnUNet) [24] is built on UNet and proposed an automatically configurable segmentation architecture. It can configure data pre-processing, network design, and post-processing for many medical image segmentation datasets. An overview is provided in Section 3.1.

3

ResUNet [25] is an encoder-decoder architecture based on the UNet model and also incorporates knowledge of residual connections [26], atrous convolutions [27], and pyramid scene parsing (PSP) pooling [28]. Each convolution block is replaced with a residual block to achieve consistent training with the increased network depth, atrous convolutions help increase the receptive field, and PSP pooling enhances the network's performance by including background context information.

Inspired by the emergence of vision transformers [29] in computer vision regimes [30], many hybrid architectures that utilize multi-head self-attention (MHSA) [31] have been proposed for medical image segmentation. TransUNet [32] is a UNet architecture that utilizes both CNN and self-attention. This includes a transformer-based encoder that extracts features from images and a CNN-based decoder that upsamples the encoded features. UTNet [33] is also a hybrid architecture integrating transformer and CNN for medical image segmentation. It proposes a revised MHSA mechanism to reduce the complexity of the model. In addition, a hybrid layer utilizing CNN and revised MHSA is incorporated into the encoder and decoder stages. The Multi-Compound Transformer (MCTrans) [34] aims to combine rich features and semantic structures into multi-scale convolutional features using self-attention. The MCTrans transforms convolutional features as a sequence of tokens to perform intra- and inter-scale self-attention across multiple scales. A multi-view and transformer-based architecture named Transfusion was proposed by [35] to correlate and fuse data coming from SA/LA views. It proposed Divergent Fusion Attention (DiFA), which combines features from different views using multi-scale self-attention. Al Khalil et al. [36] proposed a three-stage approach: firstly, the region of the heart is detected using a regression model; secondly, a GAN-based augmentation technique is used for image synthesis to increase the diversity of the training data for segmentation tasks. More specifically, their approach generates more examples of pathologies to balance instances of pathological and normal cases. Lastly, the late-fusion segmentation approach combined with intensity transformations is utilized to generate the final segmentation map.

Sun et al. [37] utilized labels from the end-diastolic and end-systolic phases through an intensity-based image registration approach. These registered labels increase the size of the training set. Arega et al. [38] relied on the MRI-specific based, intensity, and spatial data augmentation techniques to improve the generalization and robustness of their segmentation models. In [39], a multi-view SA-LA Network was proposed to simultaneously segment

the RV blood pools in both the SA and LA views. It merged the bottleneck features from both the SA and LA and combined the labels of the left ventricle (LV) and myocardium (MYO) to generate a label that aids with contextual information to better segment the right ventricle (RV). Another multi-encoder-decoder network (xUnet) is proposed by [40] to simultaneously process the SA and LA views. It utilizes a pre-processing step where both views are centered and rotated to match their axes. A spatial transformer multi-pass feature pyramid (Tempera) [41] segments the RV in both SA and LA cardiac MR images. Tempera is based on the multi-scale feature pyramid network from [42] and transforms the SA features to LA via a geometric target spatial transformer. InfoTrans [43] proposed a nnUNet-based architecture, where the first 2D-nnUNet segments the LA views and then utilizes the LA prediction to crop the region of interest (ROI) from SA views. The Refined Deep Layer Aggregation (RDLA) [44] proposed a two-stage 2D architecture, using DLA-34 stride-2 network [45] as the backbone. The LA and SA images are segmented independently, followed by a refinement step by utilizing the complementary information of another view along with the images.

Our proposed approach unites the strengths of [41], [42], and [44], as shown Fig. 1 and introduces a sequential approach in which each network benefits from the previous one, without introducing additional parameters compared to other multi-stage approaches such as [36]. InfoTrans [42] performs information transition only from an LA to an SA network and utilizes the transformed SA map (LA2SAmap) only to extract the region of interest (ROI) from original SA images. Compared to this method, we introduce the transformation from LA to SA to obtain the LA2SAmap and from SA to LA to obtain the SA2LAmap. Additionally, we also take advantage of these transformation maps and utilize them as segmentation priors for anatomically plausible predictions. The Tempera [41] architecture only transforms SA to LA and segments each view twice using a two-stage methodology. Compared to Tempera, the proposed pipeline performs a complete two-way transformation, using the SA2LAmap and LA2SAmap to achieve 3D-to-2D-to-3D segmentation. Furthermore, Tempera does not utilize the transformed map to localize and crop the heart region in intensity images, as proposed by the HLC module. Hence, both stages in Tempera utilize full-scale images. Compared to [44], which is implemented using four 2D networks, the proposed refinement strategy uses three networks, where we use a 3D network for SA images to effectively exploit the 3D spatial context and a single 2D

network can achieve state-of-the-art performance using the transformation from TriggerNet as a segmentation prior. In addition, in RDLA architectures, the second-stage networks have the same complexity as the first-stage networks without taking advantage of alignment and first-stage predictions. However, in the proposed setting, the second-stage networks, i.e., LA-SegNet and SA-SegNet, receive images with a lower in-plane spatial resolution (images containing only heart regions) and have fewer encoder-decoder stages, which results in reduced computational complexity, as shown in Fig. 2 (further details in Section 5). Al Khalil et al. [36] employed a regression-based neural network as a crucial component for heart region detection within SA and LA images. While practical, this regression-based neural network brings additional trainable parameters, thereby increasing the complexity of the overall pipeline. In contrast, our proposed HLC module achieves the exact heart region localization and cropping task as their regression-based model without introducing additional trainable parameters. Instead, the HLC module utilizes segmentation maps from the previous network and crops the heart regions, thus offering a streamlined alternative to the regression-based approach. Our contributions can be summarized as follows:

1. We propose a sequential 3D-to-2D-to-3D approach for multi-view cardiac image segmentation by effectively utilizing the trans-dimensional segmentation priors (TDSP), which transform a segmentation from one view into another and serve as guidance.

2. The TDSP provides a robust anatomical reference at the network's input and encourages the network to produce anatomically plausible segmentation maps.

3. We leverage the TDSP and introduce a Heart Localization and Cropping (HLC) module to focus the segmentation on the heart region only. This strategy reduces the computation for the second and third-stage segmentation network and eliminates false positive predictions.

4. Extensive experiments are conducted to showcase the efficacy of the proposed pipeline utilizing the HLC module and TDSP, where our proposed method outperforms the state-of-the-art as well as methods on the M&Ms-2 challenge leaderboard.

