

MCAL: An Anatomical Knowledge Learning Model for Myocardial Segmentation in 2D Echocardiography

Journal:	<i>Transactions on Ultrasonics, Ferroelectrics, and Frequency Control</i>
Manuscript ID	TUFFC-11452-2021.R1
Manuscript Type:	Papers
Date Submitted by the Author:	09-Feb-2022
Complete List of Authors:	cui, xiaoxiao; Shandong University, School of Information Science and Engineering Zhang, Pengfei; Shandong University Qilu Hospital, Department of Cardiology Li, Yujun; Shandong University, School of Information Science and Engineering Liu, Zhi; Shandong University, School of Information Science and Engineering Xiao, Xiaoyan; Shandong University Qilu Hospital, Department of Nephrology Zhang, Yang; Shandong University Qilu Hospital, Department of Radiology Sun, Longkun; Shandong University Qilu Hospital, Department of Cardiology Cui, Lizhen; Shandong University, Joint SDU-NTU Center for Artificial Intelligence Research (C-FAIR) Yang, Guang; Imperial College London, National Heart and Lung Institute Li, Shuo; Western University, Department of Medical Imaging
Keywords:	Medical Imaging < Medical Ultrasonics:, Medical Signal and Image Processing < Medical Ultrasonics:

MCAL: An Anatomical Knowledge Learning Model for Myocardial Segmentation in 2D Echocardiography

Xiaoxiao Cui, Pengfei Zhang, Yujun Li, Zhi Liu, Xiaoyan Xiao, Yang Zhang, Longkun Sun, Lizhen Cui, Guang Yang, and Shuo Li

Abstract—Segmentation of the left ventricular (LV) myocardium in 2D echocardiography is essential for clinical decision making, especially in geometry measurement and index computation. However, segmenting the myocardium is a time-consuming process as well as challenging due to the fuzzy boundary caused by the low image quality. Previous methods based on deep Convolutional Neural Networks (CNN) employ the ground-truth label as class associations on the pixel-level segmentation, or use label information to regulate the shape of predicted outputs, works limit for effective feature enhancement for 2D echocardiography. We propose a training strategy named multi-constrained aggregate learning (referred as MCAL), which leverages anatomical knowledge learned through ground-truth labels to infer segmented parts and discriminate boundary pixels. The new framework encourages the model to focus on the features in accordance with the learned anatomical representations, and the training objectives incorporate a Boundary Distance Transform Weight (BDTW) to enforce a higher weight value on the boundary region, which helps to improve the segmentation accuracy. The proposed method is built as an end-to-end framework with a top-down, bottom-up architecture with skip convolution fusion blocks, and carried out on two datasets (our dataset and the public CAMUS dataset). The comparison study shows that the proposed network outperforms the other segmentation baseline models, indicating that our method is beneficial for boundary pixels discrimination in segmentation.

Index Terms—Boundary distance transform weight, multi-constrained aggregate learning, myocardial segmentation.

I. INTRODUCTION

ECHOCARDIOGRAPHY is routinely used in the diagnosis and management of cardiovascular disease because it can provide real-time images of a beating heart, combined with its availability and portability [1]. Heart function assessment, such as diastolic analysis, calculation of the cardiac output, and ejection fraction (EF), are key determinants of clinical decisions. And segmentation of the left ventricular (LV) myocardium helps accurate quantification of these indexes in the clinical workflow. Thus developing an automatic approach for

This work was supported in part by the National Natural Science Foundation of China under Grant 91846205, the Joint fund for smart computing of Shandong Natural Science Foundation under Grant ZR2020LZH013, the Major Fundamental Research of Natural Science Foundation of Shandong Province under Grant ZR2019ZD05, the Fundamental Research Funds of Shandong University under Grant 2018JC009, and the General Financial Grant from the China Postdoctoral Science Foundation under Grant 2017M622215.

Xiaoxiao Cui and Pengfei Zhang contributed equally to this paper.

Corresponding author: Yujun Li (liyujun@sdu.edu.cn), Zhi Liu (lizhi@sdu.edu.cn) and Lizhen Cui (clz@sdu.edu.cn).

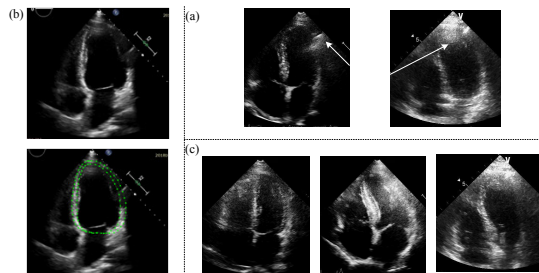


Fig. 1. Typical images extracted from our dataset. (a) Examples of samples with an ambiguous border (Left: Fuzzy chamber border; Right: Fuzzy apical border); (b) Illustration of the annotation of the myocardium. Up: the original ultrasound image; Down: the annotation of the myocardium with a green dot; (c) Different image quality (Left: good image quality; Middle: medium image quality; Right: Poor image quality).

accurate myocardial segmentation liberates radiologists from manual annotation.

Several research works have been performed efficiently on the segmentation in B-mode echocardiography in the past few decades [2]–[4]. With the combination of various feature enhancement modules [5], [6] and different deep network architectures [7]–[10], the ground-truth is applied as a class associate or shape regulation by minimizing the loss function. However, these methods still have scope for improvement. First, methods that focus on feature enhancement to achieve a better result still work limit for echocardiography. Since the limitations of the ultrasound image due to resolution, the presence of speckle noise, and artifacts caused by the complex interaction between the tissue and ultrasound, usually lead to an ambiguous border between the myocardium and chamber (Fig. 1-(a)), making it difficult for an accurate delineation of the myocardium (Fig. 1-(b)) by feature enhancement modules. Second, works that regulate the segmented output with some constraint strategy [9], [10] are much like post-procession and global constraint. But boundary pixels of echocardiography are hard to capture by shape constraints. Because imaging quality varies from subject to subject (Fig. 1-(c)), giving rise to difficulty in capturing the intensity change on the boundary.

To fully use the annotations to address the limitations, we propose a novel training strategy named Multi-Constrained Aggregate Learning (MCAL). Specifically, we force the distributions divergence of the latent features of the input and the ground-truth to be close in the training process. This helps to infer anatomical structure from the deeper layer of the

encoder. Furtherly adjusting the scale and offset of the learned anatomical information by a Feature-wise Linear Modulation (FiLM) [11], the segmented relevance feature information is enhanced. Finally, upon observing that the boundary is hard to detect, we further enhance a higher weight on the border neighborhood pixels by proposing a Boundary Distance Transform Weight (BDTW) for the segmentation loss, which acts as guidance for penalizing the learning process.

The main contributions of our proposed framework are:

- Our method derives segmented-relevance information by narrowing distribution divergence between the latent space of input and label, which exploits anatomical knowledge to guide feature enhancement. FiLM is applied to enhance segmented-relevant features with the guidance of the anatomical information, which restrains the irrelative features under low image quality.
- A novel Boundary Distance Transform Weight (BDTW) is applied to the cross-entropy loss. It forces the network to focus on the boundary region pixels in each training batch and improves the discrimination on boundary pixels, which is useful in cases with low image quality.

II. RELATED WORKS

There have been many works on the segmentation of B-mode echocardiography, which mainly fall into two categories: the traditional methods and the deep learning methods. Most of the solutions based on traditional methods need prior information such as the appearance or shape of the LV [2], [4], [12]–[14], which presents an assumption that the border between myocardium and blood pool is accessible, and therefore possible to achieve good segmentation results based on prior knowledge. As these studies are dependent on the predefined knowledge, so they may fail if data vary from the information stored in the priors. The other methods aim to minimize the energy function by tuning a large number of parameters [3], [15], [16].

Recently, with the development of deep learning in medical image analysis [17], [18] segmentation methods based on deep CNN learned the features with different convolutional kernels and connection methods to obtain accurate and robust results [19]–[24]. Two publicly available datasets in echocardiography CAMUS [25] and Dynamic-Echonet [26] are researched, proved that the deep learning algorithm outperformed in the tasks of segmenting the left ventricle, especially the encoder-decoder based architectures [25]. Several works deal with echocardiographic sequence segmentation by incorporating temporal information such as optical flow [27] and hierarchical convolution aggregated with temporal relevance [28]. However, the temporal information may deteriorate significantly in a low-quality frame because of the high noise. With insights from shape regularization on the prediction, the anatomically constrained neural networks (ACNNs) [9] and shape reconstruction neural network [29] have worked to maintain a realistic shape of the resulting segmentation without post-processing.

VAE [30] approximates posterior distribution via a parameterized variational inference. The distribution is enforced to

be close to a normal distribution as a regularization, which is applied in the cross-modality image segmentation [31], [32]. By regularization, the model can learn a shared domain-invariant latent space with the same distribution. In this work, we applied regularization to narrow the distribution divergence between the latent features of the input and the label, which helps to infer the segmented information.

