

Multi-scale, Data-driven and Anatomically Constrained Deep Learning Image Registration for Adult and Fetal Echocardiography

Md. Kamrul Hasan^a, Haobo Zhu^a, Guang Yang^{b,c,d,e}, Choon Hwai Yap^{a,*}

^a*Department of Bioengineering, Imperial College London, London SW7 2AZ, UK*

^b*Bioengineering Department and Imperial-X, Imperial College London, London W12 7SL, UK*

^c*National Heart and Lung Institute, Imperial College London, London SW7 2AZ, UK*

^d*Cardiovascular Research Centre, Royal Brompton Hospital, London SW3 6NP, UK*

^e*School of Biomedical Engineering & Imaging Sciences, King's College London, London WC2R 2LS, UK*

Abstract

Temporal image registration of echocardiography images is a basic building block that can be used for important clinical quantifications such as cardiac motion estimation, myocardial strain measurements, and stroke volume quantifications. Deep learning image registration (DLIR) is a promising approach as it can be consistently accurate and precise, requiring substantially lower computational time and showing promising results in past implementations. We propose that a greater focus on the warped moving image's anatomic plausibility and image quality can support robust DLIR performance. Further, past implementations have focused on adult echocardiography, and there is an absence of DLIR implementations for fetal echocardiography. We propose a framework combining three strategies for DLIR for both fetal and adult echo: (1) an anatomic shape-encoded loss to preserve physiological myocardial and left ventricular anatomical topologies in warped images; (2) a data-driven loss that is trained adversarially to preserve good image texture features in warped images; and (3) a multi-scale training scheme of a data-driven and anatomically constrained algorithm to improve accuracy. Our experiments show that the shape-encoded loss and the data-driven adversarial loss are strongly correlated to good anatomical topology and image textures, respectively. They improve different aspects of registration performance in a non-overlapping way, justifying their combination. We show that these strategies can provide excellent registration results in both adult and fetal echocardiography using the publicly available CAMUS adult echo dataset and our private multi-demographic fetal echo dataset,

despite fundamental distinctions between adult and fetal echo images. Our approach also outperforms traditional non-DL gold standard registration approaches, including Optical Flow and Elastix. Registration improvements could also be translated to more accurate and precise clinical quantification of cardiac ejection fraction, demonstrating a potential for translation. The source code of our data-driven and anatomically constrained DLIR methods will be made publicly available at <https://github.com/kamruleee51/DdC-AC-DLIR>.

Keywords: Echocardiography registration, Deep learning image registration, Adult and fetal echocardiography, Anatomical and data-driven constraints.

Contents

1	Introduction	3
1.1	Background and Motivation	3
1.2	Related Works	6
1.2.1	Non-Echocardiographic DLIRs	6
1.2.2	DLIRs for Echocardiography	7
1.3	Our Contributions	7
2	Methods	9
2.1	Basic Framework	10
2.2	Anatomic Constraints (AC)	11
2.3	Adversarial Data-driven Constraint (DdC)	13
2.4	Multi-Scale Training	14
2.5	DLIR Configuration Summary	14
3	Experimental Setup and Datasets	14
3.1	Elastix and Optical Flow Control Experiments	14

*Corresponding author

Email addresses: k.hasan22@imperial.ac.uk OR kamruleeekuet@gmail.com (Md. Kamrul Hasan), haobo.zhu18@imperial.ac.uk (Haobo Zhu), g.yang@imperial.ac.uk (Guang Yang), c.yap@imperial.ac.uk (Choon Hwai Yap)

3.2	Training Protocol	15
3.3	Evaluation Criterion	15
3.4	Experimental Datasets	17
3.4.1	CAMUS Adult Dataset	18
3.4.2	Private Fetal Dataset	18
4	Experimental Results	19
4.1	Non-Deep Learning Approaches Compared to VanDLIR	19
4.2	Benefits of Anatomic Constraints	22
4.3	Benefits of Adversarial Data-driven Constraints	24
4.4	Effect of Multi-scale Learning	26
4.5	Effect of Data Augmentation	27
4.6	Temporal Registration	28
5	Estimating Ejection Fraction (EF)	31
6	Discussion and Conclusion	32
Appendix A	Structure of Networks	41
Appendix B	Additional Results	43

1. Introduction

1.1. Background and Motivation

Accurate estimations of cardiac motion from clinical echocardiography images are essential and can aid in the evaluation of cardiac function [1] and the detection of dysfunction in conditions such as cardiomyopathy [2], chemotherapy toxicity [3], and infraction [4]. Motion estimation is also a good way to obtain myocardial strains with accuracy on par with speckle tracking [5], and it can be used to estimate stroke volume and ejection fraction (EF) accurately. The extracted motion information is also important for biomechanical simulations of cardiac function to elucidate intracardiac flow and myocardial (MYO) stresses

[6, 7]. However, echo measurements, especially fetal echo, have precision issues due to subjective manual inputs, inter-observer variability, and vendor-specific ultrasound machines and processing algorithms. For instance, clinical measurements of fetal left ventricular (LV) strains can vary more than one-fold between reputable groups [8]. It is thus essential to improve cardiac motion estimation algorithms that could be achieved by temporal echo image registration [9], and machine learning is an ideal way to do so.

Cardiac motion estimation is typically achieved via non-rigid pair-wise temporal echo image registration, coupled with regularizations. Such pair-wise image registration determines the spatial transformation that will enable optimally aligned pixel- or voxel-wise correspondence between two images from different time points, thereby retrieving the deformations of the heart between these two time points. Methodologies previously established for registration included free-form deformation [10], demons [11], and optical flow [12], and regularizations such as cyclic cardiac motion and spatial consistency have previously been proposed [9]. However, these traditional approaches take substantial computational time, especially with high dimensionality and extensive regularizations, which can be improved with deep learning (DL)-based approaches.

Recently, DL-based image registration (DLIR) showed the ability to perform accurate registration that effectively avoids ill-conditioned and highly non-convex cost functions [13, 14]. DLIR can be supervised or unsupervised. For example, Dosovitskiy et al. [15], Rohé et al. [16], Sokooti et al. [17], Yang et al. [18] developed supervised DLIR, which are easy to train and have lower computational costs. However, obtaining the ground-truth motion field for training is often challenging, and non-DL approaches to seeking ground-truths are often imperfect. For this reason, investigators have attempted synthetic images and ground truths, for example, Østvik et al. [19]. For echocardiography images, however, despite advances in generating synthetic images, they may still have differences from natural clinical images. For this reason, unsupervised approaches are of interest, for example, those by Balakrishnan et al. [13], Mansilla et al. [20], Ali and Rittscher [21], Hu et al. [22], Xu and Niethammer [23], which use constraints enforcing a similarity between the moved and fixed images. Many of such works are based on the VoxelMorph algorithm [13], which is successful

due to additional considerations of label shape matching between fixed and warped moving images, demonstrating that a focus on the physical plausibility of the moving image after deformation is important to DLIR.

In echocardiography, a few strategies have been proposed for DLIR, including supervised training based on an optical flow estimator [19], unsupervised patch-based MLP and transformer networks [24], and unsupervised networks utilizing advanced multi-scale correlation [25]. The strategy of improving performance by imposing physiological and physical plausibility on the warped moving image, such as employed by VoxelMorph, has not been fully explored. Further, VoxelMorph imposes an overall dice loss, thus imposing pixel-based local constraints, and has not utilized global latent variables for constraints. We hypothesize that greater emphasis on the quality of the warped image can improve DLIR performance. We have investigated this here.

Further, the majority of current DLIR work has used adult echocardiography datasets. Cardiac motion estimation and myocardial strain measurements are equally important for fetal cardiology, which is important to detect abnormalities and guide management and treatment [26]. Fetal echo images are substantially different from postnatal echo images. The fetal heart is smaller, and transducers are separated further from the heart by maternal tissues, leading to more challenging image quality issues. Consequentially, fetal echo measurements have serious issues with precision [8]. It is unclear if our DLIR approaches work well for fetal images. Thus, we investigate both fetal and adult echo images here.

Here, we propose a new framework for echocardiography DLIR focused on the physical and physiological plausibility of warped moving images. Our framework involves enforcing the anatomic topology of the heart in warped images, as well as warped image quality, using a multi-scale training approach. We show that a focus on the warped image quality is sufficient for robust performance compared to current DLIR reports and can be successful for both adult and fetal echo.

1.2. Related Works

1.2.1. Non-Echocardiographic DLIRs

Several past DLIR studies estimate the spatial transformation parameters under the supervision of the ground-truth deformation vector fields [14]. Conventional registration techniques [27–29] and simulated images with known ground-truth fields [16, 17, 30] are two common methods for obtaining such ground truths. However, the use of conventional registration as ground truth may imply that the DLIR will have accuracy limited by that of the conventional approach, while synthetic images may have differences from natural clinical images.

Unsupervised DLIRs require additional constraints or regularizations for network training [13, 20, 22, 31–35]. Balakrishnan et al. [13] developed VoxelMorph, which uses an unsupervised DLIR regularized by a segmentation-aware auxiliary loss (dice score) to enforce anatomical plausibility in the deformed images. The authors experimentally demonstrated that this improves registration results for brain MR images compared to conventional methods. Similar auxiliary loss concepts were utilized in [22] for T2-weighted transrectal MR and ultrasound scans. VoxelMorph’s idea was further extended by He et al. [34] and Mansilla et al. [20]. He et al. [34] proposed a DLIR method based on the attention-guided fusion of multi-scale deformation fields by having a separate segmentation network for regions of interest to distinguish them from the uninterested areas and successfully tested it on chest X-ray images. Mansilla et al. [20] added a new loss function to VoxelMorph’s segmentation-aware loss to enforce global non-linear representations of image anatomy from the segmentation mask, which was also tested on chest X-ray images. Xu and Niethammer [23] proposed DeepAtlas to perform both weakly supervised registration and semi-supervised segmentation for knee and brain 3D MR images.

Further, Kim et al. [32] and Lian et al. [31] combined Cycle-GAN [36] and VoxelMorph [13] for image registration, using bidirectional deformation to retain cycle consistency and implicit regularization to improve performance, validating the approach with 2D facial expression images, 3D brain MR images, and multi-phase liver CT images. Adversarial constraints were also utilized by Yan et al. [33], Dey et al. [35], and Mahapatra et al. [37] for MR

and transrectal ultrasound scan (TRUS) image fusion, brain MR images, and retinal and cardiac MR images, respectively. Here, an adversarially-trained network was used to classify between the moved and fixed images [38]. Such adversarial training is based on image intensities and their spatial properties, which allows for improved image texture in deformed intensity images. Czolbe et al. [39] developed an unsupervised DLIR network for 2D and 3D MR images of different modalities, using dataset-specific learned spatial invariance features rather than raw intensity images to estimate the loss value. They further showed that deep semantic features can provide DLIR improvements.

1.2.2. DLIRs for Echocardiography

A number of supervised and unsupervised DLIRs have been applied to echocardiography [19, 24, 25, 40, 41]. Østvik et al. [19] proposed a supervised PWC-Net [42]-based framework for motion estimation, which is based on a feature pyramid extractor and an optical flow estimator, and showed that the resulting myocardial strains compared well with truths and current commercial machine quantifications. Wang et al. [24] proposed an unsupervised registration method employing a patch-based MLP and transformer network. The patch-based learning features from moving and fixed images were fitted to the cross-feature block for block matching. MLP and/or Swin Transformer were utilized as cross-feature blocks. Wei et al. [41] proposed co-learning of registration and segmentation using a 3D UNet structure and showed that registration and segmentation benefited from each other. Fan et al. [25] presented an unsupervised multi-scale correlation iterative registration network (SearchMorph) with a correlation layer to improve feature relevance and a correlation pyramid to offer multi-scale relevance information. They additionally developed a deformation field iterator to let the search module register details and large deformations through the model.

1.3. Our Contributions

Existing DLIRs applied to echocardiography typically focus on a match between the warped moving image and the fixed image and do not focus on constraints involving the quality of the warped moving images or the anatomic plausibility of the heart in the warped

image. Anatomic constraints are considered by Voxelmorph, but the constraints are pixel-wise and localized instead of relying on global latent variables. For these reasons, we propose a DLIR framework with a more robust preservation of the anatomic topology of the heart and the quality of the image after registration warping. Further, to date, DLIR has not been applied to fetal echo images. We thus test our framework on both adult and fetal echo. Our specific contributions in this article are:

- We propose a DLIR framework that adds on three strategies to the vanilla DLIR to ensure robust DLIR performance: (1) an anatomic shape-encoded loss to preserve physiological myocardial and cardiac anatomical topologies after warping; (2) a data-driven loss that is trained adversarially to maintain good image texture after warping and to preserve a high degree of perceptual similarity between the fixed and warped intensity images; and (3) a multi-scale training scheme of data-driven and anatomically constrained algorithm to improve accuracy.
- We demonstrate that there are non-overlapping benefits in concurrently using the shape-encoded loss and the data-driven adversarial loss, as the former is correlated and likely responsible for improved anatomical topology, while the latter is correlated to and likely responsible for better image textures in warped moving images.
- We demonstrate that our approach has robust performance in both adult and fetal echocardiography, despite essential differences between the two.
- We demonstrate that our proposed framework offers robust performance in comparison to other DLIR methods and can support accurate clinical measurements such as ejection fraction.

The remainder of the article is prepared as follows: Section 2 details the proposed DLIR incorporating different regularizations, including their implementations. The experimental setup, training protocols, evaluation criteria, and datasets are described in Section 3. Section 4 describes the experimental results with critical observations and a complete ab-

lation study, demonstrating the advantages of different integral components and their non-overlapping benefits. The use of registration for estimating ejection fraction has been shown in Section 5. Finally, Section 6 concludes the article.

2. Methods

Fig. 1 illustrates the proposed DLIR pipeline. The training module consists of the Vanilla DLIR (VanDLIR) (block A in Fig. 1) enhanced with three specific strategies: (1) anatomical constraints (block B in Fig. 1), (2) data-driven constraints (block C in Fig. 1), and (3) multi-scale training. In contrast, the inference module (Fig. 1 right panel) requires only a pair

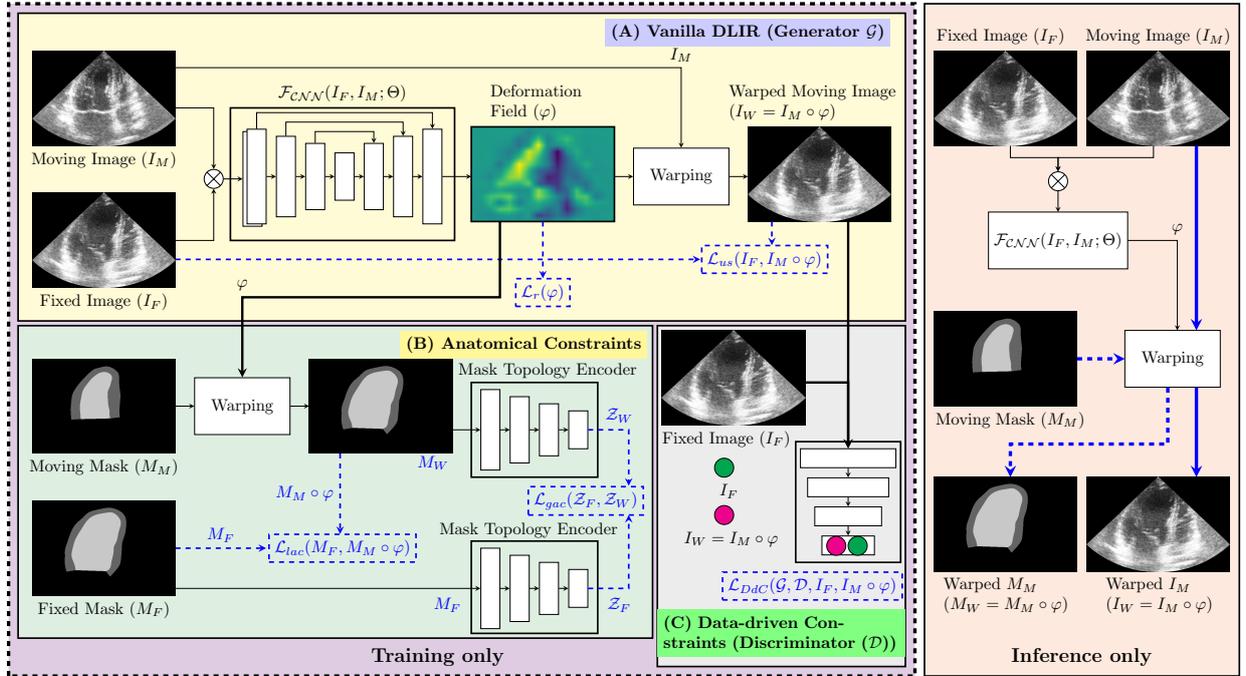


Figure 1: The proposed DLIR model combines different loss functions (\mathcal{L}), such as unsupervised intensity image-based loss (\mathcal{L}_{us}) (in block A), anatomical shape-, size-, and orientation-based loss (\mathcal{L}_{gac}) (in block B), and data-driven adversarial loss (\mathcal{L}_{DdC}) (in block C) through a discriminator. The sum of those losses teaches and constrains the deformation field to minimize the differences between fixed and moving images and/or masks to generate more anatomically plausible and realistic intensity and mask images after the deformation.

of intensity images to predict a deformation field (DDF), which can be used to warp the

intensity and mask moving images to match the fixed images.