## 3. Proposed Framework

Fig. 1(d) depicts our proposed framework, where the pipeline starts with trigger network (TriggerNet), followed by the transformation of its predictions
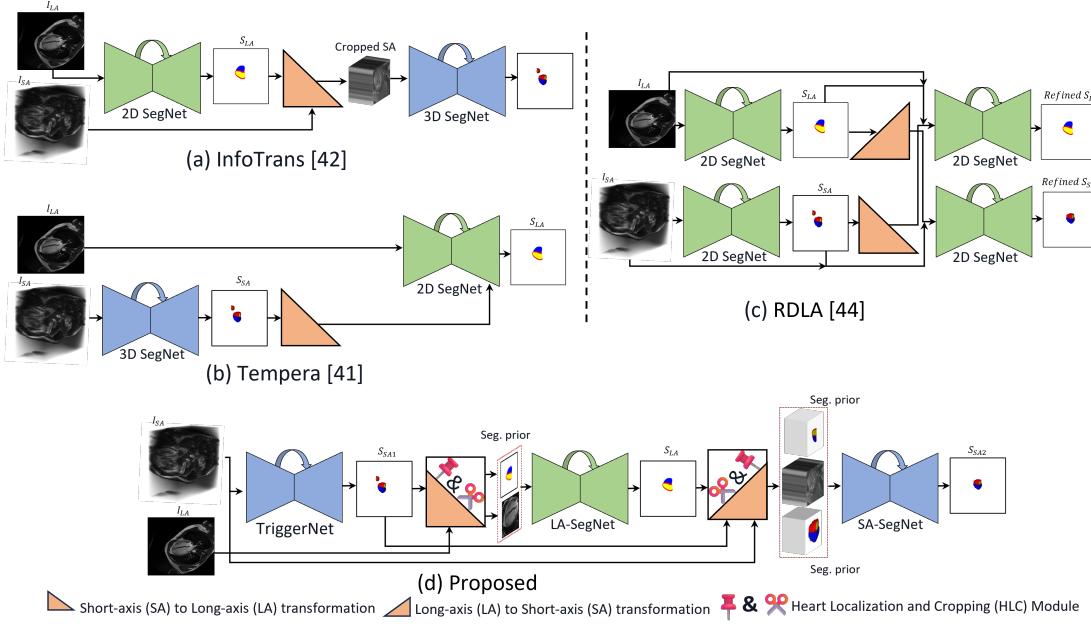
6

Figure 1: Overview of other related methods and our proposed pipeline. The proposed framework (d) segments $I_{SA}$ using TriggerNet to produce $S_{SA1}$. It then generates $S_{LA}$ using LA-SegNet and refines $S_{SA2}$ using SA-SegNet, along with relevant segmentation prior and image. Compared to (a), (b), and (c), our method (d) performs a two-way transformation (LA2SA and SA2LA) along with the utilisation of transformed maps as guidance for the HLC module.

$S_{SA1}$ to the LA view, resulting in SA2LAMap. As a pre-processing step for SA to LA transformation, the header information of the original SA image ($I_{SA}$) is applied to $S_{SA1}$ to ensure the matching of all properties of the $I_{SA}$ and $S_{SA1}$. The SA2LAMap is used to remove the unrelated non-cardiac areas of the original LA image ($I_{LA}$) utilizing the HLC module and as input to the LA-SegNet through concatenation with the cropped $I_{LA}$ as a segmentation prior. The output from LA-SegNet is restored to its original size, followed by copying all the metadata information from the $I_{LA}$ to preserve the header information of the $I_{LA}$ in $S_{LA}$. The final segmentation for LA ($S_{LA}$) is transformed to the SA view using the $I_{SA}$ to obtain LA2SAMap. Following the same process, the HLC module utilizes the LA2SAMap and $I_{SA}$ to localize and crop the heart in a full-scale image. Here, we also cropped and concatenated the $S_{SA1}$ from the TriggerNet (further details are provided in the ablation in Section 5). Finally, the $S_{SA2}$ is restored to its original size.

All three networks (TriggerNet, LA-SegNet, and SA-SegNet) are trained

independently. The downstream tasks, such as LA-2-SA transformations and vice versa, and the HLC module are applied when the previous network outputs are available. However, at inference, the full framework is used sequentially to perform 3D-to-2D-to-3D segmentation output results for both the LA and SA images.

The following subsections will list the details of each step, segmentation networks, HLC module, and transformation process in the proposed pipeline shown in Fig. 1.

### 3.1. Segmentation Networks

Segmentation networks used within the proposed framework, i.e., Trigger-Net, LA-SegNet, and SA-SegNet, are implemented using nnUNet [24]. The nnUNet is built upon the original U-Net architecture with modifications and improvements adopted for medical image segmentation tasks. It stands out from UNet due to its ability to configure the architecture and hyperparameters automatically during training. Moreover, nnUNet achieved the best results in the challenge cohort of M&Ms-2 [20], and the proposed pipeline is built using its architecture as a baseline. The nnUNet architecture has three major components: (i) Encoder, (ii) Decoder, and (iii) Skip-Connections.

The encoder extracts features from the input data by gradually increas-
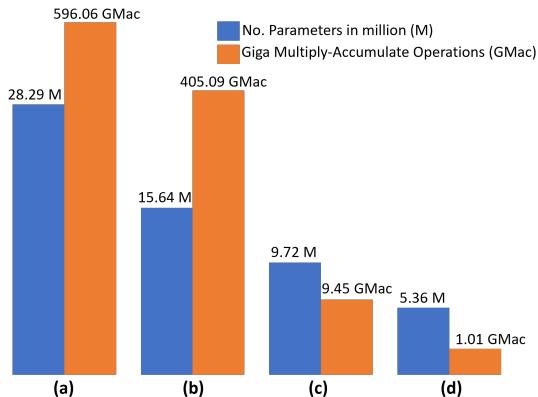


Figure 2: Comparison of segmentation network complexities regarding the number of parameters and Multiply-Accumulate (MAC) operations with and without using the HLC module. (a) SA-SegNet without HLC Module (TriggerNet), and (Depth, Height, Width = 64,192,192), (b) SA-SegNet With HLC Module (SA-SegNet) and (Depth, Height, Width = 112,128,112), (c) LA-SegNet without HLC Module and (Height, Width = 384,384), and (d) LA-SegNet with HLC Module and (Height, Width = 128,128).

8

ing the number of features while reducing the spatial dimensions as it goes deeper. Each encoder block consists of two consecutive convolutions with a kernel size 3 (3×3 for 2D/LA, and 3×3×3 for 3D/SA network). To reduce the spatial resolution, the features are again convolved with a kernel size of 3 and stride of 2. Each convolutional layer is followed by LeakyReLU activation and instance normalization.

The decoder reconstructs the segmentation map by progressively increasing the spatial dimension and reducing the number of features from the bottleneck layer. Each decoder block has two consecutive convolutions with a kernel size of 3, followed by a transpose convolution layer with a kernel size of 2 and stride 2. Similar to the encoder, each convolution layer is followed by LeakyReLU activation and instance normalization. The final convolution layer utilizes a sigmoid activation function with four kernels of size 1, where each kernel generates segmentation output for four classes, i.e., MYO, LV, RV, and background.

The skip-connections have been shown to improve segmentation methods [46] for medical image segmentation tasks, and hence, we utilize skip connections in our proposed architecture. These skip-connections copy and concatenate the features from the contracting path from the encoder to the expanding path in the decoder for a better gradient flow during backpropagation and to recover the lost spatial information.

The number of encoder-decoder blocks can be different for each segmentation network shown in Fig. 1(d), and the corresponding computational complexity in Fig. 2, depending upon the spatial dimension of the input data. The TriggerNet gets spatial dimension images of $64 \times 192 \times 192$, and it has five downsampling (encoder's block) and upsampling blocks (decoder's blocks). The LA-SegNet network gets the cropped input of spatial size $128 \times 128$ and has the four encoder and corresponding decoder blocks. The spatial dimension of input data for SA-SegNet is $112 \times 128 \times 112$, and it also has four blocks for downsampling and upsampling. We also trained the LA-SegNet network on LA images without utilizing the HLC module for the ablation studies mentioned in Section 5. For this network, the nnUNet configures six encoder-decoder stages due to its large spatial dimension of $384 \times 384$.