III. METHODOLOGY

In this paper, the MCAL leverages anatomical and spatial knowledge learned through ground-truth labels on the myocardial segmentation in 2D echocardiography, on the backbone of an encoder-decoder architecture, as shown in Fig. 2. A latent representation encoder maps the input to a latent space by learning the high-level semantical information. For a raw input image, the spatial space contains the segmented anatomical information mixed with the other contexture. So we apply the Kullback-Leibler (KL) divergence loss to learn the distribution difference between the spatial space and the anatomical space of the ground-truth label. On the other hand, FiLM highlights the feature responses in relevant segmentation regions by modulating the spatial space. Further, the skip convolution fusion blocks are applied to discriminate more fine-grained details from the intermediate feature maps through the encoder. Finally, the BDTW focuses on the border neighborhood pixels in each training batch and improves the segmentation accuracy. Fig.3 shows the detailed architecture of each module.

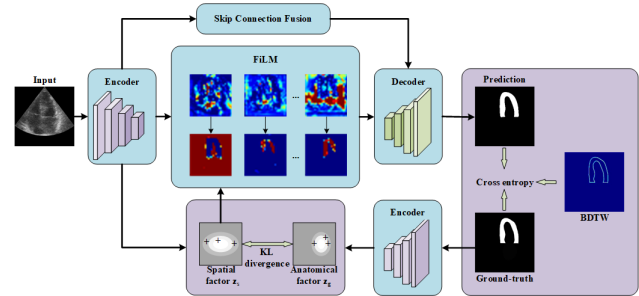


Fig. 2. Block diagram of the proposed model. An input image and its corresponding annotation are encoded to a spatial space s and anatomical space g respectively using an encoder f_s and f_g . Then s and the spatial factor z_s are combined as an input to a decoder f_h to produce a myocardial segmentation prediction. The spatial factor z_s is constrained to learn the distribution, which is close to that of anatomical factor z_g . The parameters of the whole model (the black line and the red line) propagate to achieve an optimal result in the training stage, and the parameters among the black lines are loaded in the test stage.

A. Anatomical Information Derivation with Distribution Divergence Regularization

VAE is a generated model based on samples from a latent variable, of which the posterior distribution is approximated from the input. For an input x , the approximate posterior distribution of its latent variable z can be estimated by an encoder $p(\cdot|\cdot)$. The encoded distributions are set to be an isotropic multivariate Gaussian $N(\mu, \sigma)$ with mean μ and variance σ . Specifically, the encoded feature space produces dimensional mean and diagonal co-variance by dimensional

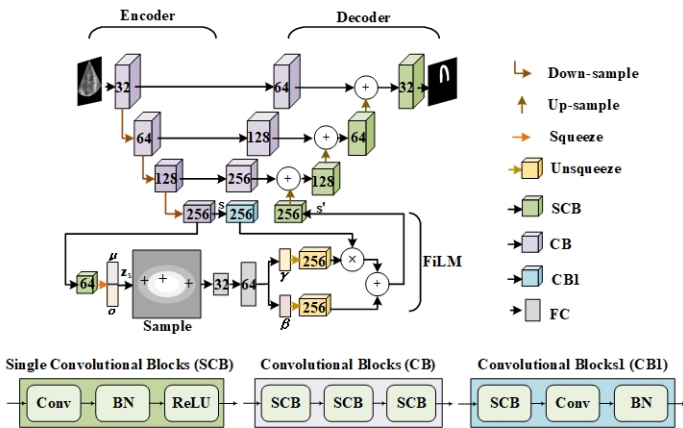


Fig. 3. The architectures of the encoder-decoder that make up the MCAL network. The spatial encoder module is constructed of four convolutional blocks and produces spatial space s for the input image. Then it is modulated by the spatial factor z_s with a FiLM. Finally, the decoder combines the bridge layer from the encoder with a skip convolution fusion block, to produce a segmentation prediction of the myocardial. The number on each box represents the channels of feature maps. Here the “Conv”, “BN”, “ReLU”, and “FC” represent the convolution layer, the batch normalization layer, the rectified linear unit activation layer, and the fully connected layer, respectively.

squeeze, then they are sampled to be an axis-aligned Gaussian distribution to yield the final latent variable. And a decoder $q(\cdot|\cdot)$ converts the samples from z back to the input space.

Motivated by the distribution regularization in VAE, we aim to apply the regularization to narrow distribution divergence between the latent feature space of input and label. Specifically, we first transform the input and label into a latent feature variable by a latent representation encoder, separately. Then we estimated the approximate posterior distribution by a parameterized variational form. The divergence of the distribution from the input and label is regularized to infer segmented-relevance information in segmentation.

The latent representation encoder transforms the input into a spatial representation in our model. The encoded spatial representation is a group of feature maps that contain the spatial information of the input in different channels, so we define spatial feature space s as $f_s(x)$, specifically for an input x . Considering that the encoded latent space of the annotation ground truth y mainly includes the anatomic information, we define the anatomical feature space g as $f_g(y)$. And the latent factor from the encoded feature space is denoted as z_s and z_g for the input and ground-truth, respectively. Their corresponding approximate probability distribution is denoted as $p_\theta(z_s|x, s)$ and $p_\phi(z_g|y, g)$, respectively. The distance between $p_\theta(z_s|x, s)$ and $p_\phi(z_g|y, g)$ is then used as an effective regularization for segmentation directly. The distribution discrepancy convergences gradually during the network training. The distribution difference is directly penalized by the KL divergence:

$$L_{kl} = D_{KL}(p_\theta(z_s|x, s)||p_\phi(z_g|y, g)). \quad (1)$$

B. Feature Enhancement with FiLM

A FiLM constrains the information stored in the spatial space by adjusting the scale and offset of the sampled data

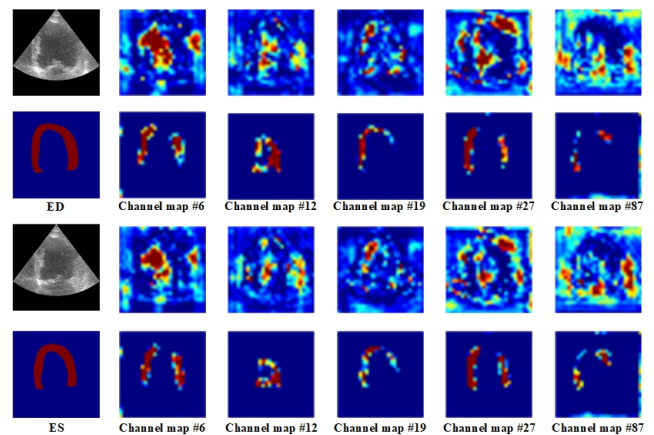


Fig. 4. Visualization of the learned feature maps of some selected channels before and after FiLM. The input images and the corresponding prediction results of our method are in the first column. The five most relevant channel maps of the spatial information before and after modulation are in the first and the third rows, the second and the fourth rows, respectively.

over the spatial factor z_s . The re-scale and offset coefficients are predicted from spatial factor s , which, after a series of convolutions, are conditioned by z_s samples. Specifically, z_s is sampled and then fed into two fully connected layers to obtain the scale γ and offset β , as shown in Fig.3. To modulate each feature map in the spatial space s , γ and β are un-squeezed by the dimensional expanding. Then the spatial space s is passed through a convolution layer, and the modulated output of each channel is formulated by an element-wise multiplication \odot and addition operation as follows:

$$F'_c = \gamma_c \odot F_c + \beta_c. \quad (2)$$

Each feature map is affined to learn from the sampled data, here F_c represents the feature map in each channel c .

To verify the effectiveness of FiLM, the feature maps of some selected channels relevant to the segmentation before and after FiLM are visualized in Fig. 4. The first column represents the input images and predicted segmentation results of the MCAL. The feature maps in the first and third rows demonstrate that the spatial space contains anatomical information drowned in the complicated semantical and spatial information. The second and the fourth rows display the corresponding output of the FiLM. It is obvious that with the modulation of the spatial factor, the network has learned critical information of the segmented structure in some channels. However, we observed that the feature map of channel 12th is weakly associated with the myocardial segmentation, and the structure in the 27th and 87th channels is incomplete. So we use the skip convolution fusion block, which combines the semantic information of the encoder and the relative rich anatomical information in the decoder to solve the problem and fulfill segmentation.

C. Segmentation with Decoder

The architecture of the decoder shown in Fig. 3 is a bottom-up structure with a skip connection with the encoder. Each up-sampling layer adds the feature maps of the bridge pathway

that makes input of the Single Convolutional Block (SCB). The bridge pathway between the encoder and the decoder consists of a convolutional block with three successive convolution layers.

Formally, we formulate the bridge pathway as follows: let F_i^e and F_i^d denote the learned feature maps with the same size before the i^{th} down-sampling layer along the encoder and the i^{th} up-sampling layer along the decoder, respectively. The stack of feature maps represented by F_i^s is computed as

$$F_i^s = H_1(F_i^e) + U(F_{i+1}^d) \quad (3)$$

where function $H_1(\cdot)$ is a convolutional block operation, $U(\cdot)$ denotes an up-sampling layer. Then we obtain the output by $F_i^d = H_2 F_i^s$, where function $H_2(\cdot)$ is a convolution operation followed by batch normalization and ReLU activation.