2.1. Basic Framework

The basic form of a DLIR is the unsupervised VanDLIR network, which consists of a CNN-based deformation network (\mathcal{F}_{CNN}), containing six layers with (128, 256, 512, 1024, 2048, and 4096) channels in each layer, and a differentiable warping module. \mathcal{F}_{CNN} predicts a deformation field (DDF), $u(x) = \mathcal{F}_{CNN}(I_F, I_M; \Theta)$, where Θ is a \mathcal{F}_{CNN} 's training parameters and $u(x) : \mathcal{R}^N \mapsto \mathcal{R}^N$, where N is an image dimension. The warping transformation (φ) is thus defined as $\varphi(x) = x + u$.

For generating DDF from the inputted intensity image pair (I_F and I_M), we used MONAI's [43] RegUNet¹, as illustrated in Fig. 2, which exploits the benefits of the skip connections from the UNet [44] to recover the spatial information lost in the encoder pooling. \mathcal{F}_{CNN} processes the two intensity images, and generates a 2-channel filter map representing

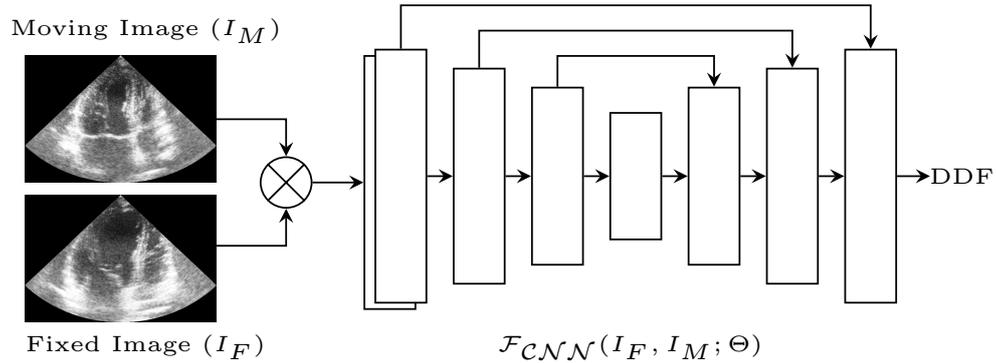


Figure 2: The \mathcal{F}_{CNN} network structure for generating a DDF from the inputted intensity image pair (I_F and I_M), which follows that UNet [44] and DeepReg [45] structures.

the 2D deformation field (3 channels for 3D images), and has a random initialization. For warping the moving image and mask, we used MONAI's [43] Warp² block ($I_W = I_M \circ \varphi$ and $M_W = M_M \circ \varphi$).

¹<https://docs.monai.io/en/stable/networks.html#regunet>

²<https://docs.monai.io/en/stable/networks.html#warp>

The registration between fixed (I_F) and moving (I_M) images is formulated as an optimization problem, as in Eq. 1, which uses pixel displacement fields to represent the spatial transformation,

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmin}} \left\{ \mathcal{L}_{us}(I_F, I_M \circ \varphi) + \lambda_r \times \mathcal{L}_r(\varphi) + \lambda_{lac} \times \mathcal{L}_{lac}(M_F, M_M \circ \varphi) \right. \\ \left. + \lambda_{gac} \times \mathcal{L}_{gac}(M_F, M_M \circ \varphi) + \lambda_{DdC} \times \mathcal{L}_{DdC}(\mathcal{G}, \mathcal{D}, I_F, I_M \circ \varphi) \right\}, \quad (1)$$

where \mathcal{L}_{us} is the mutual information (MI) (dis)similarity measure between I_F and $I_W (= I_M \circ \varphi)$, defined as in Eq. 2.

$$\mathcal{L}_{us}(I_F, I_W = I_M \circ \varphi) = - \sum_{i_f, i_w} p(i_f, i_w) \log \frac{p(i_f, i_w)}{p(i_f)p(i_w)}, \quad (2)$$

where $p(i_f)$ represents the probability that a pixel (or voxel) in image I_F has the intensity of i_f , $p(i_w)$ represents the same probability for image I_W , and $p(i_f, i_w)$ represents the joint probability distribution of the intensities of two images, I_F and I_W . The regularization loss (a bending energy loss [45]), \mathcal{L}_r on $u(x)$ is a smooth regularization on the warping transformations (φ), and $\lambda_r \geq 0$ is the smoothness weight.

\mathcal{L}_{lac} is the local anatomic similarity constrain, \mathcal{L}_{gac} is the global anatomic similarity constraint, and \mathcal{L}_{DdC} is the adversarial data-driven image similarity constrain, while λ_{lac} , λ_{gac} , and λ_{DdC} are their respective weights (≥ 0). These loss terms are explained below.

2.2. Anatomic Constraints (AC)

Several previous studies Balakrishnan et al. [13], Mansilla et al. [20], Hu et al. [22], Xu and Niethammer [23], He et al. [34], Czolbe et al. [39] show that the VanDLIR often fails to align anatomical correspondences between I_F and I_W , and produces anatomically non-plausible shape, size, and orientation of the moved masks. For this reason, we inserted anatomic constraints during training (block B in Fig. 1). First, we modeled a local anatomic similarity loss as proposed in VoxelMorph by Balakrishnan et al. [13],

$$\mathcal{L}_{lac}(M_F, M_M \circ \varphi) = \frac{2}{K} \sum_{c=1}^K \frac{|M_F^c \cap (M_M^c \circ \varphi)|}{|M_F^c| + |M_M^c \circ \varphi|}, \quad (3)$$

where $\mathcal{C}\forall K$ is a number of anatomical regions. Here, $K=2$, as we modeled the LV chamber and myocardium regions. However, this is unlikely to ensure that the warped masks will have sufficient anatomical accuracy. As such, we further applied a variational autoencoder (VAE) [46], $\mathcal{E}\mathcal{N}$, to encode the anatomic topology. This shape encoder is trained to convert the input LV chamber and myocardium mask (I) into latent vectors and to decode them to reconstruct the same mask (I') closely. The lower-dimensional latent representation (\mathcal{Z}) retains significant global information of shape, size, and location. A global anatomic similarity loss, $0 \leq \mathcal{L}_{gac} \leq 1$, is utilized to enforce an anatomic similarity between I_F and I_W ,

$$\mathcal{L}_{gac}(M_F, M_M \circ \varphi) = \|\mathcal{Z}_F - \mathcal{Z}_{W=M \circ \varphi}\|_2^2, \quad (4)$$

To train the VAE network (illustrated in [Appendix A](#) (Table A.4)), we propose using a hybrid loss function, a combination of DSC and structural similarity index (SSIM) [47], rather than the traditional distribution-based cross-entropy loss function to estimate the reconstruction error. The segmentation-aware DSC measures the pixel-level similarity values between I and I' , whereas SSIM analyzes the structural differences between them using a $N \times N$ pixel window. This compound loss aims to produce a better-reconstructed mask. The overall loss function (\mathcal{L}_{VAE}) is defined as,

$$\mathcal{L}_{VAE} = \frac{2}{K} \sum_{\mathcal{C}=1}^K \frac{|I^{\mathcal{C}} \cap I'^{\mathcal{C}}|}{|I^{\mathcal{C}}| + |I'^{\mathcal{C}}|} + \left(\frac{2\mu_I\mu_{I'} + C_1}{\mu_I^2 + \mu_{I'}^2 + C_1} \right) \left(\frac{2\sigma_{II'} + C_2}{\sigma_I^2 + \sigma_{I'}^2 + C_2} \right) - \sum_{i \in N} p(I_i) \log \left(\frac{p(I_i)}{q(I'_i)} \right), \quad (5)$$

where $\mathcal{C}\forall K$ is a number of anatomical organs (which is 2 in our experiments, for the LV chamber and myocardium). C_1 and C_2 are empirically selected scalar parameters; μ_I and $\mu_{I'}$ are the mean values of a neighborhood around I and I' with the individual variances of σ_I^2 and $\sigma_{I'}^2$, respectively; and $\sigma_{II'}$ is their covariance (see details in [47]). The third term in Eq. 5 is the Kullback-Leibler divergence loss [48], which forces the latent vector, \mathcal{Z} , to approximate the normal distribution ($\mathcal{N}(0, I)$). $q(I')$ is the reconstructed distribution of the mask, while $p(I)$ is the actual distribution of the mask. In our DLIR framework, we use the encoder of our pre-trained VAE encoder to extract the latent vector ($\mathcal{Z} \in \mathcal{R}^d$) for a given mask ($\mathcal{Z} = \mathcal{E}\mathcal{N}(I)$), which are then used in Eq. 4 as topological shape features.

2.3. Adversarial Data-driven Constraint (DdC)

In addition, a data-driven constraint was imposed. Similarity constraints imposed by \mathcal{L}_{US} enforce local, pixel-wise similarity and tend to work well with well-aligned images with highly correlated intensities [39], but not with low contrast, noisy images with ambiguous matches, which are typical of echo images. They further lack neighborhood context, such as texture attributes involving neighboring pixels’ correlation and consistency. Enforcing texture similarity is known to give superior network performance and preserve the perceptual similarity of images [33, 49, 50]. To enforce texture similarity, we incorporate a discriminator (\mathcal{D}) in block C of Fig. 1, using a classifier of the warped image output from the VanDLIR block. This data-driven loss function is,

$$\begin{aligned} \mathcal{L}_{DdC} &= \left\langle \mathcal{D}(I_F), \mathcal{D}(I_M \circ \varphi) \right\rangle \\ &= \mathbb{E}_{I_F}[\log \mathcal{D}(I_F)] + \mathbb{E}_{I_W=I_M \circ \varphi}[\log(1 - \mathcal{D}(\mathcal{G}(I_M)))] \end{aligned} \quad (6)$$

The loss function is enforced with adversarial learning [38] by considering the VanDLIR or AC-DLIR as a generator \mathcal{G} , and is trained concurrently with VanDLIR or AC-DLIR. The training procedure for \mathcal{G} maximizes the probability of \mathcal{D} making a mistake in such a way that \mathcal{D} ’s gain is \mathcal{G} ’s loss, and vice versa. As the training goes on, both \mathcal{G} and \mathcal{D} are iteratively updated. The feedback from \mathcal{D} via data-driven loss \mathcal{L}_{DdC} will be used to improve \mathcal{G} so that eventually, \mathcal{G} will be well-trained to generate I_W close to I_F .

The discriminator (\mathcal{D}) is an image classifier that classifies images into fixed or moved echo images. It’s architecture is demonstrated in Appendix A (Table A.5). The discriminator follows the ResNet structure [51] but with modifications. Instead of batch normalization in ResNet, we use instance normalization [52] followed by the 10% dropout layer in each residual block. Instance normalization normalizes across each channel instead of normalizing across input features in a training example. In contrast to batch normalization, the instance normalization layer is also implemented during testing due to the mini-batches independence. Further, we replaced the rectified linear unit (ReLU) activation in ResNet with a leaky ReLU [53] in order to avoid discarding potentially vital information due to negative input values.

2.4. Multi-Scale Training

Coarse-to-fine multi-scale training has shown effectiveness in many previous studies for flow estimation or image registration [17, 54–57] by giving the network more contextual information and sensitivity of features at varying size scales and avoiding local minima [17, 55]. We implemented this for all networks within the three blocks in Fig. 1. During the training phase, trained parameters (Θ) from the lower scale, i.e., resolutions, are used to initialize the training for the following higher scale.

2.5. DLIR Configuration Summary

Our proposed approach is a combination of various strategies and can be summarized in table 1.

Table 1: Summary of different model configurations to validate the role of different integral components in our proposed DLIR.

Model configurations	Anatomical constraints (AC) (W/O shape encoding)	Anatomical constraints (AC) (W/ shape encoding)	Data-driven constraints (DdC)	Multi-scale learning (MS)
VanDLIR	✗	✗	✗	✗
VoxelMorph [13]	✓	✗	✗	✗
AC-DLIR	✓	✓	✗	✗
DdC-DLIR	✗	✗	✓	✗
DdC-AC-DLIR	✓	✓	✓	✗
MS-DdC-AC-DLIR	✓	✓	✓	✓

3. Experimental Setup and Datasets

3.1. Elastix and Optical Flow Control Experiments

Non-deep learning registration approaches are performed for controlled comparisons. The Elastix [58] registration is performed using the Python SimpleITK-Elastix package³ with the settings suggested by Chan et al. [59], which was initially initialized by a rigid Euler transformation followed by an affine and b-spline transformation. Advanced mean squares, transform bending energy penalty, and advanced mattes mutual information are

³<https://simpleelastix.github.io/>

the penalty functions. The dark regions outside the conical field of view of the ultrasound images are masked out to allow image deformation beyond these boundaries. The grid was configured to 16×16 for image warping.

For optical flow (OF) [12] implementation, we have used the iterative Lucas-Kanade solver [60] at each level of the image pyramid. The skimage registration package⁴ has been used for this OF implementation.

3.2. Training Protocol

All the DLIR networks were implemented using Python in the Pytorch⁵ [61] and MONAI⁶ frameworks [43]. The batch size for training DLIRs was eight, for a total of one hundred epochs. All DLIRs were trained to utilize an Adam optimizer without AMSGrad, with a learning rate of 0.0002, betas of 0.9 and 0.999, and a weight decay rate of 0. For multi-scale analysis, all intensity images were resized to 32×32 , 64×64 , 128×128 , 256×256 , and 512×512 using bicubic interpolation, whereas these sizes were obtained from the manual MYO and LV masks using nearest neighbor interpolation. The mask labels for the background, MYO, and LV chambers were 0, 1, and 2, and the intensity images were normalized to the range 0–1. During the validation phase, we set the weighting factors to $\lambda_r = 1$, $\lambda_{lac} = 2$, $\lambda_{gac} = 2$, and $\lambda_{DaC} = 0.001$ by trial-and-error. The machine (running on Ubuntu 20.04 LTS) used for the experiments was equipped with 4 Nvidia[®] GeForce RTX[®] 3090Ti PCIe 3.0 GPUs with 24GB GDDR6X memory with Intel[®] Xeon[®] W-2295 x18 @ 3.00GHz CPU and 128 GB of memory.

3.3. Evaluation Criterion

The quality of intensity image deformation is measured using mean squared error (MSE) and learned perceptual image patch similarity (LPIPS) [49], considering I_F and I_W (see Eq. 7). The mean dice similarity coefficient (DSC) and Hausdorff distance (HD) [62] are

⁴<https://scikit-image.org/docs/stable/api/skimage.registration.html>

⁵<https://pytorch.org/get-started/locally/>

⁶<https://docs.monai.io/en/stable/apps.html>

used to evaluate the quality of mask image deformation, considering M_F and M_W (see Eq. 7).

$$\begin{aligned} MSE &= \frac{1}{N \times M} \sum_{i=1}^N \sum_{j=1}^M (I_{Fij} - I_{Wij})^2, \\ DSC &= \frac{1}{N} \sum_{i=1}^N \frac{2 \times |M_{Fi} \cap M_{Wi}|}{|M_{Fi}| + |M_{Wi}|}, \end{aligned} \quad (7)$$

$$HD = \frac{1}{2 \times N} \sum_{i=1}^N \left(\frac{1}{M_F} \sum_{x \in M_F} \min_{y \in M_W} d(x, y) + \frac{1}{M_W} \sum_{y \in M_W} \min_{x \in M_F} d(x, y) \right),$$

where M , N , and d denote the number of pixels, sample numbers, and shortest distance, respectively. The MSE estimates the local pixel-level similarity, while LPIPS measures the perceptual similarity and computes the similarity between the activations of two image patches for some pre-defined network (VGG [63] in our experiments). A low LPIPS score means that image patches are perceptually similar. The DSC and HD quantify the amount of similarity between M_F and M_W and edge roughness (or irregularity), respectively.

For the clinical evaluation [64–66], we use the correlation between the true and predicted EFs from the end-diastolic (ED) and end-systolic (ES) volumes. The CAMUS dataset (explained below) contains true ED and ES volumes and EFs, and an assumption was made presuming that the number of pixels of the LV area correlates with chamber volume in accordance with Simpson’s rule, as previously proposed [66]. Consequently, using Eq. 8, the predicted ED- and ES-LV volumes (in ml) and corresponding EF are calculated.

$$\begin{aligned} LV_{EDV-pred} &= LV_{EDV-true} \times \frac{LV_{EDV-pred-pxls}}{LV_{EDV-true-pxls}}, \\ LV_{ESV-pred} &= LV_{ESV-true} \times \frac{LV_{ESV-pred-pxls}}{LV_{ESV-true-pxls}}, \\ LV_{EF-pred} &= 1 - \frac{LV_{ESV-pred}}{LV_{EDV-pred}}, \end{aligned} \quad (8)$$

where $LV_{EDV-pred-pxls}$ and $LV_{EDV-true-pxls}$ are the predicted pixel numbers, i.e., LV ED predicted and true areas, respectively. Similarly, $LV_{ESV-pred-pxls}$ and $LV_{ESV-true-pxls}$ are for the ESs. $LV_{EDV-true}$ and $LV_{ESV-true}$ are the true ED and ES volumes that are provided in the CAMUS dataset [64].