### 3.2. Heart Localization and Cropping (HLC) Module

The foreground-background imbalance of pixels has been a fundamental issue for accurately segmenting medical images [47]. The foreground pixels
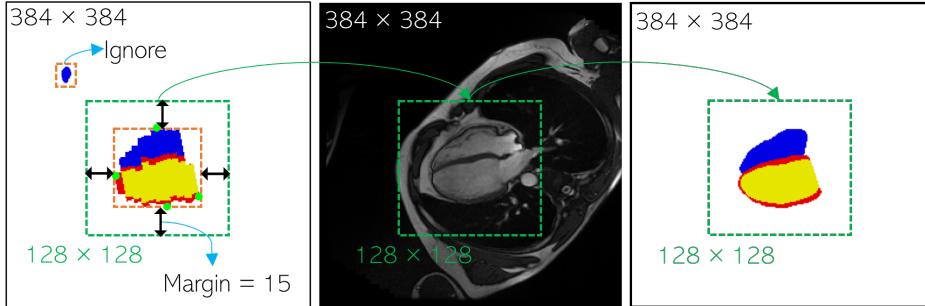
9

Figure 3: Proposed Heart Localization and Cropping (HLC) module. The heart region is localized and cropped in both intensity and label images.

occupy a smaller proportion of the image than the background objects. It can significantly degrade the segmentation performance by forcing the model to focus more on the background pixels due to their majority compared to the foreground pixels [48]. A common way of solving this issue is by designing the loss functions for segmentation, which can weigh the foreground pixels more than the background pixels [49], [50]. However, in the proposed framework, this problem is solved intrinsically to some extent, as we cropped the original full-scale image using the segmentation prior, resulting in a reduction of background search space. The segmentation networks utilizing full-scale images frequently confuse the background and cardiac tissue, resulting in considerable false positives. Using the proposed HLC module, the network only focuses on a smaller area containing heart tissues. Moreover, as the HLC module reduces the spatial resolution, the next stage network will run on a low-resolution image, leading to reduced computational complexity. The HLC module is implemented by extracting the heart region containing the LV blood pool (LV), RV blood pool (RV), and left ventricular MYO from the original full-scale images. This can be achieved using the SA2LAmap for $I_{LA}$ and either LA2SAmap or $S_{SA1}$ for $I_{SA}$.

To find the bounding box across three labeled regions (LV-blood pool, RV-blood pool, and left ventricular MYO) within a segmentation prior, we'll use transformation maps (SA2LAmap, LA2SAmap, or $S_{SA1}$) as binary masks, where nonzero values represent the regions of interest. To ensure that the cropped region occupies the entire region and does not miss any pixels of the foreground regions, we defined a parameter named *margin*, which can help to safeguard the edges and provide a margin that assists in preserving the

10

entire heart in the cropped region, as shown in Fig. 3.

We further confirmed that the obtained cropped region from the HLC module safely encloses the entire heart by cropping the ground truth, restoring the original size, and finding the Dice score between the original size and the restored ground truth. We ensured a Dice score of 1 and Hausdorff Distance (HD) of 0 for the resized ground truth segmentation against the original size ground truth. Empirically, we found that a *margin* of 15 pixels perfectly fits this purpose for all images in the training, validation, and test sets. We applied the same protocols to crop the intensity image and bring back the prediction from the cropped image to the original spatial dimension.

### 3.3. The Transformation Process

The M&Ms-2 dataset is novel in terms of providing the images/labels of the LA view along with the SA view to give detailed information for the apical and basal slices of short-axis views [20]. We have utilized this information more efficiently in the proposed framework and generated the transformations for each axis. More specifically, we transformed one view's physical coordinates into the other view's image coordinates system and vice versa.

The trans-dimensional segmentaton prior **SA2LAmap** is obtained using the pseudo-code shown in Algorithm 1. The prediction from TriggerNet $S_{SA1}$ has different metadata information than the original $I_{SA}$. This metadata information is essential for the conversion between physical coordinates and image coordinate systems and includes additional information like image orientation, voxel size, and origin. We used the CopyInformation function from SimpleITK to inherit all the relevant metadata from the original SA image to the $S_{SA1}$. This ensured that $S_{SA1}$ was

---
**Algorithm 1** Transformation of the SA segmentation ($S_{SA1}$) to the viewpoint of the LA image ($I_{LA}$) using $T_{SA \to LA}$.

---
**Input** $I_{LA}, I_{SA}, S_{SA1}, T_{SA \to LA}$
**Output** $SA2LAmap$

    Initialize output to zero, $SA2LAmap = 0$
    **for** each point ($\mathbf{p}$) in $I_{LA}$ **do**
        Use $T_{SA \to LA}$ to transform $\mathbf{p}$ into $S_{SA1}$, producing $\mathbf{q}$
        $SA2LAmap(\mathbf{p}) = S_{SA1}(\mathbf{q})$
    **end for**
**Return** $SA2LAmap$

---

 

---
**Algorithm 2** Transformation of the LA segmentation ($S_{LA}$) to the viewpoint of the SA image ($I_{SA}$) using $T_{LA \to SA}$.

---
**Input** $I_{LA}, I_{SA}, S_{LA}, T_{LA \to SA}$
**Output** $LA2SAmap$

    Initialize output to zero, $LA2SAmap = 0$
    **for** each point $\mathbf{p}$ in $I_{SA}$ **do**
        Use $T_{LA \to SA}$ to transform $\mathbf{p}$ into $S_{LA}$, producing $\mathbf{q}$
        $LA2SAmap(\mathbf{p}) = S_{LA}(\mathbf{q})$
    **end for**
**Return** $LA2SAmap$

---

aligned correctly with the original $I_{SA}$ and could be used to transform the coordinate systems with the original $I_{LA}$.

The physical coordinates of the $I_{LA}$ are obtained from the image coordinate system, followed by finding the corresponding index in the $S_{SA1}$ and mapping these physical coordinates to the $S_{SA1}$ image coordinate system. Finally, it assigns the voxel value from the $S_{SA1}$ to the corresponding voxel in the SA2LAmap. In this way, the SA2LAmap is populated with all transformed values from $S_{SA1}$ to the SA2LAmap.

The trans-dimensional segmentation prior **LA2SAmap** is generated similarly, as shown in Algorithm 2. In this case, the metadata information from the original $I_{LA}$ is copied to the $S_{LA}$ from LA-SegNet. The physical coordinates of the $I_{SA}$ are extracted from its image coordinates system, followed by transforming this physical point to an index in the coordinate system of

LA2SAmap. Finally, each voxel value from the $S_{LA}$ is copied to the corresponding location in the LA2SAmap.

*3.4. Implementation Details*

The proposed architecture is implemented using a single NVidia A100 GPU with 40GB RAM. The SA and LA MRI scans are resampled to a voxel size of $1\times1\times1\ mm^3$. Dice loss and cross-entropy loss are used as loss functions to train the segmentation networks. Stochastic gradient descent is used as an optimizer with an initial learning rate of 0.01 and a Nesterov momentum of 0.99. We utilize a polynomial learning rate scheduler [51] with a weight decay of 0.0005 to decrease the learning rate after each training epoch. All networks, i.e., TriggerNet, LA-SegNet, SA-SegNet, and the nnUNet baseline, are trained independently for 1000 epochs (nnUNet default), where each epoch has 250 training iterations. The pre and post-processing steps, such as LA-2-SA transformation and vice versa, and the heart localization and cropping are performed in succession after the previous network predictions are available. All three segmentation networks and their respective pre and post-processing steps are carried out sequentially in the inference phase.