The decoder recovers the features in each up-sampling layer by fusing corresponding semantical feature information of the encoder. Specifically, the modulated spatial output s' is firstly up-sampled by a bilinear operation, which restores the dimension of the feature maps gradually in each up-sample layer and finally achieves pixel-level prediction. Mainly, the skip convolution introduces the feature maps of the encoder selectively, which contains more semantic information closer to that of the feature maps in the decoder. Then outputs of the skip convolution fusion block are fused with the previous bilinear layer output of the lower skip convolution fusion block. The fusion of the enhanced semantical feature maps with the same size between the encoder and the decoder by the skip convolution fusion block can facilitate optimizer during network training.

D. Boundary Distance Transform Weight Loss

Some works have proposed to solve the label imbalance by multiplying the class weight and tuning the weight value of hard examples iteratively [33]. The boundary is hard to detect in segmentation, especially for ultrasound images with artifacts and speckle noise. Errors located on boundaries affect further index calculation and analysis. A simple and straightforward way to solve this problem is to assign higher weights to the adjacent pixels of the boundary. More weight assigned only to the boundary pixels leads the network to strengthen the boundary information. However, determining the weight value of the boundary region is difficult.

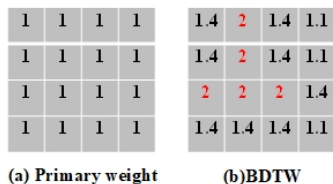


Fig. 5. Illustrations of (a) the primary weight and (b) the boundary distance transform weight (BDTW), the weight values are different according to the distance to the boundary pixels (the red). The BDTW assigns a higher weight to the boundary region, which directly learns the desirable result and therefore helps reduce the prediction errors.

In this paper, we apply the boundary information to calculate the new weight, which displays the distance of each

pixel to the boundary. The weighted distance transform map is decreased exponentially according to the Hausdorff Distance (HD) to the boundary. It forces the network to pay more attention to the boundary, which is formulated as follows:

$$W_{i,j} = \exp(-HD[i,j]) = \exp(-\min_{(k,l) \in q} d([i,j],[k,l])) \quad (4)$$

where d is the distance of each pixel $[i,j]$ of the ground-truth to $[k,l]$, which belongs to a boundary set q . The standard Euclidean distance $d([i,j],[k,l]) = \sqrt{(k-i)^2 + (l-j)^2}$ is used to calculate the distance between pixels. As shown in Fig.5(b), the BDTW assigns the pixels of the boundary region to higher values. Such a mechanism penalizes the hard prediction boundary pixels and therefore helps to reduce the overall prediction errors.

Finally, the boundary distance transform weight (BDTW) loss is obtained by:

$$L_{bdtw} = (\lambda W_{i,j} + 1) \odot L_{ce} \quad (5)$$

where λ is a hyperparameter, L_{ce} is the cross-entropy, and \odot is the Hadamard product. Since the weight of pixels that are far from the boundary is small, hence, to mitigate the vanishing gradient issue, all the weight value is increased by 1. Moreover, the BDTW can be computed in the dataset only once, which does not burden the calculations.

We also adopt the Lovasz-Softmax loss [34] as the loss function to measure the result of the segmentation, as the framework fulfills the segmentation based on the pixel classification problem. The Lovasz-Softmax loss function L_{ls} directly optimizes the mean intersection-over-union loss in the context of semantic image segmentation:

$$L_{ls} = \frac{1}{|C|} \sum_{c \in C} \overline{\Delta_{J_c}}(m(c)) \quad (6)$$

where Δ_{J_c} is the loss surrogate, $m(c)$ corresponds to the vector of pixel identification errors, and C is the class number.

In general, the final loss function is the weighted sum of the segmentation loss and KL loss, and the trainable parameters θ_w are regulated with the L2 paradigm by factor η , which is formulated as:

$$L = \lambda_l L_{ls} + L_{bdtw} + \lambda_v L_{kl} + \eta \|\theta_w\|_2^2 \quad (7)$$

where the λ_l and λ_v are hyper-parameters to allocate the corresponding loss.

IV. MATERIAL AND EXPERIMENTAL RESULTS

A. Dataset Information and Annotations

Our Echocardiography Dataset: We evaluated the proposed MCAL on our dataset, which contains a total of 1472 frames of 11 healthy subjects, and is collected from two hospitals with different devices by Philips and GE, with ethics approval from the Clinical Medical Research Ethics Board. The privacy information of patients is erased at the workstation. The temporal rate is 65–70 Hz among frames. [The pixel resolution of images from the devices are \$0.353 \times 0.353 \text{ mm}^2\$.](#) We research the apical 4-chamber view (A4C) of each examination of those subjects in the experiment. From these images,

we randomly select these subjects to train the model and test. Considering the inter-subjective appearance difference and the intra-subjective frame-relevance, we split the dataset into cross-validation set for training and a dependent test set with a ratio of 9:2 on a subject basis, based on a ratio of 8:2 on an image basis, to illustrate the robustness of our proposed method. Besides, we abandon the first and last frames of the echo cine due to poor image quality. Each subject contains at least one temporally cropped sequence that captures one complete cardiac cycle from ES to ED. An expert annotates the myocardium of each image in the dataset manually according to [1]. The other expert confirms the inconsistency of the annotation and the labels. The annotation masks are considered as ground truth to train our model.

CAMUS Dataset: To comprehensively evaluate the performance of our model, we also use the public CAMUS dataset for verification in the experiment. The CAMUS dataset contains 500 patients with an apical 2-chamber view (A2C) and an A4C view, acquired from a single vendor and center. The pixel resolution of the image is $0.154 \times 0.154 \text{ mm}^2$. Only the annotations of the ES and ED frames are available. Because the annotations of the final 50 are not given in the training data, we adopted 10 folds cross-validation for the evaluation on the CAMUS dataset on the left 450 patients.

B. Data Preprocessing

The raw image is preprocessed to keep the cardiac part only before feeding into the model. We randomly apply rotation augmentation to avoid overfitting, and the rotation angle is between -5 and 5 degrees randomly according to the real echo cine. These images are resized to 224×224 , and the gray value has been normalized to the range $[0, 1]$. For the comparison study and ablation study on the CAMUS with different views and phases, We apply the same preprocessing as [35]. Please refer to [35] for a more detailed pre-processing of the CAMUS dataset. Since we preprocessed the images by resizing the image to 224×224 , the pixel distance in the resized image is scaled down from the original image. We calculated the distance metrics by multiplying the rescaled pixel distance specified on the resized image. More importantly, we set the length and width of the image to be the same by filling zeros before resizing the image. So the aspect ratio of pixel distance is unchanged before and after the preprocessing.

C. Experimental Setup

We use 10-fold cross-validation to train the model. Since our model is based on an encoder-decoder backbone with skip concatenate fusion, we adopted U-net with the same architecture design in [24]. In detail, the number of filters is the same as the U-net in the encoder and decoder. And we set the number of latent vectors to 32 due to the computation efficiency. The normal distribution is applied to initialize the parameters of the network at first. The Adam optimizer applies an initial learning rate of 0.0001 and a weight decay of 0.9 in initialization. The batch size was set to 7 for our dataset. For each fold cross-validation, we trained 100 epochs. The Dice coefficient [36] is used to assess the accuracy of the segment

model. We performed our model on a NVIDIA GeForce RTX 2080Ti GPU in Pytorch.

We explored the impact of hyper-parameters of the loss function on the behavior of MCAL. Because L_{ls} and L_{ce} are fundamental in the segmentation loss, we set the λ_l to be 1 in our experiment. In addition, the λ and λ_v are set the same value to evaluate their performance on segmentation. The performance of MCAL was evaluated under different parameter settings $\lambda \in \{1, 3, 5, 10, 15\}$. The results are shown in Table I. Totally, the best performance is obtained when $\lambda = 15$. We also observed that the Dice dropped obviously when $\lambda \in \{3, 5\}$, which was the worst performance. Table I illustrates that although the performance varies when $\lambda \in \{1, 10, 15\}$, the magnitude of the variation is not very large, so the value of $\lambda = 10$ is adopted for its second-best performance among the three. We trained and test the model with the same parameters settings on our dataset.

D. Comparison With Existing Methods

Comparison on Our Dataset: The comparison study is carried out to evaluate the effectiveness of the network. We compared our method with the UNet, ACNN, and the effectiveness of BDTW on our dataset in this paper. We used the geometrical metric for a comprehensive evaluation of the method: three area error metrics (precision, recall, Dice) and two distance error metrics (absolute surface distance (ASD) and HD). The mean and standard deviation values of each metric were obtained from the cross-validating on the test dataset. We selected the best model on each fold validation set for the test.