In addition to DSC and HD, we propose a new metric for assessing the anatomic plausibility of the MYO, which measures the MYO’s thickness uniformity (TU) by estimating

the thickness variance, as illustrated in Fig. 3. To calculate TU, we consider N number of

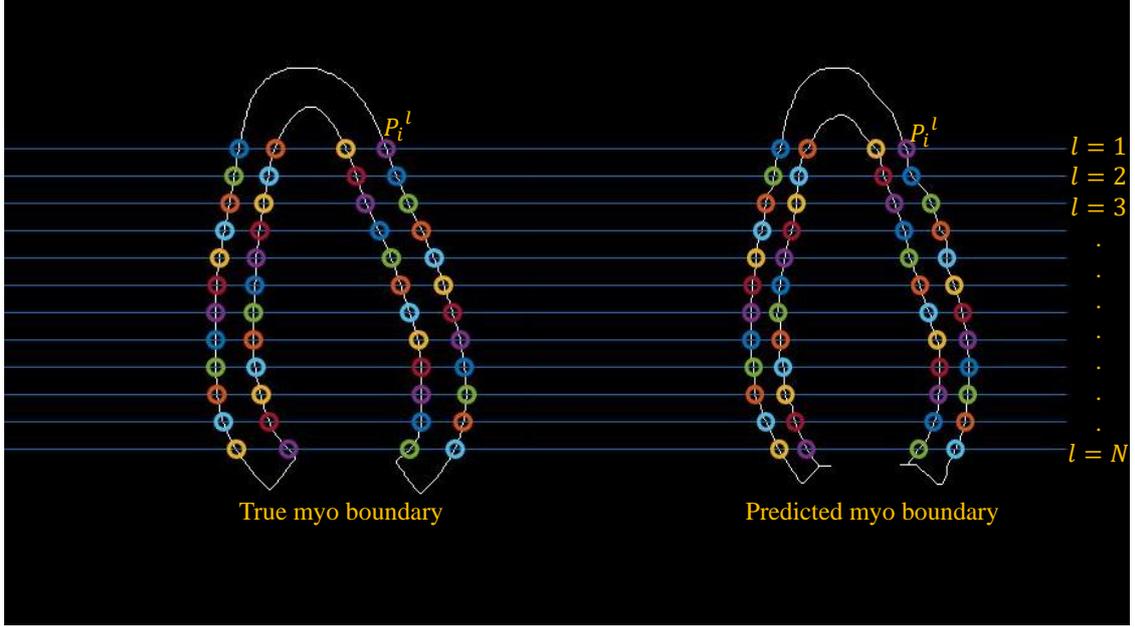


Figure 3: The proposed variance estimation of the MYO thickness to measure its thickness uniformity (TU) from the true and predicted MYO boundaries.

horizontal lines (l) and determine the intersection points (P_i^l) between MYO's inner and outer boundaries (see Fig. 3). Each true and predicted MYO boundary returns four ($C = 4$) spatial coordinates ($i \forall C$) for each line l . Then, the estimated true and predicted TUs are estimated from the N-dimensional distance vector d_N using Eq. 9.

$$d_N = \left[\frac{1}{2} (|d_R^l(P_{i=1}^l, P_{i=2}^l)| + |d_L^l(P_{i=3}^l, P_{i=4}^l)|) \right]_{N \times 1}, \quad (9)$$

$$TU = \mathbb{E}[d_N^2] - \mathbb{E}[d_N]^2$$

where d_R^l and d_L^l are the l^{th} distances between a pair of spatial coordinates (MYO's right- and left-side pairs, i.e., (P_1^l, P_2^l) and (P_3^l, P_4^l)). The smaller value of TU denotes the less thickness variance in the MYO boundary.

3.4. Experimental Datasets

Comprehensive ablation studies are conducted using two different datasets: the publicly available adult echo dataset (CAMUS [64]) and our private fetal echo dataset.

3.4.1. CAMUS Adult Dataset

This adult dataset includes 500 patients with an apical two-chamber view (A2C) and an apical four-chamber view (A4C). It was acquired from GE Vivid E95 ultrasound scanners (GE Vingmed Ultrasound, Horten, Norway), with a GE M5S probe, with a pixel resolution of $0.154m \times 0.154m$. Further descriptions of the CAMUS dataset can be found in [64].

3.4.2. Private Fetal Dataset

This fetal dataset is a collection of 4D (3D over time) fetal echocardiography images, consisting of 105 2D echo videos obtained at various planes from the 3D echo videos of 15 patients. Images were acquired with the GE Volusion 730 ultrasound machine with the RAB 4 – 8L transducer and have an in-plane image resolution of $0.95\mu m \times 0.90\mu m$ (see details in [9]). The dataset consists of fifteen healthy cases and nine disease cases (with fetal aortic stenosis), with the training dataset involving thirteen healthy and seven disease cases (5813 2D images and masks) and the validation set containing two of them each (795 2D images and masks). The dataset was labeled manually at the ES and ED time points of each 2D video, followed by a temporal registration to generate the masks of the other time points using a validated cardiac motion estimation algorithm [9, 59], which imposes a spatial b-spline of temporal Fourier registration over pair-wise Elastix image registration. Fig. 4 shows two examples of our private fetal echo dataset.

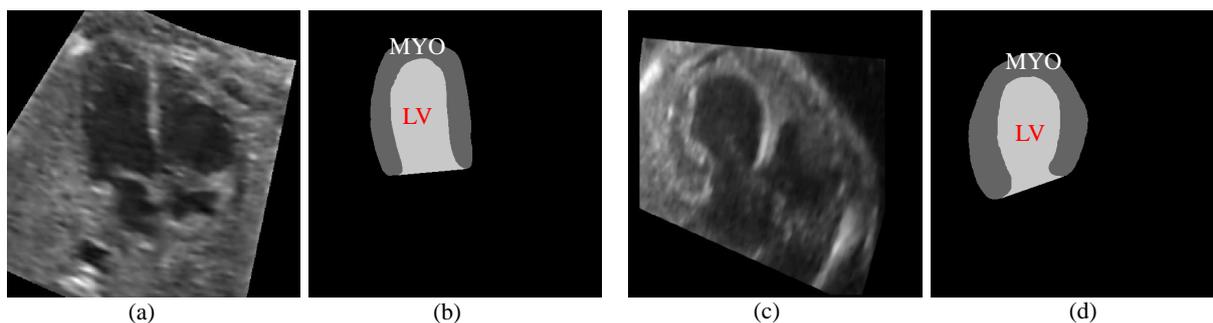


Figure 4: Examples of manually segmented LV and MYO for our private fetal echo dataset. (a) is for healthy patients with a corresponding mask in (b), and (c) is for patients with the disease with a corresponding mask in (d).

4. Experimental Results

In Sections 4.1, 4.2, 4.3, and 4.4, we undertake extensive ablation experiments, where we systematically remove specific components to reveal the effects of various components in our proposed DLIR method and to show that our proposed approach gives the best results. The performance of our ablation experiments can be summarized in Table 2, Fig. 5, and Fig. 6. These experiments utilized the ED time frame as the fixed (I_F and M_F) image and the ES time frame as the moving (I_M and M_M) image from the adult CAMUS dataset for both A2C and A4C views. In Section 4.5, we demonstrate that image augmentation further improves this best-performing DLIR model. Finally, in Section 4.6, we proceed to test the registration algorithm for a wider range of image pairs between the first time frame and all other time frames and to test the algorithm for both adult echo images and fetal echo images.

4.1. Non-Deep Learning Approaches Compared to VanDLIR

Without any registration (N/A), the DSCs between ED and ES time frames are 0.8024 and 0.7923 for CAMUS A2C and A4C views, respectively, while their HD are 5.73 mm and 5.96 mm (Table 2). This is thus the no-registration control. All registration algorithms we tested have performances above this, with higher DSC and lower HD between the warped moving image and the fixed image. The average TU of the fixed images is 0.76 mm and 0.72 mm for A2C and A4C views, and for some algorithms, the warped moving image has TU close to this baseline. In terms of MSE, the baseline no registration control is 0.0026 and 0.0031 for A2C and A4C views, but not every algorithm improves on this.

Between the two gold-standard non-deep learning algorithms, Elastix and OF, which underwent tuning, we observe that OF outperformed Elastix in terms of DSC, HD, TU , and LPIPS, but not MSE. Better DSC, HD, and TU indicate that OF has a better overall warped myocardium shape, which is also observable in Fig. 5. A better LPIPS demonstrates that there is a better perceived visual similarity between the warped moving image and the fixed image. However, OF’s MSE is very high, indicating insufficient similarity in local pixel-level intensity between warped moving and fixed images. On the other hand, Elastix performs

Table 2: Registration results of the CAMUS adult echo dataset from different techniques, employing ES (moving, I_M and M_M) and ED (fixed, I_F and M_F) images of cardiac A2C and A4C views and demonstrating the non-overlapping benefits of data-driven and anatomical constraints. All the metrics are estimated using fixed images (and masks) and warped moving images (and masks). The DSC and HD are the averages of backgrounds, MYO, and LV metrics, whereas the class-wise metrics of MYO and LV are shown in Table B.6 of Appendix B. Bold fonts denote the best-performing metrics for the A2C echo view, while the best-performing metrics for the A4C view are underlined.

Methods		DSC (\uparrow)	HD (mm) (\downarrow)	TU ‡ (mm)	MSE (\downarrow)	LPIPS (\downarrow)	Time (\downarrow)
N/A	A2C	0.8024 \pm 0.0485	5.73 \pm 1.62	–	0.0026 \pm 0.0010	–	–
	A4C	0.7923 \pm 0.0492	5.96 \pm 1.64	–	0.0031 \pm 0.0014	–	
Elastix [58]	A2C	0.8825 \pm 0.0477	5.23 \pm 2.63	3.83 \pm 6.03	0.0033 \pm 0.0019	0.2264 \pm 0.0495	4270 ms
	A4C	0.8892 \pm 0.0388	4.75 \pm 2.69	2.92 \pm 3.21	0.0049 \pm 0.0030	0.2388 \pm 0.0501	
OF [60]	A2C	0.8980 \pm 0.0323	4.48 \pm 1.61	3.14 \pm 3.16	0.0590 \pm 0.0145	0.1098 \pm 0.0151	937 ms
	A4C	0.8995 \pm 0.0358	4.48 \pm 3.0	2.75 \pm 3.91	0.0588 \pm 0.0139	0.1202 \pm 0.0314	
VanDLIR	A2C	0.8511 \pm 0.0456	5.44 \pm 1.70	2.76 \pm 2.11	0.0013 \pm 0.0005	0.1049 \pm 0.0113	85 ms
	A4C	0.8484 \pm 0.0421	5.51 \pm 1.71	2.17 \pm 1.92	<u>0.0017 \pm 0.0007</u>	0.1083 \pm 0.0132	
VoxelMorph [13]	A2C	0.8698 \pm 0.0408	4.66 \pm 1.36	1.80 \pm 1.90	0.0026 \pm 0.0009	0.1120 \pm 0.0126	82 ms
	A4C	0.8520 \pm 0.0448	4.71 \pm 1.54	1.54 \pm 1.58	0.0031 \pm 0.0014	0.1115 \pm 0.0135	
AC-DLIR	A2C	0.8854 \pm 0.0360	4.13 \pm 1.30	1.85 \pm 2.06	0.0027 \pm 0.0010	0.1204 \pm 0.0141	91 ms
	A4C	0.8867 \pm 0.0353	3.82 \pm 1.33	1.64 \pm 1.56	0.0033 \pm 0.0014	0.1239 \pm 0.0151	
DdC-DLIR	A2C	0.8726 \pm 0.0447	4.78 \pm 1.97	1.90 \pm 1.41	0.0015 \pm 0.0006	0.0951 \pm 0.0126	92 ms
	A4C	0.8669 \pm 0.0407	5.35 \pm 1.95	1.34 \pm 1.28	0.0019 \pm 0.0008	<u>0.1005 \pm 0.0138</u>	
DdC-AC-DLIR	A2C	0.9107 \pm 0.0282	3.43 \pm 1.32	1.73 \pm 2.37	0.0016 \pm 0.0005	0.1024 \pm 0.0126	92 ms
	A4C	<u>0.9136 \pm 0.0238</u>	3.99 \pm 1.08	1.16 \pm 1.13	0.0020 \pm 0.0008	0.1046 \pm 0.0139	
MS-DdC-AC-DLIR	A2C	0.9105 \pm 0.0293	3.50 \pm 1.32	1.36 \pm 1.49	0.0016 \pm 0.0006	0.1001 \pm 0.0127	93 ms
	A4C	0.9126 \pm 0.0216	<u>3.34 \pm 1.20</u>	<u>0.98 \pm 0.98</u>	0.0019 \pm 0.0008	0.1054 \pm 0.0139	

‡ An accurate TU should be as close as the fixed MYO’s TU, which is 0.76 mm for adults in A2C and 0.72 mm for adults in A4C views.

very poorly in LPIPS, showing significant perceived differences between warped moving and fixed images.

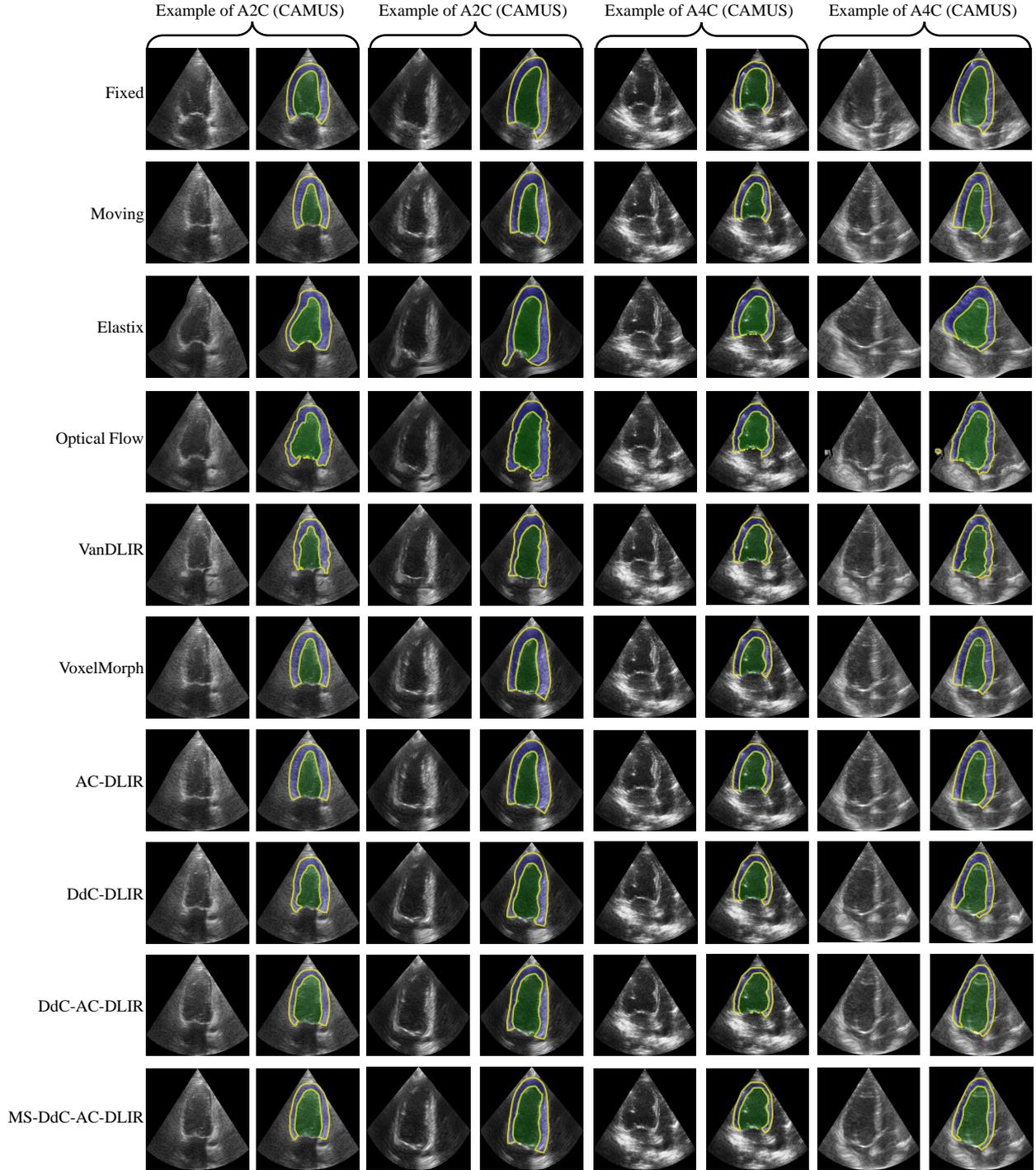


Figure 5: Example of qualitative results of the deformed intensity ($I_W = I_M \circ \varphi$) and mask ($M_W = M_M \circ \varphi$) images generated by different registration models for CAMUS’s A2C (first four columns) and A4C (last four columns) views. The fixed and moving examples are also displayed in the first two rows for comparison.