Different data augmentation techniques are applied during training to allow the networks to see a stream of distinct examples. Spatial transformations, including random rotation, scaling, and mirroring, provide distinct spatial perspectives from which the model can learn. Intensity adjustments, such as random brightness, contrast, and gamma variations, ensure the model's adaptability to varying acquisition settings. Additionally, additive zero-mean Gaussian noise is utilized to enhance stochasticity, and blurring techniques, such as Gaussian blur, contribute to the model's robustness against variations in image quality.

## 4. Dataset

The Multi-Disease, Multi-View, and Multi-Center Right Ventricular Segmentation challenge (M&Ms-2) was introduced in MICCAI 2021. The challenge focused on segmenting RV blood pools across cardiac imaging of multiple views and centers [21],[20]. The data includes diverse images from three clinical centers in Spain utilizing nine scanners from three vendors, including Siemens, General Electric, and Philips. It includes instances having various LV and RV pathologies as well as healthy subjects. The labels are provided for three regions of interest, including (i) LV blood pools, (ii) RV blood pools, and (iii) left ventricular MYO. It contains 360 instances from two cardiac cy-

cles, specifically the end-diastolic and end-systolic phases. The subjects are divided sequentially into 160 for training, 40 for validation, and 160 for testing, such that different patients are in each split. The validation and test set also includes patients with pathologies not included in the training set. For each individual, both SA and LA MR images are provided, having SA and LA 4-chamber views.

## 5. Ablation Studies

We study the effectiveness of our algorithmic design via different ablation studies. In particular, we evaluate the effect of utilizing the SA2LAmap and LA2SAmap/$S_{SA1}$ as a segmentation prior and the HLC module under different settings.

The SA2LAmap can be used in two ways to boost the network's performance: (i) *Localization and Cropping Guide for HLC module:* To localize and crop the heart in original full-scale $I_{LA}$, and (ii) *Segmentation Prior:* As a Segmentation prior concatenated to the $I_{LA}$. Table 1 lists the results of these experiments.

Row 1 represents the baseline where we do not use the SA2LAmap. Rows 2 and 3 show the results where we are using the SA2LAmap for the HLC module or as a segmentation prior, and row 4 depicts the performance where we first localized and cropped the heart in the $I_{LA}$ using the SA2LAmap and then concatenated the cropped image with the SA2LAmap. It can be seen that using either of the techniques (HLC module or segmentation prior) improves the results, which can be further validated with Fig. 4. In Fig. 4, rows (a) and (b) depict that using the SA2LAmap either for the HLC module or as a segmentation prior can help remove the outliers (false positive predictions outside the area of interest). Rows (c) and (d) further confirm the usage of segmentation prior. The segmentation prior aids the network

Table 1: Ablation study to evaluate the effect of utilizing the Segmentation Prior and Heart Localization and Cropping (HLC) module. The first row also represents the baseline nnUNet results for LA.

| HLC Module | SA2LAmap As Segmentation Prior | Dice Score LA ↑ | | | HD (mm) LA ↓ | | |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| | | LV | RV | MYO | LV | RV | MYO |
| ✗ | ✗ | 0.94 | 0.90 | 0.86 | 5.91 | 6.61 | 5.98 |
| ✓ | ✗ | 0.95 | 0.91 | 0.87 | 3.81 | 4.90 | 3.08 |
| ✗ | ✓ | 0.95 | 0.92 | 0.86 | 3.26 | 4.35 | 2.73 |
| ✓ | ✓ | **0.96** | **0.93** | **0.88** | **2.81** | **3.80** | **2.51** |

14

in generating anatomically plausible results. Thus, the SA2LAmap obtained from the $S_{SA1}$ helps to convert the invalid LA segmentation into close but correct shapes.

Table 2 showcases the ablation results for SA segmentation under various settings. Row 1 uses neither the HLC module nor the segmentation priors and lists the results of the TriggerNet. Row 2 lists the results of SA-SegNet if we only utilize the predictions from the TriggerNet to localize and crop the heart in the original full-scale image and use it as a concatenated form
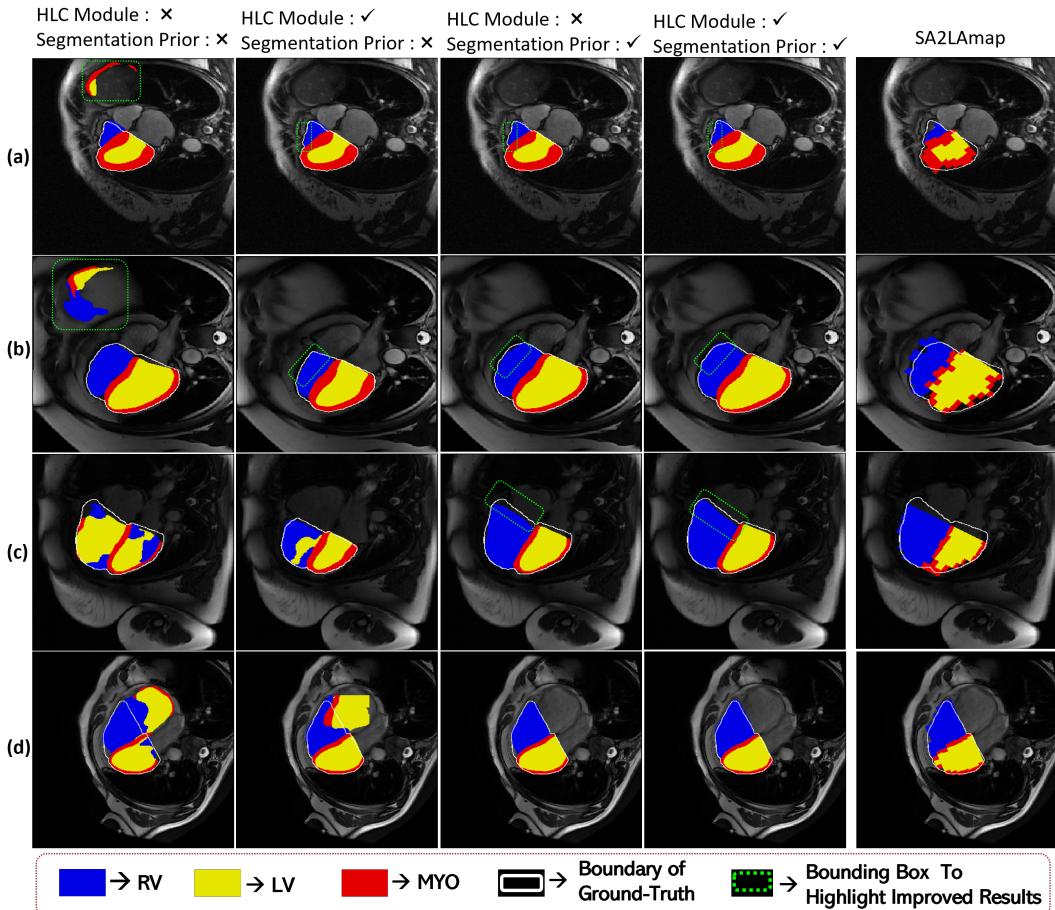


Figure 4: Comparison of visual results under different settings for LA segmentation. From the left, the first column is the baseline results, and the second and third columns utilize HLC module and segmentation priors, respectively. The fourth column is the best results using both. The fifth column represents the SA2LAmap.