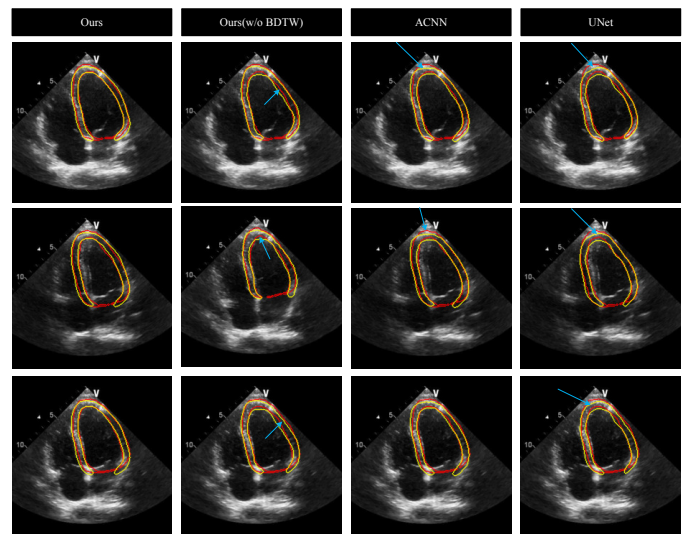


Fig. 6. Qualitative comparison of the results under four different settings on myocardial segmentation. The red and yellow colors denote the ground-truth and predict, respectively. The blue arrow indicates the wrong prediction of the boundary region. The MCAL could improve the prediction accuracy, and the precision improves more for the pixels of the boundary region, revealing that the BDTW is especially effective for the boundary region.

Table II shows the experimental results on our echo dataset. The mean and the standard deviation values are used for each metric to perform the cross-validation procedure. The bold

TABLE I
SEGMENTATION PERFORMANCE UNDER DIFFERENT HYPER-PARAMETER SETTINGS ON CAMUS.
BOLD NUMBERS REPRESENT THE BEST RESULTS OBTAINED.

Weight	A2C			A4C		
	Dice(%)	d_m (mm)	d_H (mm)	Dice(%)	d_m (mm)	d_H (mm)
1	84.15±6.70	0.97±0.44	4.34±3.34	84.69±6.45	0.77±0.31	3.63±3.14
3	83.74±7.16	0.98±0.43	4.70±3.83	83.63±7.05	0.81±0.35	4.23±3.84
5	83.96±6.67	1.00±0.61	4.75±4.03	83.83±6.95	0.80±0.35	4.17±3.85
10	84.27±6.40	0.97±0.48	4.33±3.38	84.78±6.58	0.76±0.33	3.40±2.57
15	84.39±6.60	0.95±0.42	4.36±3.54	84.92±6.50	0.75±0.30	3.46±2.89
	ED			ES		
1	83.93±6.44	0.86±0.42	3.97±3.06	84.84±7.18	0.86±0.36	3.97±3.43
3	83.26±6.67	0.88±0.40	4.36±3.53	84.10±7.49	0.91±0.40	4.55±4.10
5	83.38±6.72	0.89±0.54	4.46±3.82	84.42±6.87	0.90±0.47	4.45±4.08
10	84.13±6.32	0.84±0.44	3.88±2.95	84.93±6.64	0.87±0.40	3.81±3.09
15	84.22±6.41	0.83±0.38	3.95±3.12	85.11±6.67	0.85±0.37	3.83±3.37

TABLE II
MCAL OUTPERFORMS THE OTHER METHODS UNDER DIFFERENT CONFIGURATIONS ON OUR DATASET.
BOLD NUMBERS REPRESENT THE BEST RESULTS OBTAINED.

Method	Precision (%)	Recall (%)	Dice (%)	d_H (mm)	d_m (mm)
UNet	69.20 ±7.11	73.62 ±6.13	71.12 ±5.31	7.21 ±4.32	0.80 ±0.16
ACNN	68.58 ±6.54	76.05 ±8.15	71.76 ±5.33	13.27 ±6.42	0.90 ±0.45
MCAL (w/o BDTW)	69.62 ±6.52	80.44 ±4.54	74.42 ±4.19	7.19 ±3.14	0.70 ±0.09
MCAL	75.43 ±6.24	76.18 ±7.00	76.16 ±4.16	5.85 ±1.74	0.83 ±0.21

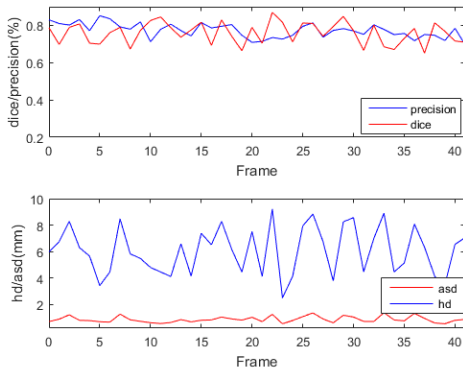


Fig. 7. Precision, Dice, HD, and ASD at different frames of the cardiac cycle of one test subject.

font indicates the best results for each metric. We observed that our framework outperforms other methods on all metrics, achieving the highest mean values of precision (75.43%) and Dice (76.16%), the lowest mean values of HD (5.85 mm), and significantly lower standard deviations of all metrics, especially with the BDTW. This finding demonstrates that a combination of shape and latent anatomical information brings improvement in myocardial segmentation. Fig.6 presents some typical segmentation results, which visually illustrate that the mentioned method obtains a more accurate segmentation, especially the BDTW keeps more fine anatomical information on the prediction.

All echocardiographic sequences from ES to ED are an-

alyzed in Fig.7 to observe the temporal performance of the proposed method. The precision, Dice, HD, and ASD are computed at each frame of a whole cardiac cycle of one test subject to assess the temporal stability. As shown in Fig.7, all the four metrics fluctuate moderately between each frame in the entire cardiac cycle, which means that the proposed method has a limitation on a single image without spatial information. This limitation could be improved by taking into account the relevance of the successive frame.

Comparison on CAMUS: Since the domain gap in our dataset may affect the performance, we conducted another comparison experiment on the public CAMUS dataset to evaluate our method intuitively. For CAMUS dataset, we compared the method with UNet++ [21], SegNet [38], CPFNet [37], HarDNet-MSEG [39] and PLANet in [35], except UNet and ACNN. While the first two are leading methods, the middle two are newly proposed public methods, and the last is a new method proposed on CAMUS. The Dice, ASD, and HD are used in the comparison study.

Geometrical results are analyzed comprehensively by performing the comparison study on the public CAMUS from the view and phase perspective, to assess the influence of the latent representations in the myocardial segmentation between different training views. We carried out the comparison study without any post-processing, such as filling the hole and removing the small area on the segmentation result. Results in Table III showed that the proposed method achieved for most of the metrics compared with other methods. Some methods, such as PLANet and ACNN in our experiments, have been integrated with the label coherence information and shape prior to the learning of anatomical structures. Notably, based on an encoder-decoder design, CPFNet outperformed the other methods and performed close to our method, demonstrating that global/multi-scale information fusion on context information can also achieve better segmentation performance. Our method applied the ground-truth to capture the anatomical information indirectly can achieve higher performance. Furthermore, we performed a statistical comparison of the Dice results using paired t-test with a confidence interval of 0.95. MCAL is compared to CPFNet for statistical significance, and the p values specified to the Dice of A4C/ES/ED are 0.004100/0.000178/0.001100. It can be seen that the proposed method

TABLE III
PERFORMANCE COMPARISON OF MCAL AGAINST EXISTING METHODS ON THE CAMUS DATASET.
BOLD NUMBERS REPRESENT THE BEST RESULTS OBTAINED.

Methods	A2C			A4C		
	Dice(%)	d_m (mm)	d_H (mm)	Dice(%)	d_m (mm)	d_H (mm)
Ours	85.33±5.65	0.92±0.37	3.54±2.45	85.85±5.59	0.73±0.30	3.61±2.98
CPFNet [37]	85.84±6.70	0.84±0.41	4.34±3.28	85.25±6.65	0.75±0.30	3.17±2.08
SegNet [38]	83.37±7.39	0.95±0.48	5.82±4.77	83.45±7.58	0.79±0.39	6.19±5.92
PLANet [35]	83.54±6.38	1.06±0.51	4.36±3.14	85.85±5.68	0.71±0.32	2.97±2.07
HarDNet-MSEG [39]	82.41±6.86	1.14±0.69	4.83±3.41	82.57±7.13	0.89±0.47	4.05±2.99
ACNN [9]	84.31±6.60	0.96±0.57	4.46±3.60	84.23±6.60	0.78±0.37	3.79±3.29
UNet [24]	79.84±8.53	1.28±0.95	6.74±5.10	81.50±7.74	0.91±0.48	5.97±5.05
UNet++ [21]	80.22±8.36	1.27±0.99	7.10±5.46	81.19±7.71	0.94±0.48	6.45±5.53
	ED			ES		
Ours	85.10±5.58	0.83±0.36	3.42±2.28	86.08±5.59	0.81±0.38	4.01±3.43
CPFNet [37]	85.08±6.56	0.78±0.35	3.94±2.86	86.00±6.77	0.85±0.34	3.27±2.25
SegNet [38]	82.97±7.13	0.86±0.46	6.13±5.43	83.86±7.80	0.88±0.43	5.88±5.31
PLANet [35]	83.68±6.16	0.92±0.51	4.14±3.25	85.71±5.96	0.85±0.40	3.18±2.02
HarDNet-MSEG [39]	81.81±6.87	1.03±0.71	4.72±3.54	83.16±7.06	1.00±0.47	4.16±2.85
ACNN [9]	83.81±6.56	0.87±0.57	4.22±3.58	84.73±6.60	0.87±0.38	4.03±3.34
UNet [24]	79.74±8.27	1.14±0.91	6.55±5.10	81.60±7.99	1.05±0.60	6.16±5.07
UNet++ [21]	79.68±8.17	1.17±0.95	6.91±5.50	81.63±7.93	1.05±0.62	7.10±5.46