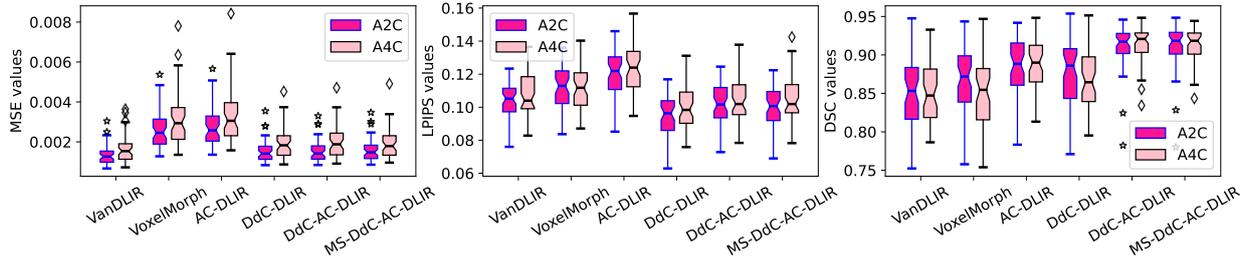


Figure 6: Demonstration of the non-overlapping benefits of anatomical and data-driven constraints, as well as the effects of their combination. We use the CAMUS A2C and A4C views to make this figure.

Comparing VanDLIR to these gold-standard non-deep learning methods, VanDLIR had poorer results in terms of DSC and HD, demonstrating that its warping produces poorer myocardium and LV chamber shapes with lower anatomic plausibility and uneven topology. The greater TU of VanDLIR than Elastix and OF indicates a more uniform MYO thickness. Its visual perception quality, LPIPS, is also better than Elastix and OF. However, VanDLIR has a very low MSE, indicating an excellent local-pixel intensity match between moving and fixed images. Since the VanDLIR is designed to seek only local intensity matching and does not enforce matching of neighboring pixels or global attributes, it is unsurprising that MSE is over-emphasized and very low at the expense of other performance indicators.

For all three methods, VanDLIR, Elastix, and OF, the shapes and edges of the myocardium and LV are very irregular and coarse (see Fig. 5). Their TU is several times higher than the fixed image TU of 0.76 mm. It is thus important to add anatomic regularization to improve anatomic plausibility and topology.

4.2. Benefits of Anatomic Constraints

Our results show that the addition of shape constraints indeed improved registration performance (see Table 2). First, we investigate adding an auxiliary local loss function \mathcal{L}_{lac} (Eq. 3) to VanDLIR, as proposed in VoxelMorph [13], and find that this improves DSC, HD, and TU by 2.2%, 14.3%, and 34.8% for the A2C view and 0.4%, 14.7%, and 29.0% for the A4C view, respectively, all of which are statistically significant ($p \ll 0.001$). In Fig. 5, the improved thickness uniformity and anatomy overlapping between fixed and warped masks

can be seen in comparison with VanDLIR. However, LPIPS and MSE worsened, and visual similarity is poorer than VanDLIR (see Fig. 5 and Fig. 6). This is again unsurprising, as \mathcal{L}_{lac} in VoxelMorph seeks only the local pixel-level similarity between the warped and fixed masks. It does not enforce matching neighboring pixels’ intensity and/or global semantic attributes.

For further anatomical plausibility improvement, we pursued a second shape constraint on the latent space of the VAE for the masks (the AC-DLIR approach). This gives the network global anatomical knowledge and topology, enabling more stringent matching of fixed and warped moving masks. The pre-training of the VAE was successful, producing high DSC between input and output myocardial and LV masks in both the CAMUS and fetal echo datasets, and high visual similarity between input and output masks (see Appendix B for Fig. B.13 and Fig. B.14). This demonstrates that the latent space design is sufficient to capture the global attributes of both the myocardium and LV chamber masks.

Results show that the AC-DLIR improved on VanDLIR more than VoxelMorph (see Fig. 6), where DSC, HD, and TU have been improved by substantial margins, which are statistically significant ($p \ll 0.001$). The DSC and HD of AC-DLIR also outperformed VoxelMorph by 1.8% and 11.4% for the A2C view and 4.1% and 18.9% for the A4C view, respectively, while TU and MSE were very close to VoxelMorph. A noteworthy 33.0% for A2C and 24.4% for A4C improvements in TU from VanDLIR indicates much better warped anatomy, which can be confirmed by visual inspection of Fig. 5.

Visual inspection in Fig. 7 also revealed that the low \mathcal{L}_{gac} loss function correlated well with qualitatively physiologic-looking masks. At the same time, correlation analysis shows that \mathcal{L}_{gac} loss function has a strong correlation coefficient of 0.855 with HD, which is stronger than the correlation of \mathcal{L}_{lac} or \mathcal{L}_{DdC} with HD (see Fig. 8 (a) and Fig. 8 (b)). This indicates that the global shape constraint in AC-DLIR is important for good warped anatomical shapes and DSC outcomes (see Fig. 6). However, Fig. 8 (d) and Fig. 8 (e) also show that there are poor correlations between LPIPS and \mathcal{L}_{gac} or \mathcal{L}_{lac} . This effect has been reflected in Table 2 as the addition of those two anatomical losses worsened the MSE and LPIPS compared to the VanDLIR (see Fig. 5 and Fig. 6). Thus, although the shape constraints can

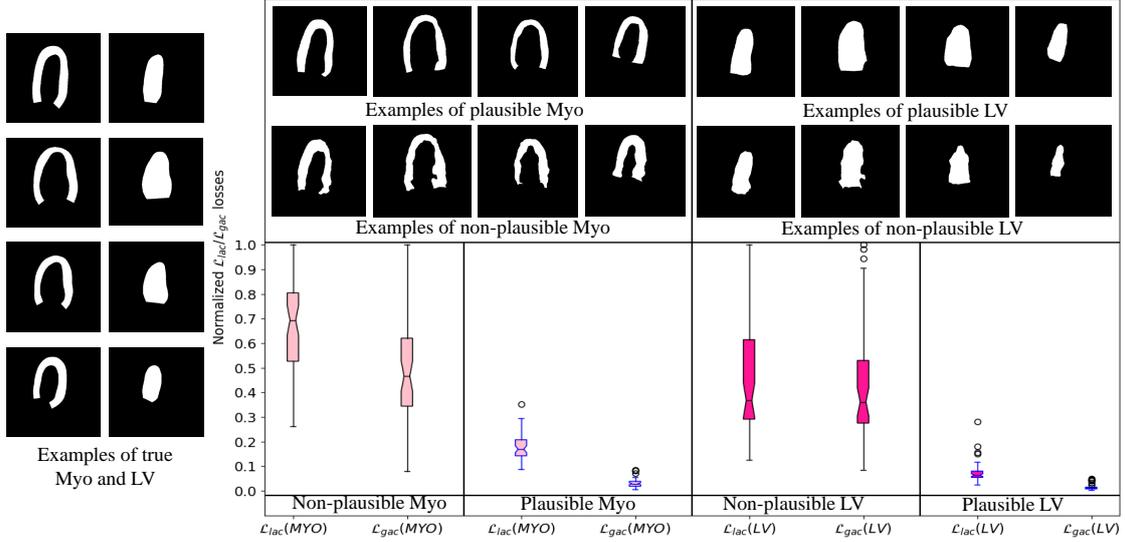


Figure 7: Normalized \mathcal{L}_{gac} and \mathcal{L}_{lac} losses for showing the effect of incorporating L_2 loss (\mathcal{L}_{gac}) in our proposed DLIR for the anatomical plausibility and realistic deformation. Non-plausible and plausible examples are selected from the VanDLIR and our proposed DLIR, respectively.

improve warped image anatomic shapes, they do not play a role in enhancing the perceptual match between the fixed and warped images.

4.3. Benefits of Adversarial Data-driven Constraints

To improve the warped image’s visual quality, we investigated the addition of the adversarial image classifier network (DdC-DLIR), using a data-driven loss function that maximizes the global semantic correspondence between the fixed and warped moving intensity images. Results in Table 2 show that adding DdC to VanDLIR significantly ($p \ll 0.001$) improves all metrics (DSC, HD, TU , and LPIPS) except for MSE, which remains similar to VanDLIR. Visually, the image has better similarity than the fixed image, but the warped label does not have a good anatomical shape in many cases (see Fig. 5).

Compared to AC-DLIR, DdC-DLIR had an almost similar DSC. However, it had lower LPIPS, indicating better visual warped image quality. This suggests that the anatomic constraints imposed in AC-DLIR enforced a better warped cardiac shape, while the semantic correspondence constraint in DdC-DLIR enforced a better visual perception for the image, and the two approaches provide non-overlapping benefits. Correlation analysis in Fig. 8

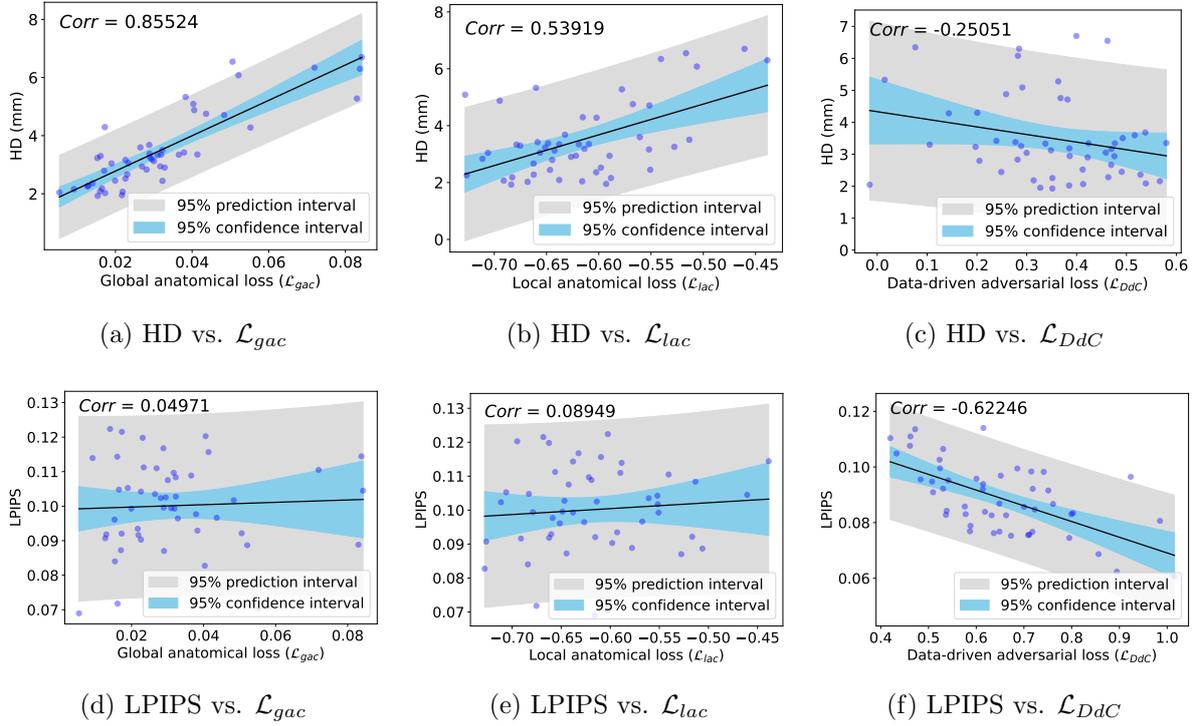


Figure 8: Degree of correlation between two anatomical losses (\mathcal{L}_{lac} and \mathcal{L}_{gac}) and data-driven loss (\mathcal{L}_{DdC}) with anatomical measures (HD) and visual similarity measures (LPIPS). The negative correlations in (c) and (f) indicate that when \mathcal{L}_{DdC} is high, the discriminator becomes more confused about classifying I_F and I_W , i.e., they are similar, yielding better LPIPS and HD.

further confirmed that the \mathcal{L}_{gac} and \mathcal{L}_{lac} loss functions in AC-DLIR correlated better with HD, while the \mathcal{L}_{DdC} loss function in DdC-DLIR correlated better with LPIPS (see Fig. 8 (f)). We thus investigated the combination of both strategies, DdC-AC-DLIR, to reap the benefits of both approaches (see Fig. 6). Results in Table 2 show that the combined approach caused improvements to DSC and TU , above AC-DLIR and DdC-DLIR, while other metrics are similar to AC-DLIR and DdC-DLIR. Average DSC has now improved to >0.91 for both A2C and A4C, while TU substantially reduces to less than 2.0 mm. Fig. 9 demonstrates that when AC-DLIR and DdC-DLIR are combined in DdC-AC-DLIR, 60.0% of testing patients have a DSC of 0.91 to 0.96, but only 30.0% of patients have DSCs in that range when AC-DLIR is used alone and only 20.0% of patients have DSCs in that range when DdC-DLIR is used alone. It has also been observed that the DdC-AC-DLIR produces DSCs of <0.86

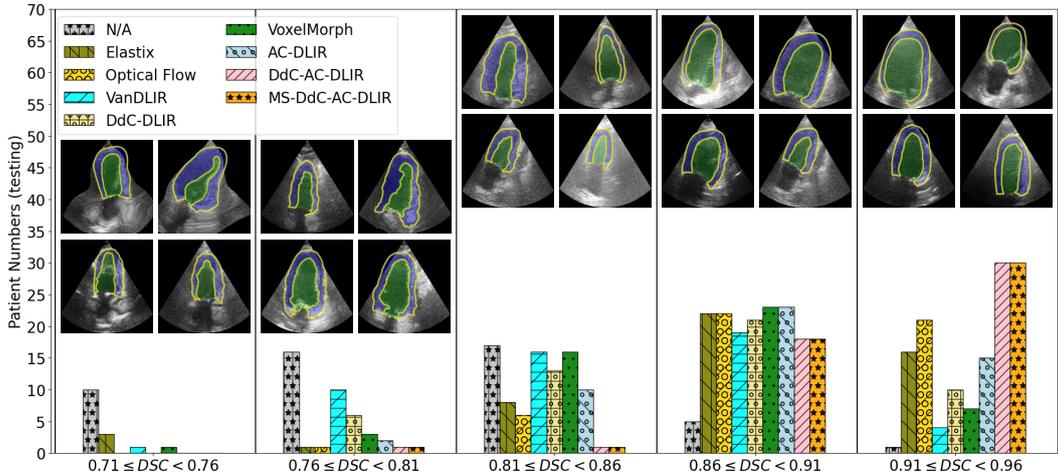


Figure 9: Five different groups of obtained DSCs, demonstrating the number of testing patients in each DSC group. This figure is for A2C of CAMUS, and a similar figure for A4C of CAMUS is in Fig. B.15 of Appendix B.

for an extremely small number of image samples, unlike AC-DLIR and DdC-DLIR, again demonstrating the non-overlapping benefits of the two strategies.

4.4. Effect of Multi-scale Learning

We further investigated the use of a multi-scale learning approach to improve performance. We first tested the performance of DdC-AC-DLIR at various image scales, from 32×32 to 512×512 . We find that increasing the image dimensions enhanced the DSC results significantly, as shown in Fig. 10. Next, we utilized a multi-scale training approach that progressively refined the training from coarse-to-fine scales. The network was first trained at the lowest scale, and trained parameters were used to initialize the training for the next image scale until the finest scale was trained. This strategy is named MS-DdC-AC-DLIR. Due to the low resolutions and comparatively poor performance of the 32×32 and 64×64 scales, we started training with 128×128 and refined it to 256×256 and 512×512 . Table 2, Fig. 5, Fig. 6, and Fig. 10 demonstrate the enhancement of DLIR results due to the multi-scale strategy. Compared to the DdC-AC-DLIR trained at 512×512 alone, TU was increased from 1.73 mm to 1.36 in the CAMUS A2C view and from 1.16 mm to 0.98 mm in the CAMUS A4C view. The other performance metrics, like DSC, HD,

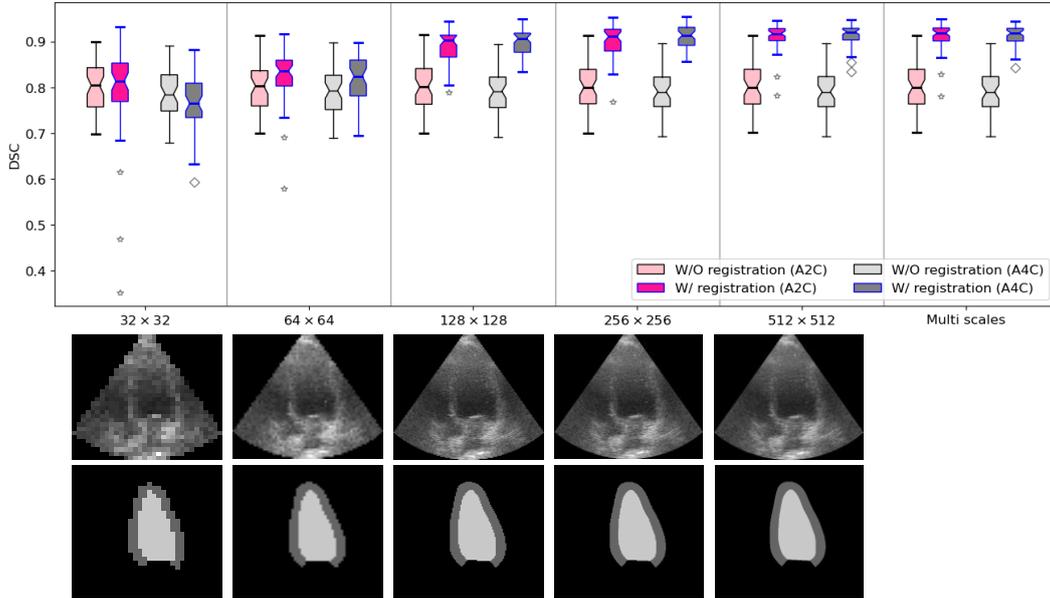


Figure 10: Advantages of multi-scale training over single-scale. The registration method barely increases or decreases DSC when image resolution is low; for example, 32×32 and 64×64 . Multi-scale training increases other measures (see Table 2) despite having a similar DSC to single 256×256 and 512×512 . CAMUS A2C and A4C views were used to create this diagram.