Table 2: Ablation study to evaluate the effect of utilizing the Segmentation Priors and Heart Localization and Cropping (HLC) module. The first row also represents the baseline nnUNet results for SA.

| HLC Module | $S_{SA1}$ As Segmentation Prior | LA2SAmap As Segmentation Prior | Dice Score LA ↑ | | | HD (mm) LA ↓ | | |
|---|---|---|---|---|---|---|---|---|
| | | | LV | RV | MYO | LV | RV | MYO |
| ✗ | ✗ | ✗ | 0.924 | 0.888 | 0.842 | 9.137 | 8.662 | 6.181 |
| ✓ | ✓ | ✗ | 0.938 | 0.912 | **0.864** | 3.809 | 5.027 | 2.671 |
| ✓ | ✓ | ✓ | **0.939** | **0.918** | 0.863 | **3.616** | **4.496** | **2.666** |

along with the image. For row 3, we also utilized the LA2SAmap. Using the LA2SAmap as segmentation prior improves all metrics except the MYO Dice score; here, we argue that the collective structure segmentation of MYO has already been improved in terms of structure variance by $S_{SA1}$ as segmentation prior, and LA2SAmap helps further eliminate errors in the boundary regions of MYO as indicated by the improved HD score.

Fig. 5 depicts the visual results of this ablation study. Row (a) and (b) validate that each segmentation prior ($S_{SA1}$ or LA2SAmap) helps to remove the outliers. However, the last two rows, (c) and (d), further confirm that using the LA-to-SA transformation as a segmentation prior can further improve the segmentation. Here, we argue that in the cases where the Trigger-Net cannot produce the accurate segmentation map for any of the regions of interest, then LA-SegNet might have already segmented that instance accurately with valid anatomical shape prior, and we can transfer that knowledge to LA2SAmap and is used it as segmentation prior for the SA-SegNet.

Fig. 2 further advocates that the HLC module not only improves the segmentation performance but also decreases the overall computational requirements of the proposed pipeline. In Fig. 2, all the networks using the HLC module to localize the heart region have fewer encoder-decoder stages than those utilizing full-scale images, resulting in fewer parameters and Giga Multiply-Accumulate Operations (GMac). This is because if the in-plane spatial resolution size for input images is $384 \times 384$, to get bottleneck features of in-plane spatial dimension $= 6 \times 6$, the encoder needs to compress the features six times; however, if the heart is cropped and localized using HLC module, it will produce the images of in-plane spatial dimension $= 192 \times 192$, and the encoder will perform the feature compression in five instances to obtain bottleneck features with in-plane spatial dimension of $6 \times 6$.

## 6. Results and Discussion

We compared the results of the proposed approach with several state-of-the-art architectures using a five-fold cross-validation method as well as comparing against the challenge leaderboard using the provided validation and test set.
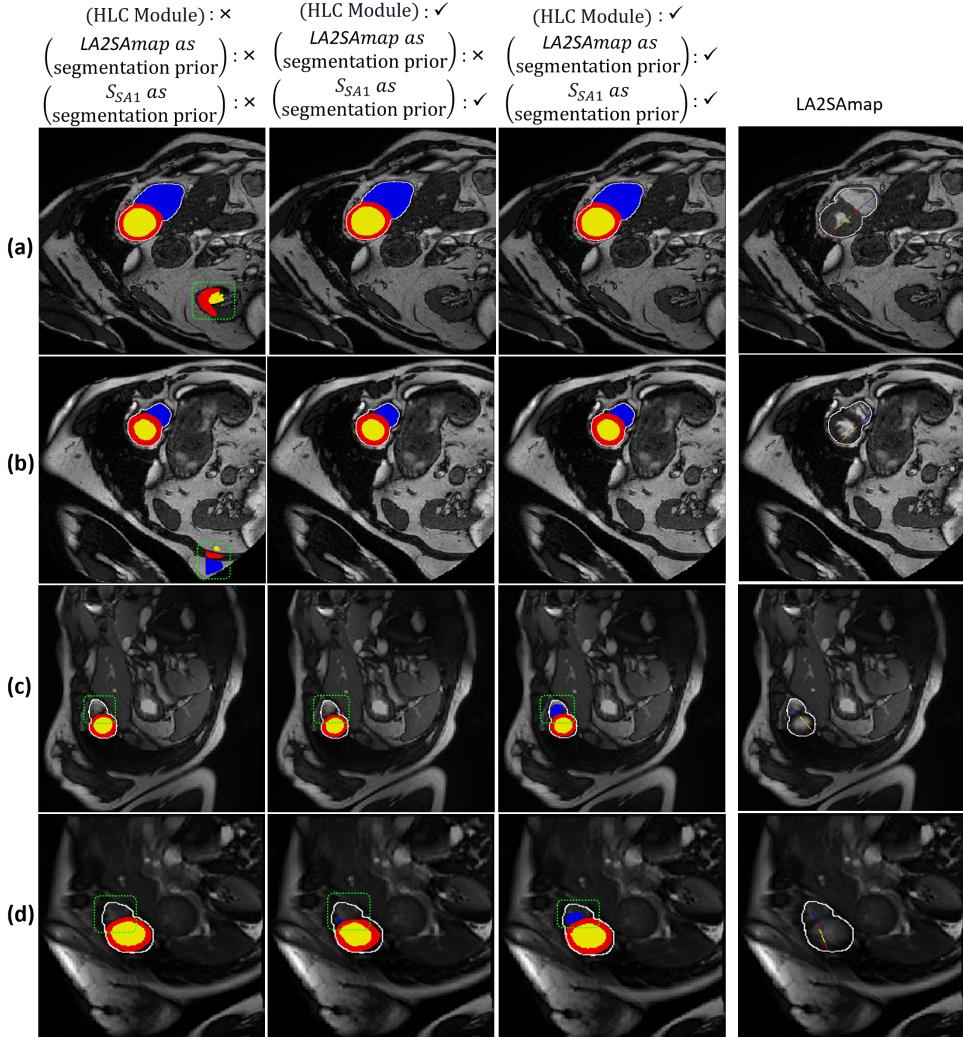


Figure 5: Comparison of visual results under different settings for SA segmentation. From the left, the first column is the baseline results, and the second utilizes the HLC module and $S_{SA1}$ as segmentation priors. The third column shows the best results using both segmentation priors. The fourth column represents the LA2SAmap.
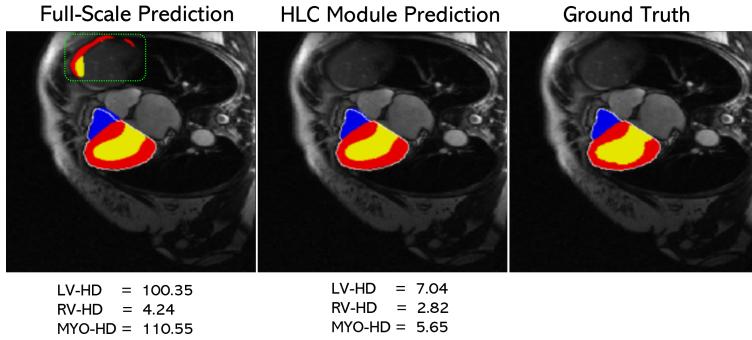
Figure 6: Visual and quantitative analysis of how the outliers contribute to the higher HD, and comparison of results obtained from methods utilizing the full-scale images vs the proposed HLC module.

Table 3 lists the comparison of the results with other segmentation architectures in terms of the Dice score and HD on a five-fold cross-validation split. The proposed method is able to produce robust segmentation across both SA and LA views. A substantial improvement can be seen in terms of reducing the HD score. This gain in performance comes from the fact that the proposed HLC module can remove the outliers efficiently, resulting in an accurate segmentation map compared to the full-scale image methods. Fig. 6 further illustrates this effect. Compared to full-scale prediction, the segmentation map generated by utilizing the proposed HLC module eliminates all the erroneous predictions outside the region of interest. We can further observe the variation in HD score for the predictions with and without the outliers. It can be seen that the difference in HD score for MYO and LV is very large compared to the RV because of the false positives.