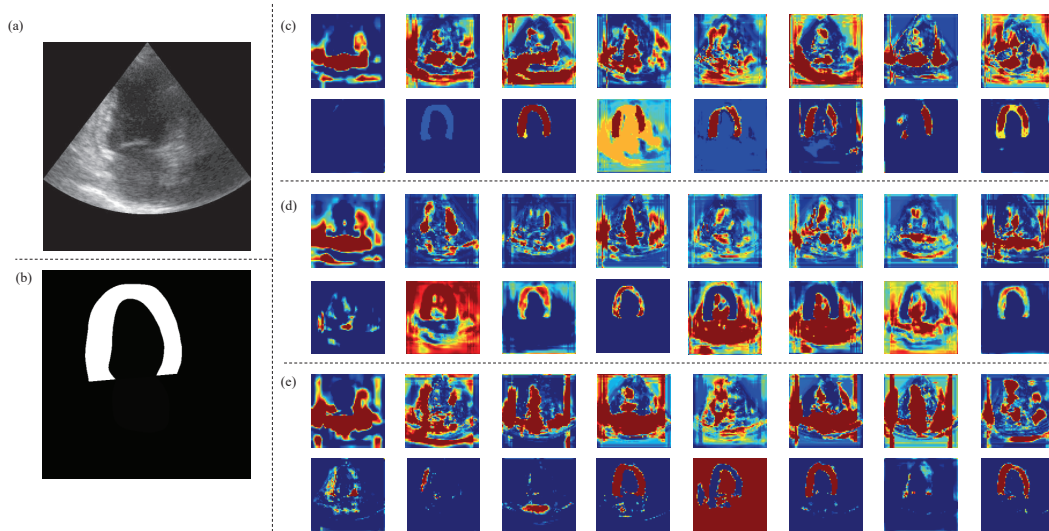


Fig. 8. The FiLM design achieved more effective feature enhancement than other settings. The input image and the corresponding label are in (a) and (b). The performance of the three settings in Table 1 is illustrated in (c), (d), and (e), respectively.

significantly outperforms CPFNet with $p < 0.05$.

E. Ablation Study

We investigated the contributions of each module of our method to the segmentation performance by different configurations in the public CAMUS. We also explored the impact of hyperparameters in the loss function on the behavior of MCAL by grid-searching. In the hyper-parameter setting experiments, the patients with annotations in CAMUS were randomly divided into training (410) and evaluation (40) datasets. In the ablation study experiments, we applied the same training strategy with the comparison study. All the experiments were conducted under the same training and evaluation methods as in [35]. We determined the model with the best performance for each group on the Dice coefficient.

Ablation for Spatial Factor and FiLM: To verify the effectiveness of FiLM design on the relative segmented features,

we replaced the re-scale and offset coefficients in FiLM design by concatenating the dimensional expanding on samples from learned spatial factors z_s with successive convolution layers, which is represented by 'C' in the Spatial factor column. Alternatively, the re-scale and offset coefficients are derived defectively from the spatial space s through successive convolution operations, which is represented by 'o' in the FiLM column. The results are shown in Table IV. We observed a performance drop when the re-scale and offset designs are replaced with successive convolution layers, indicating the effectiveness of the FiLM structure on the feature enhancement. Similarly, deriving the re-scale and offset directly from the spatial space performed worse, demonstrating the necessity of the FiLM structure. The Visualization results in Fig.8 (c) and (e) are evidence that the FiLM design performs better than the concatenation operation. It demonstrates that an affine transformation to each channel of the feature map

TABLE IV
ABLATION RESULTS FOR MCAL WITH DIFFERENT SETTINGS ON CAMUS.
BOLD NUMBERS REPRESENT THE BEST RESULTS OBTAINED.

Spatial factor	FiLM	BDTW	A2C			A4C		
			Dice(%)	d_m (mm)	d_H (mm)	Dice(%)	d_m (mm)	d_H (mm)
✓	✓	✓	85.33±5.65	0.92±0.37	3.54±2.45	85.85±5.59	0.73±0.30	3.61±2.98
<i>C</i>	✓	✓	84.11±6.73	0.97±0.50	4.82±4.18	84.51±6.25	0.76±0.30	3.98±3.74
✓	○	✓	83.83±6.59	0.99±0.50	4.80±3.74	83.30±6.98	0.83±0.32	4.19±3.37
✓	×	×	84.19±6.64	0.97±0.46	4.46±3.43	84.48±6.33	0.77±0.32	3.94±3.82
✓	✓	×	84.31±6.54	0.96±0.45	4.45±3.46	84.39±6.24	0.78±0.35	3.92±3.79
✓	○	×	84.22±6.85	0.97±0.45	4.62±3.74	84.40±6.47	0.77±0.29	3.98±3.63
<i>C</i>	✓	×	83.89±6.56	0.98±0.47	4.73±3.77	83.33±6.84	0.82±0.33	4.17±3.67
✓	×	✓	84.07±6.57	0.97±0.46	4.72±3.92	84.14±6.49	0.78±0.32	3.87±3.39
×	×	✓	84.07±6.64	0.98±0.46	4.47±3.49	83.87±7.34	0.81±0.45	3.93±3.30
×	×	×	83.69±7.05	1.01±0.59	4.45±3.41	84.02±6.93	0.79±0.32	3.69±2.74
			ED			ES		
✓	✓	✓	85.10±5.58	0.83±0.36	3.42±2.28	86.08±5.59	0.81±0.38	4.01±3.43
<i>C</i>	✓	✓	83.84±6.42	0.86±0.48	4.28±3.66	84.78±6.54	0.87±0.36	4.53±4.27
✓	○	✓	82.94±6.72	0.90±0.46	4.53±3.47	84.20±6.80	0.91±0.39	4.46±3.67
✓	×	×	83.70±6.30	0.88±0.46	4.40±3.69	84.80±6.79	0.88±0.39	4.22±3.88
✓	✓	×	83.87±6.14	0.87±0.42	4.17±3.38	84.83±6.43	0.88±0.40	4.17±3.78
✓	○	×	83.75±6.66	0.87±0.43	4.39±3.68	84.88±6.61	0.87±0.35	4.20±3.71
<i>C</i>	✓	×	82.96±6.58	0.90±0.43	4.57±3.70	84.26±6.76	0.90±0.39	4.33±3.74
✓	×	✓	83.51±6.42	0.87±0.43	4.36±3.59	84.70±6.58	0.88±0.38	4.24±3.78
×	×	✓	83.43±6.63	0.88±0.41	4.29±3.33	84.50±7.32	0.91±0.51	4.12±3.47
×	×	×	83.32±6.88	0.90±0.55	4.08±2.97	84.38±7.06	0.91±0.42	4.06±3.26

TABLE V
THE COMPARISON OF EXISTING METHODS FOR LV SEGMENTATION FOR CLINICAL DEPLOYMENT

Methods	Models	Description	Clinical Limitation
Non-deep learning	ACM	Minimizing an energy function under the influence of different forces and constraints	Require user-imposed guidance to achieve high accuracy
	AAM	Describing the image appearance and the shape as a statistical shape-appearance model	Require consistent shape prior over a large database
Deep learning	UNet	Encoder-decoder	Poor model generalization; Limit performance on the LV segmentation
	UNet++	Highly flexible feature fusion	
	HarDNet-MSE	Low memory traffic backbone	
	PLANet	Features enhancement by label coherence learning	Complicate computation; Lack of temporal information
	CPFNet	Feature enhancement by preserving abstract spatial information	Lack of temporal information; Poor model generalization; Lack of model interpretability
	SegNet	Pooling indices are applied in the max-pooling step	
	CNN	Labels are also used as anatomical prior	
	Ours	Labels are also used for feature enhancement	

helps to enhance the segmented-relevant features. However, the learned spatial factor is more effective for guiding feature enhancement (Fig.8 (d)). Obviously, performance has been greatly improved on the condition that the scale and offset parameters are derived from the learned spatial factors.

Ablation for BDTW: Based on these configurations, we investigated the performance of BDTW on segmentation. Results in Table IV displayed that improvement is limit when adding the BDTW to the baseline. However, the Dice dropped slightly when adding the BDTW to the spatial factor, while the distance error metrics improved. Because it is difficult to balance the spatial distribution and the boundary region under two different scales without intermediate operation in the training stage. An improvement was illustrated in Table IV when the FiLM design was added to the configurations. In

conclusion, the joint of FiLM structure on the spatial factors and BDTW has improved the performance of segmentation.

V. DISCUSSION

This work tackles the challenge of myocardium segmentation because of the fuzzy boundary caused by the modality imaging characteristic. This segmentation task can be solved by deep neural networks based on different settings. Since the ground-truth label as class associations on segmentation lack effective feature enhancement of segmented region. We leverage anatomical knowledge learned through ground-truth labels to infer and strengthen the segmented parts. By regulating the distribution discrepancy of two posterior probabilities, which are approximated from the sampled input and labels, we can

modulate the learned anatomical feature maps based on the FiLM. Further, we apply BDTW to assign a higher weight to boundary region pixels in each training batch. Hence, the proposed model can achieve high performance on both segmentation.