MSE, and LPIPS, either increased or remained almost similar. These enhancements can be observed visually in Fig. 5, which depicts enhanced mask overlap and a marked reduction in MYO border irregularity alongside an increase in perceptual similarity in deformed intensity images (I_W), especially for the CAMUS A2C view. The multi-scale training of DdC-AC-DLIR, as illustrated in Fig. 9 and Fig. B.15, not only results in an improvement in MYO’s thickness uniformity but also does not result in a reduction in the number of testing patients in the better DSC group.

4.5. Effect of Data Augmentation

The MS-DdC-AC-DLIR has undergone additional testing by incorporating geometric and intensity image augmentations such as motion blur, Gaussian blur, defocusing [67], contrast limited adaptive histogram equalization (CLAHE) [68], and horizontal flipping. Examples of augmented images are displayed in Fig. B.16 of Appendix B. This Aug-MS-DdC-AC-DLIR has led to a significant improvement ($p \ll 0.001$) in the thickness uniformity

of MYO, with the TU improving from 1.36 mm to 0.84 mm for A2C and 0.98 mm to 0.58 mm for A4C. The HD has decreased from 3.50 mm to 3.22 mm for A2C and from 3.34 mm to 3.16 mm for A4C, whereas the DSCs have comparable values with and without the addition of augmentations. As a result of these anatomical metric enhancements, the mask now overlaps more effectively, and there has been a significant reduction in the irregularity and coarseness of MYO’s inner and outer boundaries. Fig. B.17 of Appendix B illustrates the corresponding qualitative enhancements made to Aug-MS-DdC-AC-DLIR, comparing MS-DdC-AC-DLIR. Despite the improvement of pixel-level intensity differences in terms of MSE from 0.0016 mm to 0.0015 mm for A2C and from 0.0019 mm to 0.0018 mm for A4C, the LPIPS did not improve.

4.6. Temporal Registration

Next, we proceed to registration over a temporal range. To do this, we fixed the first time point as the fixed image and mask ($I_F(t_0)$ and $M_F(t_0)$). The other time points ($I_F(t_n)$ and $M_F(t_n)$) are registered to this first time point, i.e., $I_W(t_n \mapsto t_0) = I_M(t_n) \circ \varphi$ and $M_W(t_n \mapsto t_0) = M_M(t_n) \circ \varphi$. Results are shown in Table 3, Fig. 11, and Fig. B.18.

Generally, results were similar to those in Table 2’s ES to ED registration results. Table 3’s temporal registration results demonstrate progressive improving performance with the sequentially added strategies of anatomic constraints, data-driven constraints, and a multi-scale approach. The progressive improvements are the clearest for DSC, HD, and MSE, where MS-DdC-ACDLIR are the best in ($p \ll 0.001$ compared to VanDLIR). The progressive improvements are also evident for DSC from Appendix B (see Fig. B.18). However, for TU and LPIPS, progressive improvements are not observed. This is likely because our proposed strategies work well for large deformations, and here, many time frames have only small deformations, and the benefits are not as clearly observable.

Fig. 11 displays plausible and realistic warped moving images and masks for adults and fetal echo datasets using the MS-DdC-AC-DLIR algorithm. The figure demonstrates that the proposed DLIR provides a plausible deformed image and mask not only for the closest time point but also for distant moving time points. Supplementary videos show the warped

Table 3: Temporal image registration results for adult and fetal echo images, demonstrating the non-overlapping benefits of data-driven and anatomical constraints. All the metrics are estimated using fixed images (and masks) and warped moving images (and masks). The DSC and HD are the averages of backgrounds, MYO, and LV metrics, whereas the class-wise metrics of MYO and LV are shown in Table B.7 of Appendix B. Underlined, double-underlined, and bold fonts denote the best-performing metrics for the A2C view of adult echo, the A4C view of adult echo, and the A4C view of fetal echo, respectively.

Methods		DSC (\uparrow)	HD (mm) (\downarrow)	TU ‡ (mm)	MSE (\downarrow)	LPIPS (\downarrow)
N/A	Adult (A2C)	0.8942 \pm 0.0635	3.11 \pm 1.89	–	0.0017 \pm 0.0009	–
	Adult (A4C)	0.8878 \pm 0.0671	3.21 \pm 1.96	–	0.0021 \pm 0.0013	–
	Fetal (A4C)	0.9024 \pm 0.0665	1.20 \pm 0.64	–	0.0008 \pm 0.0008	–
VanDLIR	Adult (A2C)	0.9247 \pm 0.0470	2.83 \pm 1.61	1.04 \pm 1.02	<u>0.0011 \pm 0.0006</u>	0.0806 \pm 0.0145
	Adult (A4C)	0.9190 \pm 0.0490	3.27 \pm 1.80	0.80 \pm 0.74	0.0014 \pm 0.0009	0.0840 \pm 0.0166
	Fetal (A4C)	0.9365 \pm 0.0415	1.32 \pm 0.64	0.89 \pm 1.01	0.0008 \pm 0.0008	0.0828 \pm 0.0339
AC-DLIR	Adult (A2C)	0.9376 \pm 0.0344	2.52 \pm 1.39	<u>0.90 \pm 1.02</u>	0.0013 \pm 0.0007	0.0806 \pm 0.0150
	Adult (A4C)	0.9278 \pm 0.04304	2.80 \pm 1.57	0.72 \pm 0.63	0.0017 \pm 0.0010	0.0838 \pm 0.0167
	Fetal (A4C)	0.9425 \pm 0.0358	1.32 \pm 0.62	1.01 \pm 1.18	0.0008 \pm 0.0008	0.0868 \pm 0.0326
DdC-AC-DLIR	Adult (A2C)	0.9378 \pm 0.0356	2.51 \pm 1.41	<u>0.90 \pm 1.09</u>	0.0012 \pm 0.0006	<u>0.0803 \pm 0.0146</u>
	Adult (A4C)	0.9269 \pm 0.0437	2.73 \pm 1.51	<u>0.68 \pm 0.54</u>	0.0017 \pm 0.0011	0.0847 \pm 0.0166
	Fetal (A4C)	0.9489 \pm 0.0291	1.24 \pm 0.57	1.10 \pm 1.28	0.0006 \pm 0.0007	0.0827 \pm 0.0337
MS-DdC-AC-DLIR	Adult (A2C)	<u>0.9413 \pm 0.0262</u>	<u>2.29 \pm 1.08</u>	1.04 \pm 1.15	<u>0.0011 \pm 0.0005</u>	0.0857 \pm 0.0167
	Adult (A4C)	<u>0.9353 \pm 0.0288</u>	<u>2.38 \pm 1.16</u>	0.83 \pm 0.73	0.0013 \pm 0.0007	0.0899 \pm 0.0193
	Fetal (A4C)	0.9523 \pm 0.0261	1.17 \pm 0.54	1.11 \pm 1.08	0.0005 \pm 0.0006	0.0826 \pm 0.0304

‡ An accurate TU should be as close as the fixed MYO’s TU , which is 0.74 mm for adults in A2C, 0.67 mm for adults in A4C, and 1.18 mm for fetal in A4C views.

moving image over all the time points, demonstrating reasonable temporal consistency of the warped MYO mask (Video 1 for CAMUS-A2C⁷, Video 2 for CAMUS-A4C⁸, and Video 3 for Fetal-A4C⁹). The videos also demonstrate good anatomical plausibility and perceptual

⁷CAMUS-A2C: <https://youtu.be/EXTlyImIGgA>

⁸CAMUS-A4C: <https://youtu.be/l6ua9Qrc3JE>

⁹Fetal-A4C: <https://youtu.be/eGUU-rqWznY>

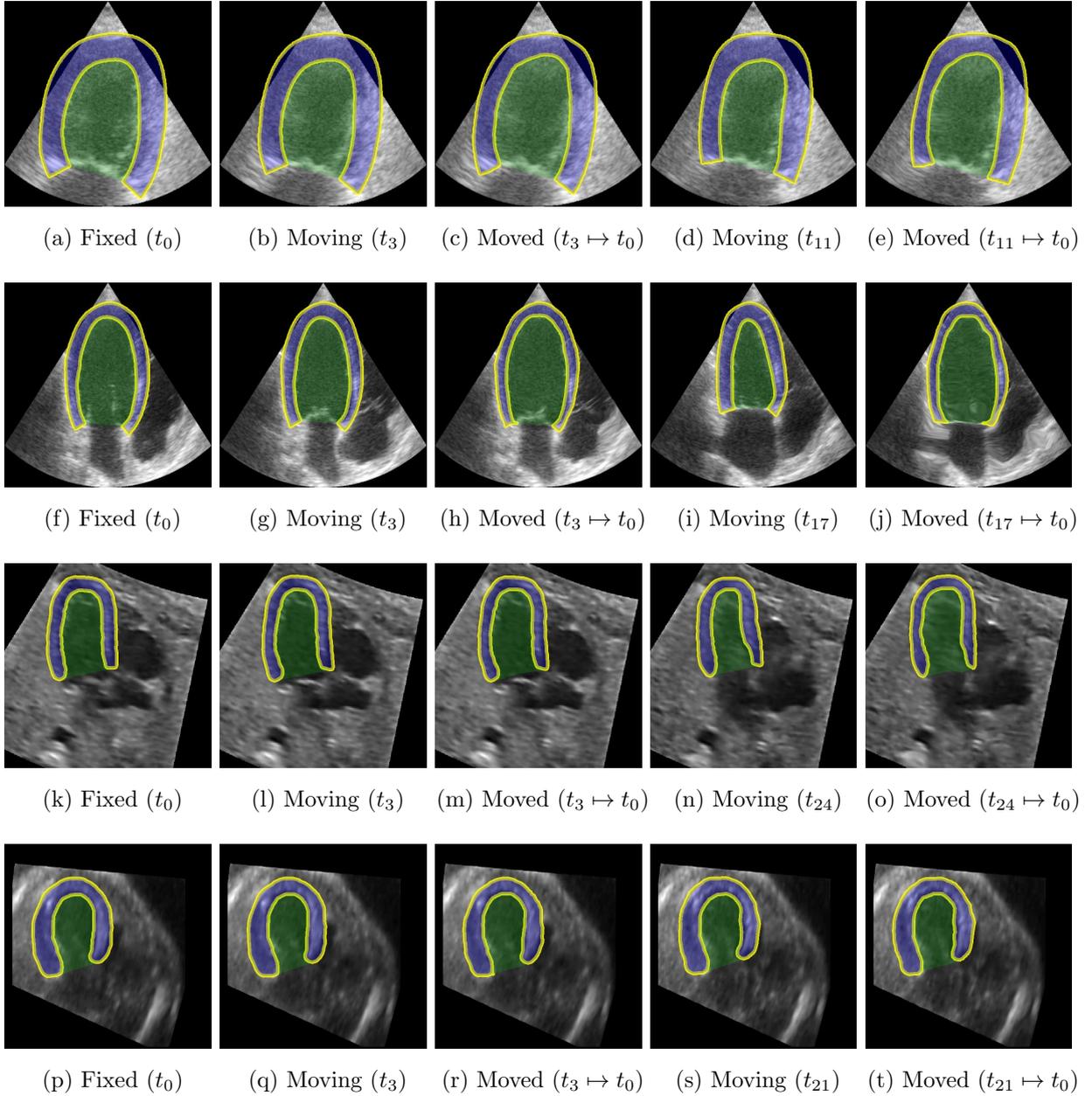


Figure 11: Example of qualitative temporal registration results for the public CAMUS echoes dataset, with the first row for A2C view (see video⁷) and the second row for A4C view (see video⁸). The last two rows are for our private fetal echo A4C view (see video⁹), where the third row is for the healthy patient and the fourth row is for the diseased patient. Two distinct time positions, one near and one far from the fixed time point, are aligned with it.

realism in the images.

5. Estimating Ejection Fraction (EF)

Finally, we evaluate the various DLIR algorithms for their ability to estimate cardiac EF using Eq. 8. Regression analysis of the predicted and actual EF is shown in Fig. 12. Here,

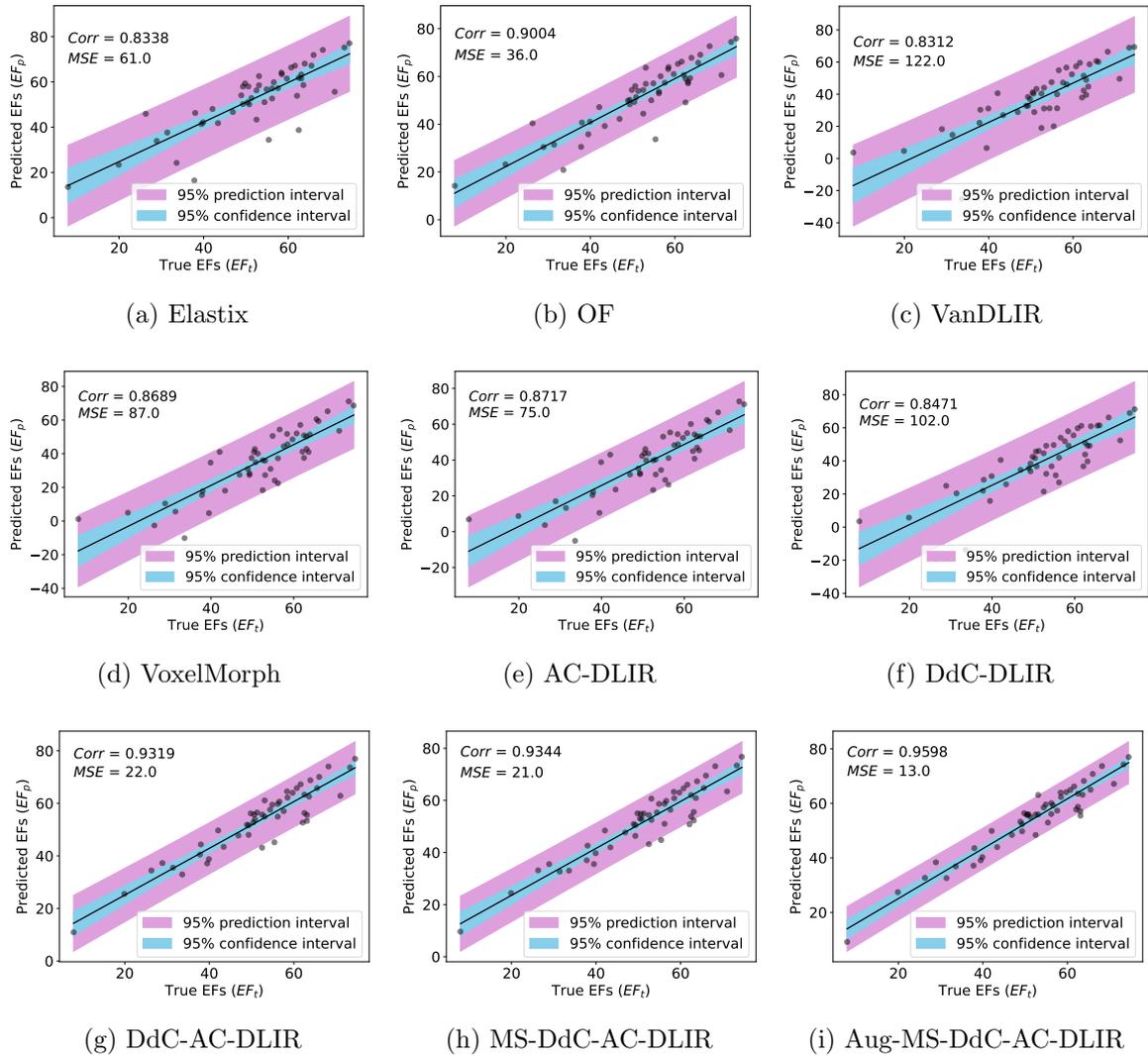


Figure 12: Regression plots between the true and estimated EFs from different registration methods to determine the degree of agreement in terms of the correlation between these two EFs. Those figures were created using the CAMUS A2C view.