Table 4 provides results compared to the challenge leaderboard on the validation set and the test set subjects for RV segmentation. This further proves the generalization ability of the proposed approach to the unseen pathologies, as the validation and test set has two pathologies (Tricuspidal Regurgitation and Congenital Arrhythmogenesis) not available in the training set.

We also compared the proposed approach with other multi-stage segmentation methods, such as those mentioned in [36], and Table 5 lists the comparison results. Similar to [36], the comparison is performed on the 200 subjects by combining the validation (40) and testing (160) examples, while the algorithm's development occurs solely based on the training data. The proposed method outperforms [36] for most of the evaluation metrics.

18

Table 3: Comparison to state-of-the-art methods for long/short-axis MRI on five-fold cross-validation split on the M&MS-2 challenge dataset. The best and second results are in **bold** and <u>underline</u>, respectively. Methods indicated with ∗ use multi-view inputs.

| Methods | Dice(%)-Short-axis ↑ | | | | HD (mm)-Short-axis ↓ | | | |
|---|---|---|---|---|---|---|---|---|
| | LV | RV | MYO | Avg | LV | RV | MYO | Avg |
| UNet | 87.02 | 88.85 | 79.07 | 84.98 | 13.78 | 12.10 | 12.23 | 12.70 |
| ResUNet | 87.98 | 89.63 | 79.28 | 85.63 | 13.80 | 11.61 | 12.09 | 12.50 |
| DLA | 87.27 | 89.88 | 80.23 | 86.12 | 13.25 | 10.84 | 12.31 | 12.13 |
| InfoTrans* | 88.24 | 90.41 | 80.25 | 86.30 | 12.41 | 10.98 | 12.83 | 12.07 |
| rDLA* | 88.64 | 90.28 | 80.78 | 86.57 | 12.74 | 10.31 | 12.49 | 11.85 |
| TransUNet | 87.91 | 88.69 | 78.67 | 85.09 | 13.80 | 10.29 | 13.43 | 12.51 |
| MCTrans | 88.52 | 89.90 | 80.08 | 86.17 | 12.29 | 9.92 | 13.28 | 11.83 |
| MCTrans* | 87.79 | 89.22 | 79.37 | 85.46 | 11.28 | 9.39 | 13.84 | 11.49 |
| UTNet | 87.52 | 90.57 | 80.20 | 86.10 | 12.03 | 9.78 | 13.72 | 11.84 |
| UTNet* | 87.74 | 90.82 | 80.71 | 86.42 | 11.79 | <u>9.11</u> | 13.41 | 11.44 |
| TransFusion* | <u>89.52</u> | <u>91.75</u> | <u>81.46</u> | <u>87.58</u> | <u>11.31</u> | 9.18 | <u>11.96</u> | <u>10.82</u> |
| **Proposed*** | **93.94** | **91.87** | **86.34** | **91.71** | **3.61** | **4.49** | **2.66** | **3.58** |
| | Dice(%)-Long-axis ↑ | | | | HD (mm)-Long-axis ↓ | | | |
| | LV | RV | MYO | Avg | LV | RV | MYO | Avg |
| UNet | 87.26 | 88.20 | 79.96 | 85.14 | 13.04 | 8.76 | 12.24 | 11.35 |
| ResUNet | 87.61 | 88.41 | 80.12 | 85.38 | 12.72 | 8.39 | 11.28 | 10.80 |
| DLA | 88.37 | 89.38 | 80.35 | 86.03 | 11.74 | 7.04 | 10.79 | 9.86 |
| InfoTrans* | 88.21 | 89.11 | 80.55 | 85.96 | 12.47 | 7.23 | 10.21 | 9.97 |
| rDLA* | 88.71 | 89.71 | 81.05 | 86.49 | 11.12 | 6.83 | 10.42 | 9.46 |
| TransUNet | 87.91 | 88.23 | 79.05 | 85.06 | 12.02 | 8.14 | 11.21 | 10.46 |
| MCTrans | 88.42 | 88.19 | 79.47 | 85.36 | 11.78 | 7.65 | 10.76 | 10.06 |
| MCTrans* | 88.81 | 88.61 | 79.94 | 85.79 | 11.52 | 7.02 | 10.07 | 9.54 |
| UTNet | 86.93 | 89.07 | 80.48 | 85.49 | 11.47 | 6.35 | 10.02 | 9.28 |
| UTNet* | 87.36 | 90.42 | 81.02 | 86.27 | 11.13 | 5.91 | 9.81 | 8.95 |
| TransFusion* | <u>89.78</u> | <u>91.52</u> | <u>81.79</u> | <u>87.70</u> | <u>10.25</u> | <u>5.12</u> | <u>8.69</u> | <u>8.02</u> |
| **Proposed*** | **95.80** | **93.07** | **87.71** | **92.19** | **2.81** | **3.80** | **2.51** | **3.04** |

*6.1. Limitations and future work*

The proposed pipeline is useful in the settings where we have cardiac MRI available for both LA and SA views. Clinically, both views are captured according to cardiac MRI acquisition protocol [18], [19]; however, most of the publicly available datasets only provide short-axis cine MR images and their labels [21]. This unavailability of long-axis MR images can be considered a limitation of the proposed approach for datasets with single-view images only. However, we expect future cardiac MRI datasets to release

Table 4: Performance comparison with the M&Ms-2 challenge leaderboard for RV segmentation.

| Quantitative comparison on validation set | | | | |
|---|---|---|---|---|
| Methods | Dice Score LA ↑ | HD-LA(mm) ↓ | Dice Score SA ↑ | HD-SA(mm) ↓ |
| [37] | 0.922 | 5.35 | 0.925 | 8.90 |
| [52] | 0.922 | 5.59 | 0.924 | 8.85 |
| [43] | 0.920 | 5.34 | 0.922 | 9.47 |
| Proposed | **0.926** | **3.49** | **0.928** | **3.72** |
| Quantitative comparison on test set | | | | |
| Methods | Dice Score LA ↑ | HD-LA(mm) ↓ | Dice Score SA ↑ | HD-SA(mm) ↓ |
| [37] | 0.919 | 6.04 | 0.925 | 10.58 |
| [52] | 0.919 | 6.10 | 0.920 | 9.94 |
| [43] | 0.916 | 6.17 | 0.920 | 10.30 |
| Proposed | **0.928** | **3.91** | **0.927** | **4.01** |

Table 5: Performance comparison with multi-stage segmentation approach [36] evaluated on 200 subjects of validation and test set.

| Methods | Dice(%)-Short-axis ↑ | | | HD (mm)-Short-axis ↓ | | |
|---|---|---|---|---|---|---|
| | LV | RV | MYO | LV | RV | MYO |
| [36] | 0.959 | **0.938** | **0.907** | 6.42 | 8.62 | 9.37 |
| Proposed | **0.963** | 0.928 | 0.870 | **3.43** | **3.87** | **3.62** |
| | Dice(%)-Long-axis ↑ | | | HD (mm)-Long-axis ↓ | | |
| | LV | RV | MYO | LV | RV | MYO |
| [36] | 0.958 | 0.924 | **0.901** | 4.07 | 5.81 | 5.27 |
| Proposed | **0.961** | **0.927** | 0.878 | **2.83** | **3.70** | **2.56** |

more complementary information, such as both views, to take advantage of their relationship. We also encourage the research community to provide a 2-chamber and 3-chamber LA view, further exploiting the multi-view aspect of cardiac MR images.