This work has limitations. The temporal relevance caused by cardiac movement is neglected in our paper. Researches on 2D echocardiography segmentation have been proven to be effective by adding temporal information [25], [40]. This makes sense as we currently studied anatomical knowledge through labels to infer segmented parts and discriminate boundary pixels. But the segmentation accuracy is limited. Our method still needs to be improved before the clinical deployments. Because we just train and test our model based on data from two different devices, the model generalization has not been fully verified. In the future, it might be expected to embed a temporal correlation of the successive frames or the model generalization to get a better result (TableV).

In general, the proposed method of myocardium segmentation can be applied to other image modalities when we analyze the medical images by segmentation. For instance, it can be used for anatomical segmentation by retraining the network based on new images and labels.

VI. CONCLUSION

Accurate LV segmentation in a 2D echocardiography is significant for cardiovascular disease diagnosis and assessment of cardiac function. However, it is difficult to discriminate between the myocardium and chamber due to the characteristic of echos. A new method named MCAL, which makes use of prior anatomical information and constraint of the predicted map, is proposed to infer the segmented structure and discriminate the boundary pixels. We apply a KL divergence to learn a spatial factor of the raw input that can account for segmentation. Furthermore, the spatial factor modulates the encoded spatial space by FiLM to strengthen the critical segmented structure information in relative channels. A skip convolution fusion block combines the semantical information from the encoder with the relatively rich anatomical information in the decoder, solves the segmentation by a bottom-up structure. In addition, we introduce the BDTW to weight the binary cross-entropy loss, to force the network to focus on the border neighborhood pixels in each training epoch. Finally, we test the proposed method on two different datasets, and experiment results reveal that the proposed method can improve the myocardial segmentation performance.

REFERENCES

- [1] R. M. Lang, L. P. Badano, V. Mor-Avi, J. Afilalo, A. Armstrong, L. Ernande, F. A. Flachskampf, E. Foster, S. A. Goldstein, T. Kuznetsova et al., "Recommendations for cardiac chamber quantification by echocardiography in adults: an update from the american society of echocardiography and the european association of cardiovascular imaging," *European Heart Journal-Cardiovascular Imaging*, vol. 16, no. 3, pp. 233–271, 2015.
- [2] T. F. Cootes, G. J. Edwards, and C. J. Taylor, "Active appearance models," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 23, no. 6, pp. 681–685, 2001.
- [3] M. Kass, A. Witkin, and D. Terzopoulos, "Snakes: Active contour models," *International journal of computer vision*, vol. 1, no. 4, pp. 321–331, 1988.

- [4] M. Rousson and N. Paragios, "Shape priors for level set representations," in *European Conference on Computer Vision*. Springer, 2002, pp. 78–92.
- [5] P. Tang, X. Yan, Y. Nan, S. Xiang, and Q. Liang, "Feature pyramid non-local network with transform modal ensemble learning for breast tumor segmentation in ultrasound images," *IEEE Transactions on Ultrasonics, Ferroelectrics, and Frequency Control*, 2021.
- [6] Y. Fang, H. Huang, W. Yang, X. Xu, W. Jiang, and X. Lai, "Nonlocal convolutional block attention module vnet for gliomas automatic segmentation," *International Journal of Imaging Systems and Technology*.
- [7] H. Ravishankar, R. Venkataramani, S. Thiruvenkadam, P. Sudhakar, and V. Vaidya, "Learning and incorporating shape models for semantic segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2017, pp. 203–211.
- [8] M. Tofighi, T. Guo, J. K. Vanamala, and V. Monga, "Deep networks with shape priors for nucleus detection," in *2018 25th IEEE International Conference on Image Processing (ICIP)*. IEEE, 2018, pp. 719–723.
- [9] O. Oktay, E. Ferrante, K. Kamnitsas, M. Heinrich, W. Bai, J. Caballero, S. A. Cook, A. De Marvao, T. Dawes, D. P. O'Regan et al., "Anatomically constrained neural networks (acnns): application to cardiac image enhancement and segmentation," *IEEE transactions on medical imaging*, vol. 37, no. 2, pp. 384–395, 2017.
- [10] J. Masci, U. Meier, D. Cireşan, and J. Schmidhuber, "Stacked convolutional auto-encoders for hierarchical feature extraction," in *International Conference on Artificial Neural Networks*. Springer, 2011, pp. 52–59.
- [11] E. Perez, F. Strub, H. De Vries, V. Dumoulin, and A. Courville, "Film: Visual reasoning with a general conditioning layer," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.
- [12] J. M. Dias and J. M. Leitao, "Wall position and thickness estimation from sequences of echocardiographic images," *IEEE Transactions on Medical Imaging*, vol. 15, no. 1, pp. 25–38, 1996.
- [13] V. Chalana, D. T. Linker, D. R. Haynor, and Y. Kim, "A multiple active contour model for cardiac boundary detection on echocardiographic sequences," *IEEE Transactions on Medical Imaging*, vol. 15, no. 3, pp. 290–298, 1996.
- [14] Z. Tao and H. D. Tagare, "Tunneling descent for map active contours in ultrasound segmentation," *Medical Image Analysis*, vol. 11, no. 3, pp. 266–281, 2007.
- [15] M. Mignotte, J. Meunier, and J.-C. Tardif, "Endocardial boundary estimation and tracking in echocardiographic images using deformable template and markov random fields," *Pattern Analysis & Applications*, vol. 4, no. 4, pp. 256–271, 2001.
- [16] T. Dietenbeck, M. Alessandrini, D. Barbosa, J. D'hooge, D. Friboulet, and O. Bernard, "Detection of the whole myocardium in 2d-echocardiography for multiple orientations using a geometrically constrained level-set," *Medical image analysis*, vol. 16, no. 2, pp. 386–401, 2012.
- [17] G. Litjens, T. Kooi, B. E. Bejnordi, A. A. A. Setio, F. Ciompi, M. Ghafoorian, J. A. Van Der Laak, B. Van Ginneken, and C. I. Sánchez, "A survey on deep learning in medical image analysis," *Medical image analysis*, vol. 42, pp. 60–88, 2017.
- [18] S. Liu, Y. Wang, X. Yang, B. Lei, L. Liu, S. X. Li, D. Ni, and T. Wang, "Deep learning in medical ultrasound analysis: A review," *Engineering*, vol. 5, no. 2, pp. 261–275, 2019.
- [19] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3431–3440.
- [20] X. Li, H. Chen, X. Qi, Q. Dou, C. Fu, and P. Heng, "H-denseunet: Hybrid densely connected unet for liver and tumor segmentation from ct volumes," *IEEE Transactions on Medical Imaging*, vol. 37, no. 12, pp. 2663–2674, 2018.
- [21] Z. Zhou, M. R. Siddiquee, N. Tajbakhsh, and J. Liang, "Unet++: A nested u-net architecture for medical image segmentation," *arXiv: Computer Vision and Pattern Recognition*, 2018.
- [22] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Deepplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 4, pp. 834–848, 2017.
- [23] D. Fourure, R. Emonet, E. Fromont, D. Muselet, A. Tremeau, and C. Wolf, "Residual conv-deconv grid network for semantic segmentation," *arXiv preprint arXiv:1707.07958*, 2017.
- [24] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.