Elastix and VanDLIR had similar performance in terms of correlation coefficients of around

0.83, but OF achieved 0.90. There is again a progressive enhancement of the correlation coefficient, moving from VanDLIR to the progressive more complex techniques. Interestingly, OF performs better than most of the basic DLIR models, and only DdC-AC-DLIR and its multi-scale and augmented multi-scale versions perform better. This thus provides another validation for our proposed approach. Additionally, it is observed that VanDLIR provides negative EFs for two out of fifty testing patients, which is physiologically impossible and likely due to poor image quality. Using only local anatomical constraints in VoxelMorph and DdC-DLIR, they remain negative, but adding global anatomical constraints in AC-DLIR results in reduced negative values for one patient. However, with DdC-AC-DLIR and its multi-scale and augmented versions, negative results are resolved. Results thus suggest that superior registration metrics, provided our proposed combination constraints, can translate to better clinical quantifications.

6. Discussion and Conclusion

Temporal registration of echo images is an important foundational step for extracting clinically relevant features from images. For example, it can be used for estimating stroke volume and ejection fraction, myocardial strains and strain rates, etc. However, the inherent high noise level, low contrast, and limited image resolution, which lead to fuzzy anatomical boundaries, make the registration process challenging, especially for fetuses whose hearts are small and far from the transducer. It is thus important to optimize registration results so that clinical measurements can have manageably low errors and reasonable precision.

Our proposed strategies can help bridge these gaps. We show that the proposed strategies can achieve very robust registration metrics, as demonstrated in Tables 2 and 3, and we show that this, in turn, can translate to better quality clinical measurements, as demonstrated in Fig. 12 on EF quantification. These good results are due to our use of a combination of strategies, where the anatomic constraint ensures better-warped anatomy while the adversarial data-driven constraint ensures better visual quality of the warped echo image. As our ablation studies show, this combination provides non-overlapping benefits to achieve

excellent final outcomes. These results also demonstrate that such combinations can be a good algorithm design approach.

Our strategy revolves around enforcing high-quality anatomic shapes and image quality in the warped images. The rationale for this is that echo images are generated from ultrasound physics and echogenic material in tissues and are governed by realistic in vivo anatomy and motions; thus, warped images should achieve similarly physiological cardiac shapes and should have realistic echo images. Requiring warped images to have physiological cardiac good echo image textures can be argued to be a way of enforcing physiological deformations and thus a way of gaining accuracy. Traditionally, anatomic constraints are used in segmentation networks, such as those by Oktay et al. [69], where anatomy latent spaces are used as constraints for segmentation training. Our study shows that it can also be a valuable strategy to improve registration.

Our Aug-MS-DdC-AC-DLIR approach provides improvements on traditional registration methods that are considered the gold standard, highlighting the benefits of adopting deep learning approaches. The deep learning approach also provides substantial computational time savings. On the same machine, DLIR took a significantly shorter time to compute a single deformation field, approximately 80–95 ms, while Elastix and OF require 4270 ms and 937 ms, respectively. Our performance compares well with good techniques in the literature. For example, Wei et al. [41] proposed the CLAS approach that concurrently performed segmentation and registration across the whole cycle with the 3D UNet and tested it on the same CAMUS dataset. They find that the mean DSC of A2C and A4C views from object tracking of the LV chamber were 0.923 and 0.903 for end-diastole and end-systole, which are lower than our DSC for the LV chamber of 0.9261 for the A2C view and 0.9360 for the A4C view. Wang et al. [24] proposed a patch-based MLP and transformers for registration and reported similar DSC scores for both the LV chamber and MYO region as our study. However, we achieve better HD results, likely due to our anatomical constraints. Fan et al. [25] used an unsupervised multi-scale correlation iterative registration network (SearchMorph) with a correlation layer. The author finds that the mean DSCs for MYO and LV are 0.880 and 0.888 for the A2C view and 0.891 and 0.919 for the A4C view, whereas

for the same CAMUS dataset, those values from our DLIR are 0.881 and 0.953 for the A2C view and 0.864 and 0.951 for the A4C view.

Thus, in conclusion, a combination of anatomical constraints via a shape encoder, an adversarially trained data-driven constraint, and multi-scale training can produce excellent image registration results for both adult and fetal echocardiography, which can translate to higher-quality clinical measurements.

Acknowledgement

Md. Kamrul Hasan was supported by the doctoral training program (DTP) studentship funds of the Engineering and Physical Sciences Research Council (EPSRC) (2022-2025).

Guang Yang was supported in part by the ERC IMI (101005122), the H2020 (952172), the MRC (MC/PC/21013), the Royal Society (IEC/NSFC/211235), the NVIDIA Academic Hardware Grant Program, the SABER project supported by Boehringer Ingelheim Ltd, and the UKRI Future Leaders Fellowship (MR/V023799/1).

References

- [1] P. Claus, A. M. S. Omar, G. Pedrizzetti, P. P. Sengupta, E. Nagel, Tissue tracking technology for assessing cardiac mechanics: principles, normal values, and clinical applications, *JACC: Cardiovascular Imaging* 8 (2015) 1444–1460.
- [2] Z. B. Popović, D. H. Kwon, M. Mishra, A. Buakhamsri, N. L. Greenberg, M. Thamarasan, S. D. Flamm, J. D. Thomas, H. M. Lever, M. Y. Desai, Association between regional ventricular function and myocardial fibrosis in hypertrophic cardiomyopathy assessed by speckle tracking echocardiography and delayed hyperenhancement magnetic resonance imaging, *Journal of the American Society of Echocardiography* 21 (2008) 1299–1305.
- [3] J. M. Balter, M. L. Kessler, Imaging and alignment for image-guided radiation therapy, *Journal of clinical oncology* 25 (2007) 931–937.
- [4] D. Seo, J. Ho, J. H. Traverse, J. Forder, B. Vemuri, Computing diffeomorphic paths with applications to cardiac motion analysis, in: *4th MICCAI Workshop on Mathematical Foundations of Computational Anatomy*, Citeseer, pp. 83–94.
- [5] B. Heyde, R. Jasaityte, D. Barbosa, V. Robesyn, S. Bouchez, P. Wouters, F. Maes, P. Claus, J. D’hooge,

- Elastic image registration versus speckle tracking for 2-d myocardial motion estimation: a direct comparison in vivo, *IEEE transactions on medical imaging* 32 (2012) 449–459.
- [6] H. S. Wong, B. Li, A. Tulzer, G. Tulzer, C. H. Yap, Fluid mechanical effects of fetal aortic valvuloplasty for cases of critical aortic stenosis with evolving hypoplastic left heart syndrome, *Annals of Biomedical Engineering* (2023) 1–14.
- [7] C. W. Ong, M. Ren, H. Wiputra, J. Mojumder, W. X. Chan, A. Tulzer, G. Tulzer, M. L. Buist, C. N. Z. Mattar, L. C. Lee, et al., Biomechanics of human fetal hearts with critical aortic stenosis, *Annals of biomedical engineering* 49 (2021) 1364–1379.
- [8] N. H. van Oostrum, C. M. de Vet, D. A. van der Woude, H. M. Kemps, S. G. Oei, J. O. van Laar, Fetal strain and strain rate during pregnancy measured with speckle tracking echocardiography: A systematic review, *European Journal of Obstetrics & Gynecology and Reproductive Biology* 250 (2020) 178–187.
- [9] H. Wiputra, W. X. Chan, Y. Y. Foo, S. Ho, C. H. Yap, Cardiac motion estimation from medical images: a regularisation framework applied on pairwise image registration displacement fields, *Scientific reports* 10 (2020) 18510.
- [10] M. De Craene, G. Piella, O. Camara, N. Duchateau, E. Silva, A. Doltra, J. D’hooge, J. Brugada, M. Sitges, A. F. Frangi, Temporal diffeomorphic free-form deformation: Application to motion and strain estimation from 3d echocardiography, *Medical image analysis* 16 (2012) 427–450.
- [11] J.-P. Thirion, Image matching as a diffusion process: an analogy with maxwell’s demons, *Medical image analysis* 2 (1998) 243–260.
- [12] S. S. Beauchemin, J. L. Barron, The computation of optical flow, *ACM computing surveys (CSUR)* 27 (1995) 433–466.
- [13] G. Balakrishnan, A. Zhao, M. R. Sabuncu, J. Guttag, A. V. Dalca, Voxelmorph: a learning framework for deformable medical image registration, *IEEE transactions on medical imaging* 38 (2019) 1788–1800.
- [14] X. Chen, A. Diaz-Pinto, N. Ravikumar, A. F. Frangi, Deep learning in medical image registration, *Progress in Biomedical Engineering* 3 (2021) 012003.
- [15] A. Dosovitskiy, P. Fischer, E. Ilg, P. Hausser, C. Hazirbas, V. Golkov, P. Van Der Smagt, D. Cremers, T. Brox, FlowNet: Learning optical flow with convolutional networks, in: *Proceedings of the IEEE international conference on computer vision*, pp. 2758–2766.
- [16] M.-M. Rohé, M. Datar, T. Heimann, M. Sermesant, X. Pennec, Svf-net: learning deformable image registration using shape matching, in: *Medical Image Computing and Computer Assisted Intervention-MICCAI 2017: 20th International Conference, Quebec City, QC, Canada, September 11-13, 2017, Proceedings, Part I* 20, Springer, pp. 266–274.
- [17] H. Sokooti, B. De Vos, F. Berendsen, B. P. Lelieveldt, I. Išgum, M. Staring, Nonrigid image registration using multi-scale 3d convolutional neural networks, in: *Medical Image Computing and Com-*

- puter Assisted Intervention- MICCAI 2017: 20th International Conference, Quebec City, QC, Canada, September 11-13, 2017, Proceedings, Part I 20, Springer, pp. 232–239.
- [18] X. Yang, R. Kwitt, M. Styner, M. Niethammer, Quicksilver: Fast predictive image registration—a deep learning approach, *NeuroImage* 158 (2017) 378–396.
- [19] A. Østvik, I. M. Salte, E. Smistad, T. M. Nguyen, D. Melichova, H. Brunvand, K. Haugaa, T. Edvardsen, B. Grenne, L. Lovstakken, Myocardial function imaging in echocardiography using deep learning, *IEEE transactions on medical imaging* 40 (2021) 1340–1351.
- [20] L. Mansilla, D. H. Milone, E. Ferrante, Learning deformable registration of medical images with anatomical constraints, *Neural Networks* 124 (2020) 269–279.
- [21] S. Ali, J. Rittscher, Conv2warp: An unsupervised deformable image registration with continuous convolution and warping, in: *Machine Learning in Medical Imaging: 10th International Workshop, MLMI 2019, Held in Conjunction with MICCAI 2019, Shenzhen, China, October 13, 2019, Proceedings 10*, Springer, pp. 489–497.
- [22] Y. Hu, M. Modat, E. Gibson, N. Ghavami, E. Bonmati, C. M. Moore, M. Emberton, J. A. Noble, D. C. Barratt, T. Vercauteren, Label-driven weakly-supervised learning for multimodal deformable image registration, in: *2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018)*, IEEE, pp. 1070–1074.
- [23] Z. Xu, M. Niethammer, Deepatlas: Joint semi-supervised learning of image registration and segmentation, in: *Medical Image Computing and Computer Assisted Intervention—MICCAI 2019: 22nd International Conference, Shenzhen, China, October 13–17, 2019, Proceedings, Part II 22*, Springer, pp. 420–429.
- [24] Z. Wang, Y. Yang, M. Sermesant, H. Delingette, Unsupervised echocardiography registration through patch-based mlps and transformers, in: *International Workshop on Statistical Atlases and Computational Models of the Heart*, Springer, pp. 168–178.
- [25] X. Fan, S. Zhuang, Z. Zhuang, Y. Yuan, S. Qiu, A. N. J. Raj, Y. Rong, Searchmorph: Multi-scale correlation iterative network for deformable registration, *arXiv:2206.13076* (2022).
- [26] L. Vasciaveo, E. Zanzarelli, F. D’Antonio, Fetal cardiac function evaluation: A review, *Journal of Clinical Ultrasound* 51 (2023) 215–224.
- [27] X. Yang, R. Kwitt, M. Niethammer, Fast predictive image registration, in: *Deep Learning and Data Labeling for Medical Applications: First International Workshop, LABELS 2016, and Second International Workshop, DLMIA 2016, Held in Conjunction with MICCAI 2016, Athens, Greece, October 21, 2016, Proceedings 1*, Springer, pp. 48–57.
- [28] X. Cao, J. Yang, J. Zhang, D. Nie, M. Kim, Q. Wang, D. Shen, Deformable image registration based on similarity-steered cnn regression, in: *Medical Image Computing and Computer Assisted Intervention-*

- MICCAI 2017: 20th International Conference, Quebec City, QC, Canada, September 11-13, 2017, Proceedings, Part I 20, Springer, pp. 300–308.
- [29] J. Fan, X. Cao, P.-T. Yap, D. Shen, Birnet: Brain image registration using dual-supervised fully convolutional networks, *Medical image analysis* 54 (2019) 193–206.
- [30] K. A. Eppenhof, J. P. Pluim, Pulmonary ct registration through supervised learning with convolutional neural networks, *IEEE transactions on medical imaging* 38 (2018) 1097–1105.
- [31] C. Lian, X. Li, L. Kong, J. Wang, W. Zhang, X. Huang, L. Wang, Cocyclereg: Collaborative cycle-consistency method for multi-modal medical image registration, *Neurocomputing* 500 (2022) 799–808.
- [32] B. Kim, D. H. Kim, S. H. Park, J. Kim, J.-G. Lee, J. C. Ye, Cyclemorph: cycle consistent unsupervised deformable image registration, *Medical image analysis* 71 (2021) 102036.
- [33] P. Yan, S. Xu, A. R. Rastinehad, B. J. Wood, Adversarial image registration with application for mr and trus image fusion, in: *Machine Learning in Medical Imaging: 9th International Workshop, MLMI 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 16, 2018, Proceedings 9*, Springer, pp. 197–204.
- [34] Z. He, Y. He, W. Cao, Deformable image registration with attention-guided fusion of multi-scale deformation fields, *Applied Intelligence* (2022) 1–15.
- [35] N. Dey, M. Ren, A. V. Dalca, G. Gerig, Generative adversarial registration for improved conditional deformable templates, in: *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 3929–3941.
- [36] J.-Y. Zhu, T. Park, P. Isola, A. A. Efros, Unpaired image-to-image translation using cycle-consistent adversarial networks, in: *Proceedings of the IEEE international conference on computer vision*, pp. 2223–2232.
- [37] D. Mahapatra, B. Antony, S. Sedai, R. Garnavi, Deformable medical image registration using generative adversarial networks, in: *2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018)*, IEEE, pp. 1449–1453.
- [38] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio, Generative adversarial networks, *Communications of the ACM* 63 (2020) 139–144.
- [39] S. Czolbe, P. Pegios, O. Krause, A. Feragen, Semantic similarity metrics for image registration, *Medical Image Analysis* (2023) 102830.
- [40] T. Haukom, E. A. R. Berg, S. Aakhus, G. H. Kiss, Basal strain estimation in transesophageal echocardiography (tee) using deep learning based unsupervised deformable image registration, in: *2019 IEEE International Ultrasonics Symposium (IUS)*, IEEE, pp. 1421–1424.
- [41] H. Wei, H. Cao, Y. Cao, Y. Zhou, W. Xue, D. Ni, S. Li, Temporal-consistent segmentation of echocardiography with co-learning from appearance and shape, in: *Medical Image Computing and Computer*