Finally, we analyze that a medical image is influenced by the anatomical features of the image and the characteristics of the imaging equipment, such as vendor information, scanner type, etc. This inspires our future work, where we will aim to design more robust pipelines to incorporate metadata along with intensity images for segmentation tasks [53]. This will enable the segmentation networks to learn not only the appearance of images but also the specific interdependence of an image structure and image-capturing

device.

## 7. Conclusion

This paper proposes a cardiac image segmentation approach relying on the trans-dimensional segmentation priors between short-axis and long-axis views. We show that the method provides a substantial improvement in the accuracy of segmentation for cardiac images in LV, RV blood pools, and left ventricular MYO. The proposed approach effectively utilizes the relationship between the SA and LA views so that a segmentation in one view informs the segmentation in the other view. The transformed maps are used to localize and crop the heart region in the original full-scale image using the HLC module of the same axis and act as a segmentation prior to the other axis. The HLC module helps to remove the outliers and improves erroneous predictions. The segmentation prior encourages anatomically plausible segmentation maps. Extensive ablation studies are conducted to show the efficacy of proposed techniques, and the results are compared with the existing state-of-the-art methods utilizing the M&Ms-2 dataset.

## References

[1] D. M. Greenfield, J. A. Snowden, Cardiovascular diseases and metabolic syndrome, The EBMT Handbook: Hematopoietic Stem Cell Transplantation and Cellular Therapies (2019) 415–420.

[2] WHO: Cardiovascular diseases (CVDs), https://www.who.int/newsroom/fact-sheets/detail/cardiovascular-diseases-(cvds), [accessed 26-Nov-2023].

[3] H. Moradi, A. Al-Hourani, G. Concilia, F. Khoshmanesh, F. R. Nezami, S. Needham, S. Baratchi, K. Khoshmanesh, Recent developments in modeling, imaging, and monitoring of cardiovascular diseases using machine learning, Biophysical Reviews 15 (1) (2023) 19–33.

[4] M. Kadem, L. Garber, M. Abdelkhalek, B. K. Al-Khazraji, Z. Keshavarz-Motamed, Hemodynamic modeling, medical imaging, and machine learning and their applications to cardiovascular interventions, IEEE Reviews in Biomedical Engineering 16 (2022) 403–423.

[5] A. Fotaki, E. Puyol-Antón, A. Chiribiri, R. Botnar, K. Pushparajah, C. Prieto, Artificial intelligence in cardiac mri: is clinical adoption forthcoming?, Frontiers in Cardiovascular Medicine 8 (2022) 818765.

[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 (2020) 25.

[7] L. Li, W. Ding, L. Huang, X. Zhuang, V. Grau, Multi-modality cardiac image computing: A survey, Medical Image Analysis (2023) 102869.

[8] K. Li, L. Yu, P.-A. Heng, Towards reliable cardiac image segmentation: Assessing image-level and pixel-level segmentation quality via self-reflective references, Medical Image Analysis 78 (2022) 102426.

[9] P. S. Rajiah, C. J. François, T. Leiner, Cardiac mri: state of the art, Radiology 307 (3) (2023) e223008.

[10] S. W. Chen, S. L. Wang, T. F. Ng, H. Ibrahim, Artificial intelligence in cardiovascular imaging, in: Cardiovascular and Coronary Artery Imaging, Elsevier, 2023, pp. 51–72.

[11] M. de Bruijne, B. van Ginneken, M. A. Viergever, W. J. Niessen, Adapting active shape models for 3d segmentation of tubular structures in medical images, in: Information Processing in Medical Imaging: 18th International Conference, IPMI 2003, Ambleside, UK, July 20-25, 2003. Proceedings 18, Springer, 2003, pp. 136–147.

[12] S. C. Mitchell, J. G. Bosch, B. P. Lelieveldt, R. J. Van der Geest, J. H. Reiber, M. Sonka, 3-d active appearance models: segmentation of cardiac mr and ultrasound images, IEEE transactions on medical imaging 21 (9) (2002) 1167–1178.

[13] H. Kirişli, M. Schaap, S. Klein, S.-L. Papadopoulou, M. Bonardi, C.-H. Chen, A. C. Weustink, N. R. Mollet, E.-J. Vonken, R. J. van der Geest, et al., Evaluation of a multi-atlas based method for segmentation of cardiac cta data: a large-scale, multicenter, and multivendor study, Medical physics 37 (12) (2010) 6279–6291.

[14] W. Baccouch, S. Oueslati, B. Solaiman, S. Labidi, A comparative study of cnn and u-net performance for automatic segmentation of medical images: application to cardiac mri, Procedia Computer Science 219 (2023) 1089–1096.

[15] W. Xu, J. Shi, Y. Lin, C. Liu, W. Xie, H. Liu, S. Huang, D. Zhu, L. Su, Y. Huang, et al., Deep learning-based image segmentation model using an mri-based convolutional neural network for physiological evaluation of the heart, Frontiers in Physiology 14 (2023) 351.

[16] X. Huang, Z. Deng, D. Li, X. Yuan, Y. Fu, Missformer: An effective transformer for 2d medical image segmentation, IEEE Transactions on Medical Imaging (2022).

[17] K. Deng, Y. Meng, D. Gao, J. Bridge, Y. Shen, G. Lip, Y. Zhao, Y. Zheng, Transbridge: A lightweight transformer for left ventricle segmentation in echocardiography, in: Simplifying Medical Ultrasound: Second International Workshop, ASMUS 2021, Held in Conjunction with MICCAI 2021, Strasbourg, France, September 27, 2021, Proceedings 2, Springer, 2021, pp. 63–72.

[18] C. M. Kramer, J. Barkhausen, C. Bucciarelli-Ducci, S. D. Flamm, R. J. Kim, E. Nagel, Standardized cardiovascular magnetic resonance imaging (cmr) protocols: 2020 update, Journal of Cardiovascular Magnetic Resonance 22 (1) (2020) 1–18.

[19] S. E. Petersen, P. M. Matthews, J. M. Francis, M. D. Robson, F. Zemrak, R. Boubertakh, A. A. Young, S. Hudson, P. Weale, S. Garratt, et al., Uk biobank's cardiovascular magnetic resonance protocol, Journal of cardiovascular magnetic resonance 18 (1) (2015) 1–7.

[20] C. Martín-Isla, V. M. Campello, C. Izquierdo, K. Kushibar, C. Sendra-Balcells, P. Gkontra, A. Sojoudi, M. J. Fulton, T. W. Arega,

K. Punithakumar, et al., Deep learning segmentation of the right ventricle in cardiac mri: The m&ms challenge, IEEE Journal of Biomedical and Health Informatics (2023).

[21] V. M. Campello, P. Gkontra, C. Izquierdo, C. Martin-Isla, A. Sojoudi, P. M. Full, K. Maier-Hein, Y. Zhang, Z. He, J. Ma, et al., Multi-centre, multi-vendor and multi-disease cardiac segmentation: the m&ms challenge, IEEE Transactions on Medical Imaging 40 (12) (2021) 3543–3554.

[22] O. Bernard, A. Lalande, C. Zotti, F. Cervenansky, X. Yang, P.-A. Heng, I. Cetin, K. Lekadir, O. Camara, M. A. G. Ballester, et al., 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.