- 1
2 652 [25] S. Leclerc, E. Smistad, J. Pedrosa, A. Østvik, F. Cervenansky, F. Es-
3 653 pinosa, T. Espeland, E. A. R. Berg, P.-M. Jodoin, T. Grenier et al.,
4 654 “Deep learning for segmentation using an open large-scale dataset in 2d
5 655 echocardiography,” *IEEE transactions on medical imaging*, 2019.
- 6 656 [26] D. Ouyang, B. He, A. Ghorbani, C. Langlotz, P. A. Heidenreich, R. A.
7 657 Harrington, D. Liang, E. A. Ashley, and J. Zou, “Interpretable ai for
8 658 beat-to-beat cardiac function assessment,” *medRxiv*, p. 19012419, 2019.
- 9 659 [27] M. H. Jafari, H. Girgis, Z. Liao, D. Behnami, A. Abdi, H. Vaseli,
10 660 C. Luong, R. Rohling, K. Gin, T. Tsang et al., “A unified framework
11 661 integrating recurrent fully-convolutional networks and optical flow for
12 662 segmentation of the left ventricle in echocardiography data,” in *Deep
13 663 Learning in Medical Image Analysis and Multimodal Learning for
14 664 Clinical Decision Support*. Springer, 2018, pp. 29–37.
- 15 665 [28] M. Li, W. Zhang, G. Yang, C. Wang, H. Zhang, H. Liu, W. Zheng, and
16 666 S. Li, “Recurrent aggregation learning for multi-view echocardiographic
17 667 sequences segmentation,” in *International Conference on Medical Image
18 668 Computing and Computer-Assisted Intervention*. Springer, 2019, pp.
19 669 678–686.
- 20 670 [29] Q. Yue, X. Luo, Q. Ye, L. Xu, and X. Zhuang, “Cardiac segmentation
21 671 from lge mri using deep neural network incorporating shape and spatial
22 672 priors,” *arXiv preprint arXiv:1906.07347*, 2019.
- 23 673 [30] D. P. Kingma and M. Welling, “Auto-encoding variational bayes,” *arXiv
24 674 preprint arXiv:1312.6114*, 2013.
- 25 675 [31] C. Ouyang, K. Kamnitsas, C. Biffi, J. Duan, and D. Rueckert, “Data
26 676 efficient unsupervised domain adaptation for cross-modality image seg-
27 677 mentation,” in *International Conference on Medical Image Computing
28 678 and Computer-Assisted Intervention*. Springer, 2019, pp. 669–677.
- 29 679 [32] F. Wu and X. Zhuang, “Unsupervised domain adaptation with variational
30 680 approximation for cardiac segmentation,” *IEEE Transactions on Medical
31 681 Imaging*, 2021.
- 32 682 [33] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, “Focal loss
33 683 for dense object detection,” in *Proceedings of the IEEE international
34 684 conference on computer vision*, 2017, pp. 2980–2988.
- 35 685 [34] M. Berman, A. Rannen Triki, and M. B. Blaschko, “The lovasz-softmax
36 686 loss: A tractable surrogate for the optimization of the intersection-
37 687 over-union measure in neural networks,” in *Proceedings of the IEEE
38 688 Conference on Computer Vision and Pattern Recognition*, 2018, pp.
39 689 4413–4421.
- 40 690 [35] F. Liu, K. Wang, D. Liu, X. Yang, and J. Tian, “Deep pyramid local
41 691 attention neural network for cardiac structure segmentation in two-
42 692 dimensional echocardiography,” *Medical Image Analysis*, vol. 67, p.
43 693 101873, 2021.
- 44 694 [36] W. R. Crum, O. Camara, and D. Hill, “Generalized overlap measures for
45 695 evaluation and validation in medical image analysis,” *IEEE Transactions
46 696 on Medical Imaging*, vol. 25, no. 11, pp. 1451–1461, 2006.
- 47 697 [37] S. Feng, H. Zhao, F. Shi, X. Cheng, M. Wang, Y. Ma, D. Xiang,
48 698 W. Zhu, and X. Chen, “Cpfnet: Context pyramid fusion network for med-
49 699 ical image segmentation,” in *IEEE TRANSACTIONS ON MEDICAL
50 700 IMAGING*, 2020.
- 51 701 [38] V. Badrinarayanan, A. Kendall, and R. Cipolla, “Segnet: A deep con-
52 702 volutional encoder-decoder architecture for image segmentation,” *IEEE
53 703 transactions on pattern analysis and machine intelligence*, vol. 39, no. 12,
54 704 pp. 2481–2495, 2017.
- 55 705 [39] C.-H. Huang, H.-Y. Wu, and Y.-L. Lin, “Hardnet-mseg: A simple
56 706 encoder-decoder polyp segmentation neural network that achieves over
57 707 0.9 mean dice and 86 fps,” 2021.
- 58 708 [40] H. Wei, H. Cao, Y. Cao, Y. Zhou, W. Xue, D. Ni, and S. Li, “Temporal-
59 709 consistent segmentation of echocardiography with co-learning from
60 710 appearance and shape,” in *International Conference on Medical Image
61 711 Computing and Computer-Assisted Intervention*. Springer, 2020, pp.
62 712 623–632.

Dear Editors,

We would like to thank you and all the reviewers for your very constructive comments and useful suggestions, which have greatly improved our manuscript. The manuscript ID is TUFFC-11452-2021. We have carefully studied each of the comments, conducted point-to-point responses, and revised the manuscript by considering all the suggestions and comments made by the reviewers.

We thank the reviewers for appreciating our work and for their constructive suggestions. We would also like to thank reviewers for deeming that this work “addressed most of the comments from the previous round of reviews”. We are also grateful to all the reviewers for their advice. In this version, we continue to improve this paper while at the same time maintaining the merits mentioned by the reviewers.

Please find one e-copy of the revised version of our submission, and a file containing our responses on how we addressed all the suggestions or comments by the reviewers. For your and the reviewer’s convenience, our responses shown in blue color are prepared based on a point-by-point response to each of the issues raised by the reviewers.

Hope you and the reviewers find the minor revision acceptable. Looking forward to hearing from you about the final decision on our submission.

Sincerely Yours,
Xiaoxiao Cui (on behalf of all the co-authors)

=====
Minor revisions we made include:

- 1) We have corrected all the tables in the manuscript and the response letter. We have added the mention that the bold numbers represent the best results into the table caption and made corresponding changes to the contents. In this way, it avoids confusion or concerns for future readers.
- 2) We have performed a statistical comparison of the Dice results by using paired t-test with a confidence interval of 0.95. MCAL is compared to CPFNet for statistical significance to prove that the proposed method significantly outperforms CPFNet with $p < 0.05$ in lines 450-456 on pages 6-7.
- 3) We have added the details of calculating distance metrics after the resizing operation in lines 348-355 on page 5. In detail, we calculated the distance metrics by multiplying the rescaled pixel distance specified on the resized image.
- 4) We have added information about the resolution of our dataset and CAMUS in lines 310-311 on page 4 and line 333 on page 5, respectively.

Comments from reviewer 1:

Major Comments:

The authors have addressed most of the comments from the previous round of reviews. However, some more explanation and corrections are required before the work can be accepted.

Response: We sincerely thank you for thinking that we “addressed most of the comments from the previous round of reviews”. In this version, according to your following constructive suggestions, we improved this paper to avoid confusion or concerns for future readers.

Some specific points:

Comment#1. Please in the table captions mention that the bold numbers represent the best results obtained. Please also correct the tables to incorporate this. For example, in Table 2 the best results are obtained for weight 15.

Response#1: We sincerely thank you for this detailed suggestion. We have corrected all the tables by adding the mention that the bold numbers represent the best results into the table caption. We have made corresponding changes to the contents of Table I on page 6 and Table III on page 7. In detail:

TABLE I
SEGMENTATION PERFORMANCE UNDER DIFFERENT HYPER-PARAMETER SETTINGS ON CAMUS.
BOLD NUMBERS REPRESENT THE BEST RESULTS OBTAINED.

Weight	A2C			A4C		
	Dice(%)	d_m (mm)	d_H (mm)	Dice(%)	d_m (mm)	d_H (mm)
1	84.15±6.70	0.97±0.44	4.34±3.34	84.69±6.45	0.77±0.31	3.63±3.14
3	83.74±7.16	0.98±0.43	4.70±3.83	83.63±7.05	0.81±0.35	4.23±3.84
5	83.96±6.67	1.00±0.61	4.75±4.03	83.83±6.95	0.80±0.35	4.17±3.85
10	84.27±6.40	0.97±0.48	4.33±3.38	84.78±6.58	0.76±0.33	3.40±2.57
15	84.39±6.60	0.95±0.42	4.36±3.54	84.92±6.50	0.75±0.30	3.46±2.89
	ED			ES		
1	83.93±6.44	0.86±0.42	3.97±3.06	84.84±7.18	0.86±0.36	3.97±3.43
3	83.26±6.67	0.88±0.40	4.36±3.53	84.10±7.49	0.91±0.40	4.55±4.10
5	83.38±6.72	0.89±0.54	4.46±3.82	84.42±6.87	0.90±0.47	4.45±4.08
10	84.13±6.32	0.84±0.44	3.88±2.95	84.93±6.64	0.87±0.40	3.81±3.09
15	84.22±6.41	0.83±0.38	3.95±3.12	85.11±6.67	0.85±0.37	3.83±3.37

TABLE III
PERFORMANCE COMPARISON OF MCAL AGAINST EXISTING METHODS ON THE CAMUS DATASET.
BOLD NUMBERS REPRESENT THE BEST RESULTS OBTAINED.

Methods	A2C			A4C		
	Dice(%)	d_m (mm)	d_H (mm)	Dice(%)	d_m (mm)	d_H (mm)
Ours	85.33±5.65	0.92±0.37	3.54±2.45	85.85±5.59	0.73±0.30	3.61±2.98
CPFNet [37]	85.84±6.70	0.84±0.41	4.34±3.28	85.25±6.65	0.75±0.30	3.17±2.08
SegNet [38]	83.37±7.39	0.95±0.48	5.82±4.77	83.45±7.58	0.79±0.39	6.19±5.92
PLANet [35]	83.54±6.38	1.06±0.51	4.36±3.14	85.85±5.68	0.71±0.32	2.97±2.07
HarDNet-MSEG [39]	82.41±6.86	1.14±0.69	4.83±3.41	82.57±7.13	0.89±0.47	4.05±2.99
ACNN [9]	84.31±6.60	0.96±0.57	4.46±3.60	84.23±6.60	0.78±0.37	3.79±3.29
UNet [24]	79.84±8.53	1.28±0.95	6.74±5.10	81.50±7.74	0.91±0.48	5.97±5.05
UNet++ [21]	80.22±8.36	1.27±0.99	7.10±5.46	81.19±7.71	0.94±0.48	6.45±5.53
	ED			ES		
Ours	85.10±5.58	0.83±0.36	3.42±2.28	86.08±5.59	0.81±0.38	4.01±3.43
CPFNet [37]	85.08±6.56	0.78±0.35	3.94±2.86	86.00±6.77	0.85±0.34	3.27±2.25
SegNet [38]	82.97±7.13	0.86±0.46	6.13±5.43	83.86±7.80	0.88±0.43	5.88±5.31
PLANet [35]	83.68±6.16	0.92±0.51	4.14±3.25	85.71±5.96	0.85±0.40	3.18±2.02
HarDNet-MSEG [39]	81.81±6.87	1.03±0.71	4.72±3.54	83.16±7.06	1.00±0.47	4.16±2.85
ACNN [9]	83.81±6.56	0.87±0.57	4.22±3.58	84.73±6.60	0.87±0.38	4.03±3.34
UNet [24]	79.74±8.27	1.14±0.91	6.55±5.10	81.60±7.99	1.05±0.60	6.16±5.07
UNet++ [21]	79.68±8.17	1.17±0.95	6.91±5.50	81.63±7.93	1.05±0.62	7.10±5.46

Comment#2. The tables in the response letter also should be corrected accordingly since they represent miss information (it gives the impression that the method is achieving the best results compared to SOTA which is not the case (see example Table 3)).