- Assisted Intervention–MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part II 23, Springer, pp. 623–632.
- [42] D. Sun, X. Yang, M.-Y. Liu, J. Kautz, Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 8934–8943.
- [43] M. J. Cardoso, W. Li, R. Brown, N. Ma, E. Kerfoot, Y. Wang, B. Murrey, A. Myronenko, C. Zhao, D. Yang, et al., Monai: An open-source framework for deep learning in healthcare, arXiv:2211.02701 (2022).
- [44] O. Ronneberger, P. Fischer, T. Brox, U-net: Convolutional networks for biomedical image segmentation, in: Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18, Springer, pp. 234–241.
- [45] Y. Fu, N. M. Brown, S. U. Saeed, A. Casamitjana, Z. Baum, R. Delaunay, Q. Yang, A. Grimwood, Z. Min, S. B. Blumberg, et al., Deepreg: a deep learning toolkit for medical image registration, arXiv:2011.02580 (2020).
- [46] D. P. Kingma, M. Welling, et al., An introduction to variational autoencoders, Foundations and Trends® in Machine Learning 12 (2019) 307–392.
- [47] S. Zhao, B. Wu, W. Chu, Y. Hu, D. Cai, Correlation maximized structural similarity loss for semantic segmentation, arXiv:1910.08711 (2019).
- [48] C. Doersch, Tutorial on variational autoencoders, arXiv:1606.05908 (2016).
- [49] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, O. Wang, The unreasonable effectiveness of deep features as a perceptual metric, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 586–595.
- [50] X. Hou, L. Shen, K. Sun, G. Qiu, Deep feature consistent variational autoencoder, in: 2017 IEEE winter conference on applications of computer vision (WACV), IEEE, pp. 1133–1141.
- [51] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 770–778.
- [52] D. Ulyanov, A. Vedaldi, V. Lempitsky, Instance normalization: The missing ingredient for fast stylization, arXiv:1607.08022 (2016).
- [53] B. Xu, N. Wang, T. Chen, M. Li, Empirical evaluation of rectified activations in convolutional network, arXiv:1505.00853 (2015).
- [54] W. Zhu, Y. Huang, D. Xu, Z. Qian, W. Fan, X. Xie, Test-time training for deformable multi-scale image registration, in: 2021 IEEE International Conference on Robotics and Automation (ICRA), IEEE, pp. 13618–13625.
- [55] Z. Jiang, F.-F. Yin, Y. Ge, L. Ren, A multi-scale framework with unsupervised joint training of

- convolutional neural networks for pulmonary deformable image registration, *Physics in Medicine & Biology* 65 (2020) 015011.
- [56] S. Shao, Z. Pei, W. Chen, W. Zhu, X. Wu, B. Zhang, A multi-scale unsupervised learning for deformable image registration, *International Journal of Computer Assisted Radiology and Surgery* (2022) 1–10.
- [57] B. D. De Vos, F. F. Berendsen, M. A. Viergever, H. Sokooti, M. Staring, I. Išgum, A deep learning framework for unsupervised affine and deformable image registration, *Medical image analysis* 52 (2019) 128–143.
- [58] S. Klein, M. Staring, K. Murphy, M. A. Viergever, J. P. Pluim, Elastix: a toolbox for intensity-based medical image registration, *IEEE transactions on medical imaging* 29 (2009) 196–205.
- [59] W. X. Chan, Y. Zheng, H. Wiputra, H. L. Leo, C. H. Yap, Full cardiac cycle asynchronous temporal compounding of 3d echocardiography images, *Medical Image Analysis* 74 (2021) 102229.
- [60] B. D. Lucas, T. Kanade, An iterative image registration technique with an application to stereo vision, in: *IJCAI'81: 7th international joint conference on Artificial intelligence*, volume 2, pp. 674–679.
- [61] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al., Pytorch: An imperative style, high-performance deep learning library, *Advances in neural information processing systems* 32 (2019).
- [62] O. U. Aydin, A. A. Taha, A. Hilbert, A. A. Khalil, I. Galinovic, J. B. Fiebach, D. Frey, V. I. Madai, On the usage of average hausdorff distance for segmentation performance assessment: hidden error when used for ranking, *European Radiology Experimental* 5 (2021) 1–7.
- [63] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, *arXiv:1409.1556* (2014).
- [64] S. Leclerc, E. Smistad, J. Pedrosa, A. Østvik, F. Cervenansky, F. Espinosa, T. Espeland, E. A. R. Berg, P.-M. Jodoin, T. Grenier, et al., Deep learning for segmentation using an open large-scale dataset in 2d echocardiography, *IEEE transactions on medical imaging* 38 (2019) 2198–2210.
- [65] W. Xue, H. Cao, J. Ma, T. Bai, T. Wang, D. Ni, Improved segmentation of echocardiography with orientation-congruency of optical flow and motion-enhanced segmentation, *IEEE Journal of Biomedical and Health Informatics* 26 (2022) 6105–6115.
- [66] C. Sfakianakis, G. Simantiris, G. Tziritas, Gudu: Geometrically-constrained ultrasound data augmentation in u-net for echocardiography semantic segmentation, *Biomedical Signal Processing and Control* 82 (2023) 104557.
- [67] D. Hendrycks, T. Dietterich, Benchmarking neural network robustness to common corruptions and perturbations, *arXiv:1903.12261* (2019).
- [68] G. Yadav, S. Maheshwari, A. Agarwal, Contrast limited adaptive histogram equalization based enhancement for real time video system, in: *2014 international conference on advances in computing*,

communications and informatics (ICACCI), IEEE, pp. 2392–2397.

- [69] 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* 37 (2017) 384–395.

Appendix A. Structure of Networks

Table A.4: The architecture of the topological shape encoder to learn the latent vectors of LV and MYO. Conv: Convolution layer with a dilation of 1; BN: Batch normalization for each mini-batch; ReLU: Rectified linear unit; FC: Fully connected layer; TConv: Transpose convolution layer for the upsampling with a dilation of 1; SIG: Sigmoid activation at the output; B: Mini-batch size; and d : Latent vector dimension (\mathcal{Z}^d).

Layer	Feature dimension ($B \times \text{n.filters} \times \text{height} \times \text{width}$)	Parameters	Kernel	Stride	Padding
Encoding mask topology to Latent vector (\mathcal{Z}^d) using Encoder (\mathcal{EN})					
Input (I)	$B \times 2^\dagger \times H \times W$	–	–	–	–
Conv	$B \times 8 \times H/2 \times W/2$	$(3 \times 3 \times 2 + 1) \times 8$	3×3	2×2	1×1
ReLU	$B \times 8 \times H/2 \times W/2$	–	–	–	–
Conv	$B \times 16 \times H/4 \times W/4$	$(3 \times 3 \times 8 + 1) \times 16$	3×3	2×2	1×1
BN	$B \times 16 \times H/4 \times W/4$	32	–	–	–
ReLU	$B \times 16 \times H/4 \times W/4$	–	–	–	–
Conv	$B \times 32 \times H/8 \times W/8$	$(3 \times 3 \times 16 + 1) \times 32$	3×3	2×2	1×1
ReLU	$B \times 32 \times H/8 \times W/8$	–	–	–	–
FC	$B \times 128$	$(32 \times H/8 \times W/8 \times 128) + 128$	–	–	–
ReLU	$B \times 128$	–	–	–	–
FC	$B \times 128$	$(128 \times 128) + 128$	–	–	–
ReLU	$B \times 128$	–	–	–	–
FC	$B \times d$	$(128 \times d) + d$	–	–	–
ReLU	$B \times d$	–	–	–	–
Reconstructing mask topology from Latent vector (\mathcal{Z}^d) using Decoder (\mathcal{DE})					
FC	$B \times 128$	$(128 \times d) + 128$	–	–	–
ReLU	$B \times 128$	–	–	–	–
FC	$B \times (32 \times H/8 \times W/8)$	$(32 \times H/8 \times W/8 \times 128) + 32 \times H/8 \times W/8$	–	–	–
ReLU	$B \times (32 \times H/8 \times W/8)$	–	–	–	–
TConv	$B \times 16 \times H/4 \times W/4$	$(3 \times 3 \times 32 + 1) \times 16$	3×3	2×2	1×1
BN	$B \times 16 \times H/4 \times W/4$	32	–	–	–
ReLU	$B \times 16 \times H/4 \times W/4$	–	–	–	–
TConv	$B \times 8 \times H/2 \times W/2$	$(3 \times 3 \times 16 + 1) \times 8$	3×3	2×2	1×1
BN	$B \times 8 \times H/2 \times W/2$	16	–	–	–
ReLU	$B \times 8 \times H/2 \times W/2$	–	–	–	–
TConv	$B \times 2 \times H \times W$	$(3 \times 3 \times 8 + 1) \times 2$	3×3	2×2	1×1
ReLU	$B \times 2 \times H \times W$	–	–	–	–
SIG	$B \times 2 \times H \times W$	–	–	–	–
Output (I')	$B \times 2^\dagger \times H \times W$	–	–	–	–

[†]This article focuses on LV's and MYO's topologies of the adult and fetal hearts.

Table A.5: Discriminator’s architecture to learn the intensity images’ attributes (I_F and $I_W = I_M \circ \varphi$) in adversarial training to retrieve the texture features in the moved intensity images ($I_W = I_M \circ \varphi$). Conv: Convolution layer with a dilation of 1; IN: Instance normalization for each mini-batch; LReLU: Leaky Rectified linear unit; FC: Fully connected layer; SIG: Sigmoid activation at the output; and B: Mini-batch size.

Layer	Feature dimension (B \times n.filters \times height \times width)	Parameters	Kernel	Stride	Padding
2D feature learning blocks					
Input	$B \times 1 \times H \times W$	–	–	–	–
Conv	$B \times 8 \times H/2 \times W/2$	$(3 \times 3 \times 1 + 1) \times 8$	3×3	2×2	1×1
Conv	$B \times 8 \times H/2 \times W/2$	$(3 \times 3 \times 1 + 1) \times 8$	3×3	2×2	1×1
IN	$B \times 8 \times H/2 \times W/2$	–	–	–	–
Dropout (10%)	$B \times 8 \times H/2 \times W/2$	–	–	–	–
LReLU	$B \times 8 \times H/2 \times W/2$	–	–	–	–
Conv	$B \times 8 \times H/2 \times W/2$	$(3 \times 3 \times 8 + 1) \times 8$	3×3	1×1	1×1
IN	$B \times 8 \times H/2 \times W/2$	–	–	–	–
Dropout (10%)	$B \times 8 \times H/2 \times W/2$	–	–	–	–
LReLU	$B \times 8 \times H/2 \times W/2$	–	–	–	–
Conv	$B \times 16 \times H/4 \times W/4$	$(3 \times 3 \times 8 + 1) \times 16$	3×3	2×2	1×1
Conv	$B \times 16 \times H/4 \times W/4$	$(3 \times 3 \times 8 + 1) \times 16$	3×3	2×2	1×1
IN	$B \times 16 \times H/4 \times W/4$	–	–	–	–
Dropout (10%)	$B \times 16 \times H/4 \times W/4$	–	–	–	–
LReLU	$B \times 16 \times H/4 \times W/4$	–	–	–	–
Conv	$B \times 16 \times H/4 \times W/4$	$(3 \times 3 \times 16 + 1) \times 16$	3×3	1×1	1×1
IN	$B \times 16 \times H/4 \times W/4$	–	–	–	–
Dropout (10%)	$B \times 16 \times H/4 \times W/4$	–	–	–	–
LReLU	$B \times 16 \times H/4 \times W/4$	–	–	–	–
Conv	$B \times 32 \times H/8 \times W/8$	$(3 \times 3 \times 16 + 1) \times 32$	3×3	2×2	1×1
Conv	$B \times 32 \times H/8 \times W/8$	$(3 \times 3 \times 16 + 1) \times 32$	3×3	2×2	1×1
IN	$B \times 32 \times H/8 \times W/8$	–	–	–	–
Dropout (10%)	$B \times 32 \times H/8 \times W/8$	–	–	–	–
LReLU	$B \times 32 \times H/8 \times W/8$	–	–	–	–
Conv	$B \times 32 \times H/8 \times W/8$	$(3 \times 3 \times 32 + 1) \times 32$	3×3	1×1	1×1
IN	$B \times 32 \times H/8 \times W/8$	–	–	–	–
Dropout (10%)	$B \times 32 \times H/8 \times W/8$	–	–	–	–
LReLU	$B \times 32 \times H/8 \times W/8$	–	–	–	–
Conv	$B \times 64 \times H/16 \times W/16$	$(3 \times 3 \times 32 + 1) \times 64$	3×3	2×2	1×1
Conv	$B \times 64 \times H/16 \times W/16$	$(3 \times 3 \times 32 + 1) \times 64$	3×3	2×2	1×1
IN	$B \times 64 \times H/16 \times W/16$	–	–	–	–
Dropout (10%)	$B \times 64 \times H/16 \times W/16$	–	–	–	–
LReLU	$B \times 64 \times H/16 \times W/16$	–	–	–	–
Conv	$B \times 64 \times H/16 \times W/16$	$(3 \times 3 \times 64 + 1) \times 64$	3×3	1×1	1×1
IN	$B \times 64 \times H/16 \times W/16$	–	–	–	–
Dropout (10%)	$B \times 64 \times H/16 \times W/16$	–	–	–	–
LReLU	$B \times 64 \times H/16 \times W/16$	–	–	–	–
Multi-layer perceptron for the classification of I_F and $I_W = I_M \circ \varphi$					
Flatten	$B \times (64 \times H/16 \times W/16)$	$64 \times H/16 \times W/16$	–	–	–
FC	$B \times 1024$	$64 \times H/16 \times W/16 \times 1024 + 1024$	–	–	–
FC	$B \times 256$	$1024 \times 256 + 256$	–	–	–
FC+SIG	$B \times 1$	$256 \times 1 + 1$	–	–	–
Output	$B \times 1$	–	–	–	–

Appendix B. Additional Results

Table B.6: Additional registration results of the CAMUS adult echo dataset from different registration techniques, demonstrating the metrics for MYO and LV, where Table 2 shows the average metrics of background, MYO, and LV. All the metrics are estimated using fixed images (and masks) and warped moving images (and masks). Bold fonts denote the best-performing metrics for the A2C echo view, while the best-performing metrics for the A4C view are underlined.

Methods		DSC (\uparrow)		HD (mm) (\downarrow)	
		MYO	LV	MYO	LV
Elastix [58]	A2C	0.7809 \pm 0.0768	0.8906 \pm 0.0669	6.04 \pm 3.10	5.07 \pm 3.47
	A4C	0.7792 \pm 0.0682	0.9096 \pm 0.0502	5.64 \pm 3.00	4.18 \pm 3.17
OF [60]	A2C	0.8140 \pm 0.0456	0.8990 \pm 0.0520	5.32 \pm 1.96	4.10 \pm 1.93
	A4C	0.8029 \pm 0.0651	0.9116 \pm 0.0425	5.49 \pm 3.42	4.05 \pm 3.77
VanDLIR	A2C	0.7294 \pm 0.0731	0.8512 \pm 0.0719	6.78 \pm 2.17	4.48 \pm 1.69
	A4C	0.7173 \pm 0.0713	0.8519 \pm 0.0596	6.42 \pm 2.02	5.05 \pm 1.68
VoxelMorph [13]	A2C	0.7838 \pm 0.0609	0.8443 \pm 0.0641	6.39 \pm 1.92	3.81 \pm 1.30
	A4C	0.7316 \pm 0.0769	0.8448 \pm 0.0619	5.79 \pm 1.96	4.23 \pm 1.48
AC-DLIR	A2C	0.8079 \pm 0.0526	0.8649 \pm 0.0581	5.40 \pm 1.93	3.49 \pm 1.17
	A4C	0.7925 \pm 0.0580	0.8832 \pm 0.0517	4.77 \pm 1.73	3.41 \pm 1.27
DdC-DLIR	A2C	0.7704 \pm 0.0674	0.8701 \pm 0.0688	5.80 \pm 2.56	4.12 \pm 1.86
	A4C	0.7525 \pm 0.0681	0.8699 \pm 0.0559	6.23 \pm 2.24	4.82 \pm 2.03
DdC-AC-DLIR	A2C	0.8247 \pm 0.0455	0.9225 \pm 0.0409	4.51 \pm 1.90	2.78 \pm 1.37
	A4C	<u>0.8175 \pm 0.0428</u>	0.9364 \pm 0.0303	6.53 \pm 1.70	2.49 \pm 1.15
MS-DdC-AC-DLIR	A2C	0.8229 \pm 0.0469	0.9249 \pm 0.0423	4.75 \pm 2.04	2.61 \pm 1.28
	A4C	0.8159 \pm 0.0398	0.9349 \pm 0.0278	4.56 \pm 1.78	2.47 \pm 1.01
Aug-MS-DdC-AC-DLIR	A2C	0.8200 \pm 0.0467	0.9261 \pm 0.0355	4.10 \pm 1.69	2.53 \pm 0.97
	A4C	0.8095 \pm 0.0390	<u>0.9360 \pm 0.0270</u>	<u>4.12 \pm 1.82</u>	<u>2.37 \pm 0.99</u>

Table B.7: Additional temporal image registration results for adult and fetal echo images, demonstrating the metrics for MYO and LV, where Table 3 shows the average metrics of background, MYO, and LV. All the metrics are estimated using fixed images (and masks) and warped moving images (and masks). Underlined, double-underlined, and bold fonts denote the best-performing metrics for the A2C view of adult echo, the A4C view of adult echo, and the A4C view of fetal echo, respectively.