[23] O. Ronneberger, P. Fischer, T. Brox, U-net: Convolutional networks for biomedical image segmentation, in: MICCAI 2015: Munich, Germany, October 5-9, 2015, Proceedings, Part III 18, Springer, 2015, pp. 234–241.

[24] F. Isensee, P. F. Jaeger, S. A. Kohl, J. Petersen, K. H. Maier-Hein, nnu-net: a self-configuring method for deep learning-based biomedical image segmentation, Nature methods 18 (2) (2021) 203–211.

[25] F. I. Diakogiannis, F. Waldner, P. Caccetta, C. Wu, Resunet-a: A deep learning framework for semantic segmentation of remotely sensed data, ISPRS Journal of Photogrammetry and Remote Sensing 162 (2020) 94–114.

[26] K. He, X. Zhang, S. Ren, J. Sun, Identity mappings in deep residual networks, in: Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14, Springer, 2016, pp. 630–645.

[27] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, A. L. Yuille, Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs, IEEE transactions on pattern analysis and machine intelligence 40 (4) (2017) 834–848.

[28] H. Zhao, J. Shi, X. Qi, X. Wang, J. Jia, Pyramid scene parsing network, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 2881–2890.

[29] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al., An image is worth 16x16 words: Transformers for image recognition at scale, in: International Conference on Learning Representations, 2020.

[30] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, S. Zagoruyko, End-to-end object detection with transformers, in: European conference on computer vision, Springer, 2020, pp. 213–229.

[31] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, I. Polosukhin, Attention is all you need, NeurIPS 30 (2017).

[32] J. Chen, Y. Lu, Q. Yu, X. Luo, E. Adeli, Y. Wang, L. Lu, A. L. Yuille, Y. Zhou, Transunet: Transformers make strong encoders for medical image segmentation, arXiv preprint arXiv:2102.04306 (2021).

[33] Y. Gao, M. Zhou, D. N. Metaxas, Utnet: a hybrid transformer architecture for medical image segmentation, in: MICCAI 2021: Strasbourg, France, September 27–October 1, 2021, Proceedings, Part III 24, Springer, 2021, pp. 61–71.

[34] Y. Ji, R. Zhang, H. Wang, Z. Li, L. Wu, S. Zhang, P. Luo, Multi-compound transformer for accurate biomedical image segmentation, in: MICCAI 2021: Strasbourg, France, September 27–October 1, 2021, Proceedings, Part I 24, Springer, 2021, pp. 326–336.

[35] D. Liu, Y. Gao, Q. Zhangli, L. Han, X. He, Z. Xia, S. Wen, Q. Chang, Z. Yan, M. Zhou, et al., Transfusion: multi-view divergent fusion for medical image segmentation with transformers, in: MICCAI, Springer, 2022, pp. 485–495.

[36] Y. Al Khalil, S. Amirrajab, C. Lorenz, J. Weese, J. Pluim, M. Breeuwer, Reducing segmentation failures in cardiac mri via late feature fusion and gan-based augmentation, Computers in Biology and Medicine 161 (2023) 106973.

[37] X. Sun, L.-H. Cheng, R. J. van der Geest, Right ventricle segmentation via registration and multi-input modalities in cardiac magnetic resonance imaging from multi-disease, multi-view and multi-center, in: Statistical Atlases and Computational Models of the Heart. Multi-Disease,

Multi-View, and Multi-Center Right Ventricular Segmentation in Cardiac MRI Challenge: 12th International Workshop, STACOM 2021, Held in Conjunction with MICCAI 2021, Strasbourg, France, September 27, 2021, Revised Selected Papers 12, Springer, 2022, pp. 241–249.

[38] T. W. Arega, F. Legrand, S. Bricq, F. Meriaudeau, Using mri-specific data augmentation to enhance the segmentation of right ventricle in multi-disease, multi-center and multi-view cardiac mri, in: International Workshop on Statistical Atlases and Computational Models of the Heart, Springer, 2021, pp. 250–258.

[39] S. Jabbar, S. T. Bukhari, H. Mohy-ud Din, Multi-view sa-la net: A framework for simultaneous segmentation of rv on multi-view cardiac mr images, in: International Workshop on Statistical Atlases and Computational Models of the Heart, Springer, 2021, pp. 277–286.

[40] S. Queirós, Right ventricular segmentation in multi-view cardiac mri using a unified u-net model, in: International Workshop on Statistical Atlases and Computational Models of the Heart, Springer, 2021, pp. 287–295.

[41] C. Galazis, H. Wu, Z. Li, C. Petri, A. A. Bharath, M. Varela, Tempera: Spatial transformer feature pyramid network for cardiac mri segmentation, in: International Workshop on Statistical Atlases and Computational Models of the Heart, Springer, 2021, pp. 268–276.

[42] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, S. Belongie, Feature pyramid networks for object detection, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 2117–2125.

[43] L. Li, W. Ding, L. Huang, X. Zhuang, Right ventricular segmentation from short-and long-axis mris via information transition, in: International Workshop on Statistical Atlases and Computational Models of the Heart, Springer, 2021, pp. 259–267.

[44] D. Liu, Z. Yan, Q. Chang, L. Axel, D. N. Metaxas, Refined deep layer aggregation for multi-disease, multi-view & multi-center cardiac mr segmentation, in: International Workshop on Statistical Atlases and Computational Models of the Heart, Springer, 2021, pp. 315–322.

[45] F. Yu, D. Wang, E. Shelhamer, T. Darrell, Deep layer aggregation, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2018, pp. 2403–2412.

[46] M. Drozdzal, E. Vorontsov, G. Chartrand, S. Kadoury, C. Pal, The importance of skip connections in biomedical image segmentation, in: International Workshop on Deep Learning in Medical Image Analysis, International Workshop on Large-Scale Annotation of Biomedical Data and Expert Label Synthesis, Springer, 2016, pp. 179–187.

[47] A. Braytee, A. Anaissi, M. Naji, A comparative analysis of loss functions for handling foreground-background imbalance in image segmentation, in: International Conference on Neural Information Processing, Springer, 2022, pp. 3–13.

[48] N. Yudistira, M. Kavitha, T. Itabashi, A. H. Iwane, T. Kurita, Prediction of sequential organelles localization under imbalance using a balanced deep u-net, Scientific reports 10 (1) (2020) 2626.

[49] S. S. M. Salehi, D. Erdogmus, A. Gholipour, Tversky loss function for image segmentation using 3d fully convolutional deep networks, in: International workshop on machine learning in medical imaging, Springer, 2017, pp. 379–387.

[50] T.-Y. Lin, P. Goyal, R. Girshick, K. He, P. Dollár, Focal loss for dense object detection, in: Proceedings of the IEEE international conference on computer vision, 2017, pp. 2980–2988.

[51] P. Mishra, K. Sarawadekar, Polynomial learning rate policy with warm restart for deep neural network, in: TENCON 2019-2019 IEEE Region 10 Conference (TENCON), IEEE, 2019, pp. 2087–2092.

[52] T. W. Arega, F. Legrand, S. Bricq, F. Meriaudeau, Using mri-specific data augmentation to enhance the segmentation of right ventricle in multi-disease, multi-center and multi-view cardiac mri, in: 12th International Workshop, STACOM 2021, Held in Conjunction with MICCAI 2021, Strasbourg, France, September 27, 2021, Revised Selected Papers, Springer, 2022, pp. 250–258.

[53] A. Lemay, C. Gros, O. Vincent, Y. Liu, J. P. Cohen, J. Cohen-Adad, Benefits of linear conditioning with metadata for image segmentation (2021).