Response#2: We sincerely thank you for this detailed suggestion. We have made corresponding changes to the table captions and contents in Table 1 and Table 3 in the response letters to avoid confusion on the results. In detail:

Table 1. Segmentation performance under different hyper-parameter settings on CAMUS.

Bold numbers represent the best results obtained.

Weight	A2C			A4C		
	<i>Dice</i> (%)	d_m (mm)	d_H (mm)	<i>Dice</i> (%)	d_m (mm)	d_H (mm)
1	84.15±6.70	0.97±0.44	4.34±3.34	84.69±6.45	0.77±0.31	3.63±3.14
3	83.74±7.16	0.98±0.43	4.70±3.83	83.63±7.05	0.81±0.35	4.23±3.84
5	83.96±6.67	1.00±0.61	4.75±4.03	83.83±6.95	0.80±0.35	4.17±3.85
10	84.27±6.40	0.97±0.48	4.33±3.38	84.78±6.58	0.76±0.33	3.40±2.57
15	84.39±6.60	0.95±0.42	4.36±3.54	84.92±6.50	0.75±0.30	3.46±2.89
	ED			ES		
1	83.93±6.44	0.86±0.42	3.97±3.06	84.84±7.18	0.86±0.36	3.97±3.43
3	83.26±6.67	0.88±0.40	4.36±3.53	84.10±7.49	0.91±0.40	4.55±4.10
5	83.38±6.72	0.89±0.54	4.46±3.82	84.42±6.87	0.90±0.47	4.45±4.08
10	84.13±6.32	0.84±0.44	3.88±2.95	84.93±6.64	0.87±0.40	3.81±3.09
15	84.22±6.41	0.83±0.38	3.95±3.12	85.11±6.67	0.85±0.37	3.83±3.37

Table 3. Performance comparison of MCAL against existing methods on the CAMUS dataset.

Bold numbers represent the best results obtained.

Methods	A2C			A4C		
	<i>Dice</i> (%)	d_m (mm)	d_H (mm)	<i>Dice</i> (%)	d_m (mm)	d_H (mm)
Ours	85.33±5.65	0.92±0.37	3.54±2.45	85.85±5.59	0.73±0.30	3.61±2.98
CPFNet	85.84±6.70	0.84±0.41	4.34±3.28	85.25±6.65	0.75±0.30	3.17±2.08
SegNet	83.37±7.39	0.95±0.48	5.82±4.77	83.45±7.58	0.79±0.39	6.19±5.92
PLANet	83.54±6.38	1.06±0.51	4.36±3.14	85.85±5.68	0.71±0.32	2.97±2.07
HarDNet-MSE	82.41±6.86	1.14±0.69	4.83±3.41	82.57±7.13	0.89±0.47	4.05±2.99
ACNN	84.31±6.60	0.96±0.57	4.46±3.60	84.23±6.60	0.78±0.37	3.79±3.29
UNet	79.84±8.53	1.28±0.95	6.74±5.10	81.50±7.74	0.91±0.48	5.97±5.05
UNet++	80.22±8.36	1.27±0.99	7.10±5.46	81.19±7.71	0.94±0.48	6.45±5.53
	ED			ES		
Ours	85.10±5.58	0.83±0.36	3.42±2.28	86.08±5.59	0.81±0.38	4.01±3.43
CPFNet	85.08±6.56	0.78±0.35	3.94±2.86	86.00±6.77	0.85±0.34	3.27±2.25
SegNet	82.97±7.13	0.86±0.46	6.13±5.43	83.86±7.80	0.88±0.43	5.88±5.31
PLANet	83.68±6.16	0.92±0.51	4.14±3.25	85.71±5.96	0.85±0.40	3.18±2.02
HarDNet-MSE	81.81±6.16	1.03±0.71	4.72±3.54	83.16±7.06	1.00±0.47	4.16±2.85
ACNN	83.81±6.56	0.87±0.57	4.22±3.58	84.73±6.60	0.87±0.38	4.03±3.34
UNet	79.74±8.27	1.14±0.91	6.55±5.10	81.60±7.99	1.05±0.60	6.16±5.07
UNet++	79.68±8.13	1.17±0.95	6.91±5.50	81.63±7.93	1.05±0.62	7.10±5.46

1
2
3 **Comment#3.** The authors should provide statistical significance test results (paired t test) for
4 results that are too close to each other. For example, in Table 3 the dice values of A4C/ES/ED
5 for the proposed method vs SOTA (CPFNet) appear to be very similar.
6
7

8 Response#3: We sincerely thank you for this detailed suggestion. We have performed a
9 statistical comparison of the Dice results by using paired t-test with a confidence interval of
10 0.95 in lines 450-456 on pages 6-7. MCAL is compared to CPFNet for statistical significance
11 to prove that the proposed method significantly outperforms CPFNet with $p < 0.05$. In detail:
12
13

14 Furthermore, we performed a statistical comparison of the Dice results using paired t-test with
15 a confidence interval of 0.95. MCAL is compared to CPFNet for statistical significance, and
16 the p values specified to the Dice of A4C/ES/ED are 0.004100/0.000178/0.001100. It can be
17 seen that the proposed method significantly outperforms CPFNet with $p < 0.05$.
18
19

20 **Comment#4.** The authors are resizing the image before processing (resizing to 224×224). This
21 would change the image resolution. How did they calculate the distance metrics reported in the
22 paper after this resizing operation? Details should be included in the main manuscript.
23
24

25 Response#4: We sincerely thank you for this detailed suggestion. We preprocessed the images
26 by resizing the image to 224×224 , and the pixel distance in the resized image is scaled down
27 from the original image. We have added the details of calculating distance metrics after the
28 resizing operation in lines 348-355 on page 5. In detail:
29
30
31

32 Since we preprocessed the images by resizing the image to 224×224 , the pixel distance in the
33 resized image is scaled down from the original image. We calculated the distance metrics by
34 multiplying the rescaled pixel distance specified on the resized image. More importantly, we
35 set the length and width of the image to be the same by filling zeros before resizing the image.
36 So the aspect ratio of pixel distance is unchanged before and after the preprocessing.
37
38
39

40 **Comment#5.** The authors should also include information about the resolution of the
41 ultrasound data (not the size of the image but the actual pixel resolution in mm).
42
43

44 Response#5: We sincerely thank you for this detailed suggestion. The pixel resolution is
45 $0.353 \times 0.353 \text{ mm}^2$ and $0.154 \times 0.154 \text{ mm}^2$ for our dataset and CAMUS, respectively. We have
46 added information about the pixel resolution of our data and CAMUS in lines 310-311 on page
47 4 and line 333 on page 5, respectively.
48
49

50 **Comment#6.** If the image resolution is around 0.2mm the improvements in dm and dH in
51 some cases would correspond to $<$ pixel resolution (for example 0.02mm for A4C and 0.04 for
52 ES). This is a change that can not be detected by the naked eye. I would like the authors to
53 explain why such a small change, which seems very insignificant, is important? Will that have
54 significant long-term effects for the patient.
55
56
57

58 Response#6: We sincerely thank you for this detailed question.
59
60

1
2
3 Mean absolute surface distance (refers to dm) and Hausdorff distance (refers to dH) are two
4 distance error metrics for the image segmentation. dm measures the average minimal distance
5 between two boundaries. dH measures the largest minimal distance between two boundaries.
6 Their values decrease with the increasing resemblance between the segmented results and the
7 ground truth.
8
9

10
11 A small improvement in dm and dH , which is smaller than the pixel resolution, is very important
12 for the myocardium segmentation. Because the myocardium segmentation aims to assess the
13 normality of myocardial movement, each segment of the myocardial is analyzed to diagnose
14 the presence or absence of viable myocardium, which are important considerations for patients
15 with chronic total occlusion (CTO) lesions when choosing a revascularization plan.
16 Especially for the situation where only a few pixels are moving in this segment, the absence of
17 these pixels could lead to the wrong diagnosis. So it has significant long-term effects on the
18 patient.
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60