Methods		DSC (\uparrow)		HD (mm) (\downarrow)	
		MYO	LV	MYO	LV
VanDLIR	Adult (A2C)	0.8590 \pm 0.0829	0.9287 \pm 0.0549	3.72 \pm 2.39	2.28 \pm 1.40
	Adult (A4C)	0.8455 \pm 0.0883	0.9243 \pm 0.0554	4.13 \pm 2.45	2.80 \pm 1.73
	Fetal (A4C)	0.8888 \pm 0.0658	0.9273 \pm 0.0583	1.56 \pm 0.90	1.26 \pm 0.77
AC-DLIR	Adult (A2C)	0.8717 \pm 0.0689	<u>0.9539 \pm 0.0320</u>	3.64 \pm 2.20	<u>1.61 \pm 1.0</u>
	Adult (A4C)	0.8505 \pm 0.0844	0.9446 \pm 0.0441	3.90 \pm 2.29	2.20 \pm 1.38
	Fetal (A4C)	0.9011 \pm 0.0570	0.9319 \pm 0.0512	1.68 \pm 0.93	1.32 \pm 0.77
DdC-AC-DLIR	Adult (A2C)	0.8738 \pm 0.0696	0.9523 \pm 0.0346	3.59 \pm 2.18	1.69 \pm 1.09
	Adult (A4C)	0.8477 \pm 0.0854	0.9445 \pm 0.0453	3.83 \pm 2.23	1.95 \pm 1.33
	Fetal (A4C)	0.9106 \pm 0.0488	0.9411 \pm 0.0383	1.57 \pm 0.92	1.22 \pm 0.72
MS-DdC-AC-DLIR	Adult (A2C)	<u>0.8814 \pm 0.0535</u>	0.9529 \pm 0.0241	<u>3.15 \pm 1.78</u>	1.77 \pm 0.86
	Adult (A4C)	<u>0.8638 \pm 0.0584</u>	<u>0.9507 \pm 0.0307</u>	<u>3.36 \pm 1.94</u>	<u>1.74 \pm 0.90</u>
	Fetal (A4C)	0.9155 \pm 0.0457	0.9463 \pm 0.0332	1.42 \pm 0.84	1.14 \pm 0.68

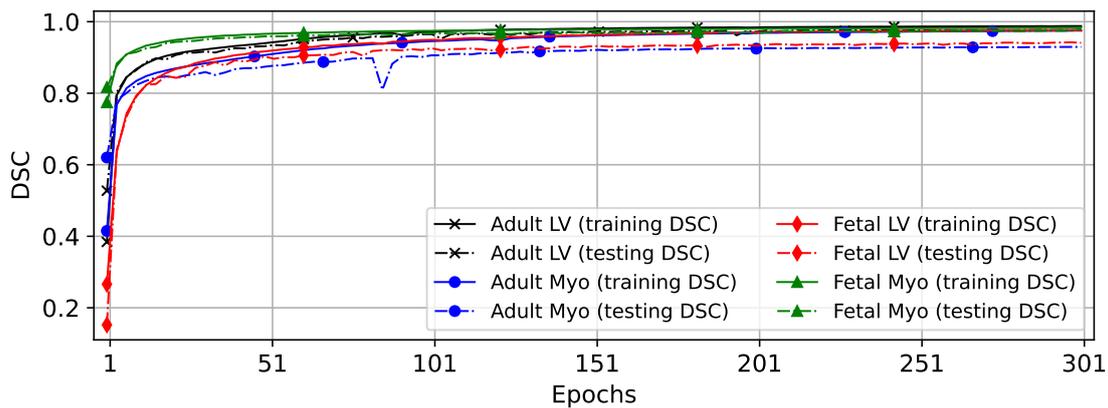


Figure B.13: Training and testing histories of the VAE to learn the latent vectors (\mathcal{Z}). This figure demonstrates that the learned \mathcal{Z} can successfully represent the anatomical MYO and LV topology of the fetal and adult datasets, as the DSCs between the inputted and reconstructed masks are high (see Fig. B.14).

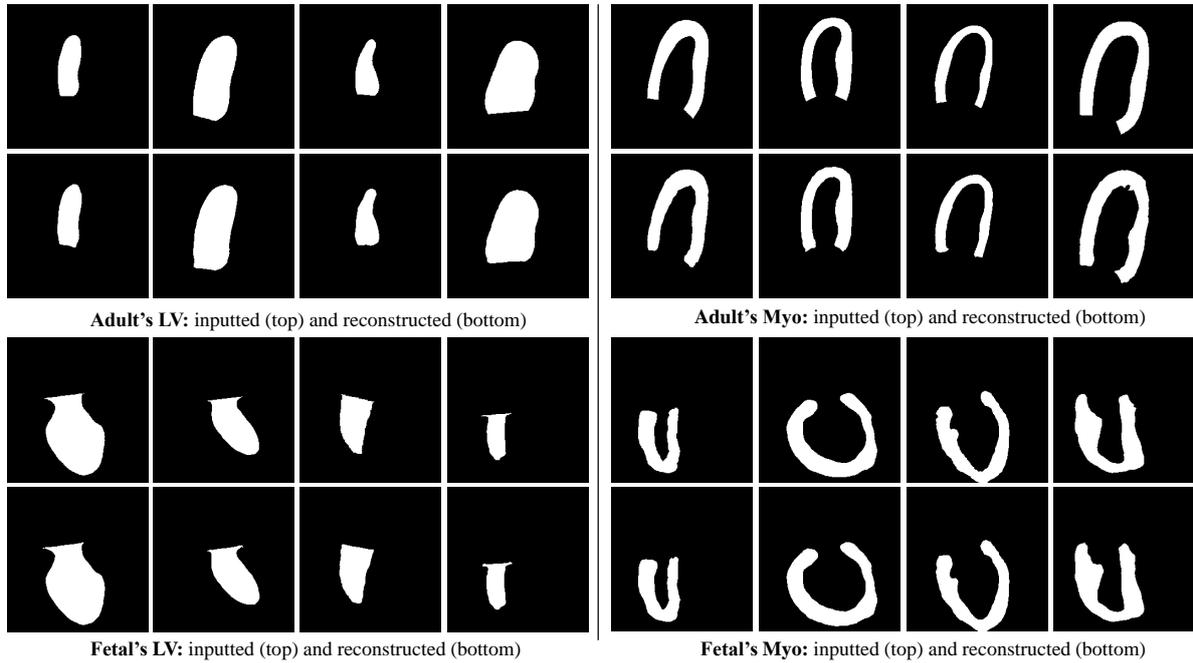


Figure B.14: Qualitative results of the VAE for both adult (two top rows) and fetal (two bottom rows) datasets. The input masks were encoded to the latent vector (\mathcal{Z}), and the reconstructed masks were decoded from the \mathcal{Z} . This figure illustrates that the inputted and reconstructed masks follow a similar anatomical topology, ensuring the latent vector's quality to represent the global attributes of MYO and LV.

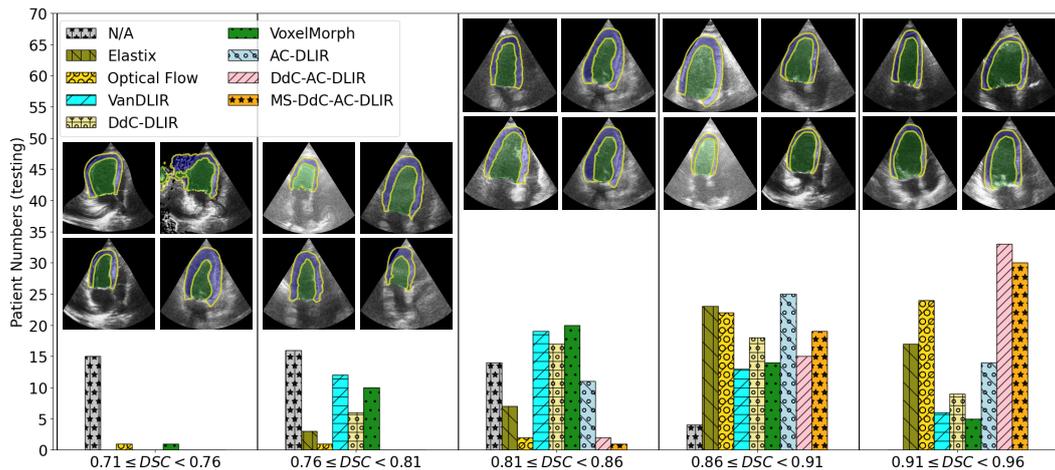


Figure B.15: Five different groups of obtained DSCs, demonstrating the number of testing patients for A4C of CAMUS in each DSC group. Similar results for A2C of CAMUS are in Fig. 9.

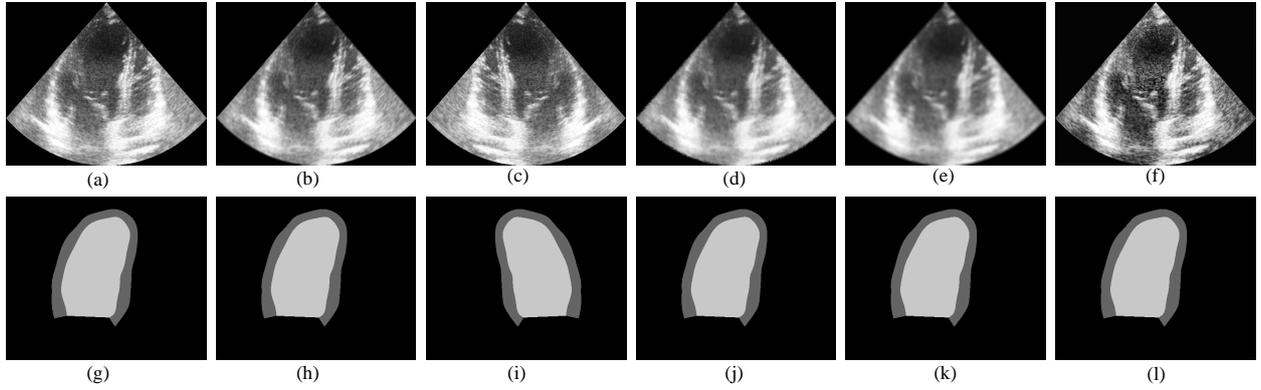


Figure B.16: Examples of image augmentations showing the intensity images with the corresponding masks, where (a,g): original pair, (b,h): motion blurred pair, (c,i): horizontal flipped pair, (d,j): Gaussian blurred pair, (e,k): defocused [67] pair, and (f,l): CLAHE [68] pair.

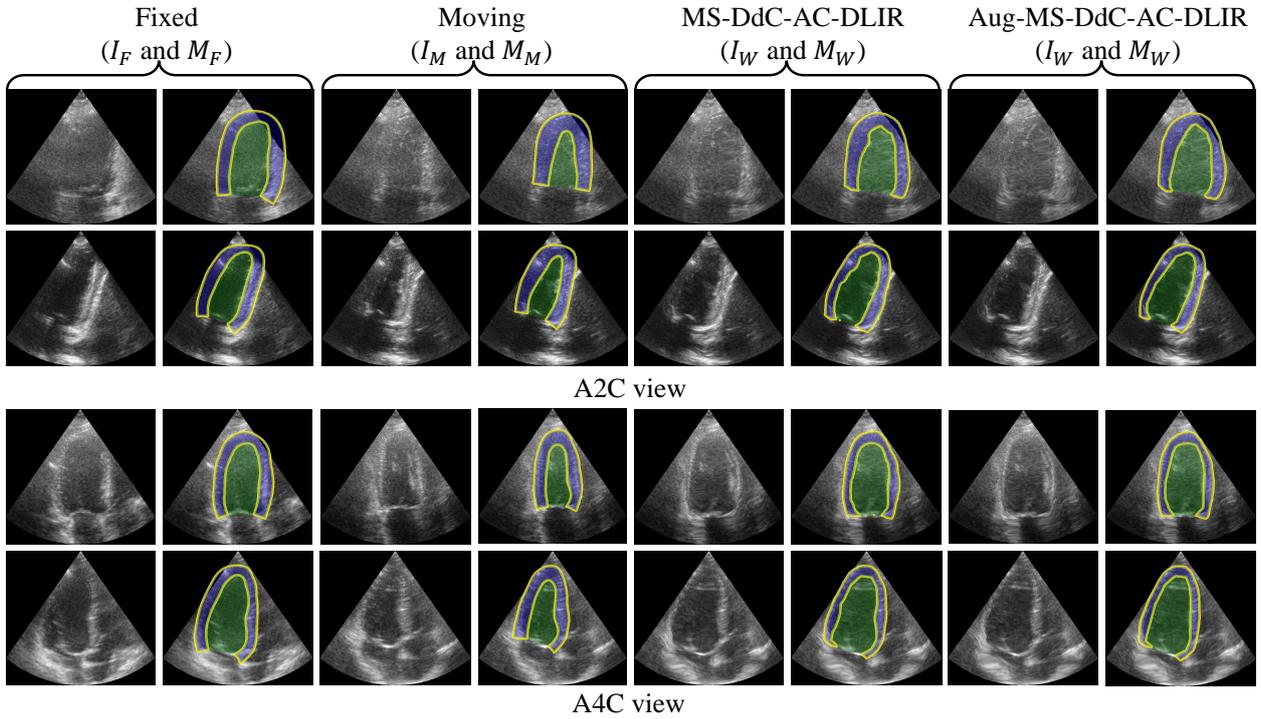
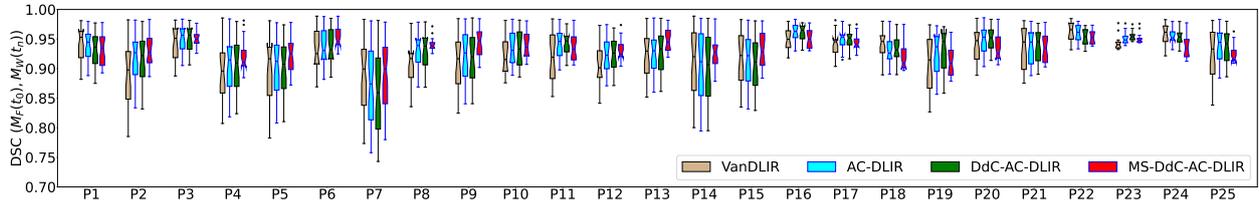
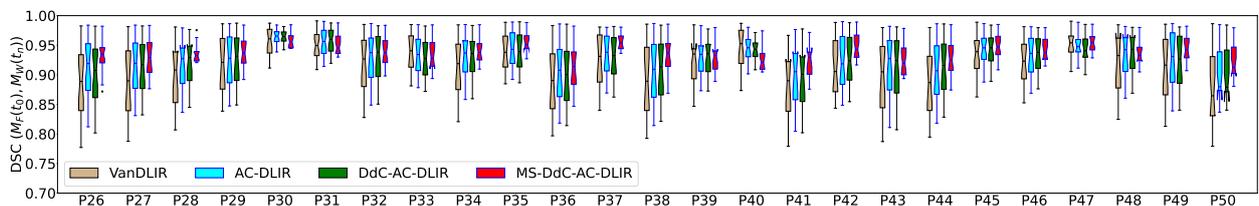


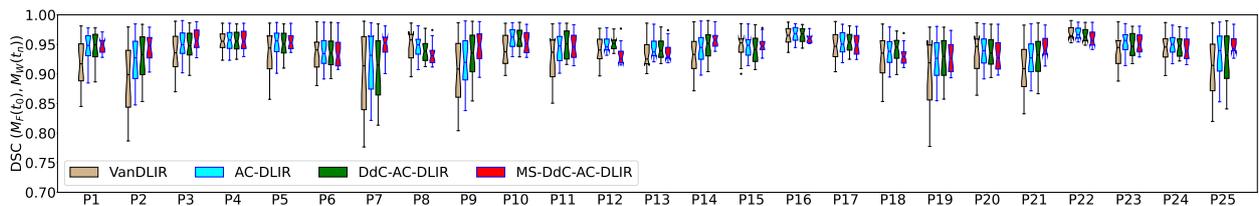
Figure B.17: Sample of qualitative results of the warped image and mask using the proposed MS-DdC-AC-DLIR and Aug-MS-DdC-AC-DLIR for A2C (top two rows) and A4C (bottom two rows) views.



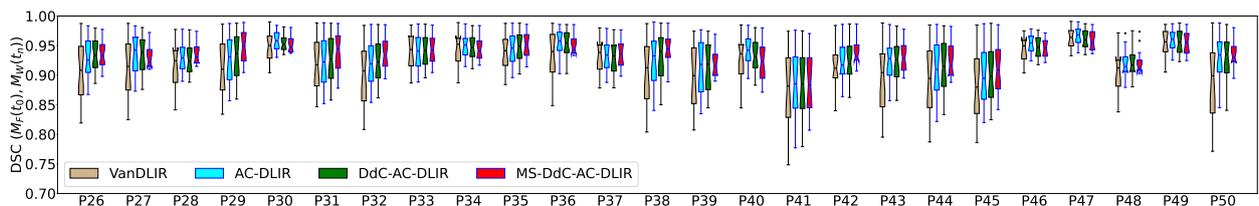
(a) For the first 25 patients (P1-P25) out of 50 testing patients (A4C view)



(b) For the second 25 patients (P26-P50) out of 50 testing patients (A4C view)



(c) For the first 25 patients (P1-P25) out of 50 testing patients (A2C view)



(d) For the second 25 patients (P26-P50) out of 50 testing patients (A2C view)

Figure B.18: Demonstration of the non-overlapping benefits of data-driven and anatomical constraints in temporal echo image registration. The lower interquartile range in the box indicates better temporal consistency in the temporal echo registration. This figure was created using the CAMUS A2C and A4C views, where the first time point is set as a fixed sample, and other time points are warped to this fixed sample.