

# CINEMORPH: LEARNING TIME-CONTINUOUS MOTION FIELD FOR MOTION TRACKING ON CINE MAGNETIC RESONANCE IMAGES

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Tracking cardiac motion using cine magnetic resonance imaging (cine MRI) is essential for evaluating cardiac function and diagnosing cardiovascular diseases. Current methods for cardiac motion tracking depend on scaling and squaring (SS) integration to learn discrete Lagrangian motion fields. However, this reliance hinders the effective exploitation of temporal continuity, leading to inadequate tracking accuracy. In this paper, we introduce a novel unsupervised learning method, CineMorph, to achieve temporally continuous cardiac motion tracking in cine MRI image sequences. Our approach integrates a frame-aware UNet with a series of time-continuous Transformer blocks to learn temporally continuous intra-frame motion fields, which are then assembled into time-continuous Lagrangian motion fields. To ensure the diffeomorphism property, we implement semigroup regularization to constrain our model, thus eliminating the reliance on SS integration. We evaluate our method on the public Automatic Cardiac Diagnostic Challenge (ACDC) dataset. The experimental results show that our method outperforms the existing state-of-the-art methods and achieves state-of-the-art performance with a mean DICE score of 83.6% and a mean Hausdorff distance of 3.4 mm.

## 1 INTRODUCTION

Cine Magnetic Resonance Imaging (cine MRI) plays a crucial role in cardiac motion tracking due to its non-invasive nature and superior imaging capabilities Bello et al. (2019); Reindl et al. (2019); Wang et al. (2023). This technique allows for detailed visualization of the heart’s anatomy and function throughout the cardiac cycle, capturing high-resolution images at multiple phases. By tracking the myocardial motion and deformation, clinicians can accurately assess cardiac function Sliman et al. (2014); Edvardsen et al. (2001), identify abnormalities in heart motion, and evaluate conditions such as myocardial infarction Reed et al. (2017), cardiomyopathies Ciarambino et al. (2021), and valvular diseases Coffey et al. (2021).

Compared to tagged MR images, cine MR images have the advantage of clearly visualizing cardiac anatomy, particularly the myocardium, as the epicardial and endocardial surfaces are distinctly visible. This makes it easier to track the radial motion of the myocardium. However, cine images fall short in accurately quantifying circumferential and longitudinal motion because there are few reliable features within the myocardium to track, and there are often insufficient long-axis images available Shi et al. (2012). Moreover, magnetic field inhomogeneities can cause variations in image brightness, especially with the balanced steady-state free precession (bSSFP) sequence, leading to dark band artifacts Ye et al. (2023).

In recent years, deep learning-based unsupervised methods have emerged as an efficient and effective design scheme for cardiac motion tracking Lu et al. (2023). These methods typically decompose the motion-tracking problem into pairwise registration processes. Using classical pairwise registration networks, such as VoxelMorph Balakrishnan et al. (2019), the motion field can be learned between two consecutive or any two images. When applied to consecutive images, the resulting motion fields need to be composed into Lagrangian motion fields to achieve motion tracking between any two images. The classical work is SequenceMorph Ye et al. (2023), which proposes a bi-directional generative diffeomorphic registration network to estimate the inter-frame motion field

054 between any two consecutive frames, and then recomposed them to the Lagrangian motion field  
055 between the reference frame and any other frame, through a differentiable composition layer. Con-  
056 sidering temporal continuity between consecutive frames, SequenceMorph shows superior tracking  
057 performance and the feasibility of the motion decomposition and recombination principle. Different  
058 from SequenceMorph, Lu *et al.* introduce the temporal relations and automatically learn spatiotem-  
059 poral information from multiple images through a bidirectional recurrent neural network to directly  
060 estimate the Lagrangian motion field between the reference image and other images. However, these  
061 methods rely on the scaling and squaring integration scheme Hernandez et al. (2007); Arsigny et al.  
062 (2006) to reconstruct the deformation field. This reliance imposes a constraint on their capacity to  
063 capture temporal continuity, particularly for large deformation motions.

064 In this paper, we introduce a novel unsupervised learning method, called CineMorph, which gener-  
065 ates time-continuous Lagrangian motion fields to facilitate smoother cardiac motion tracking. Draw-  
066 ing inspiration from Matinkia & Ray (2024), our method leverages the semigroup property Biagi &  
067 Bonfiglioli (2019) to learn the intra-frame motion field at any time and ensure diffeomorphic defor-  
068 mations without using scaling and squaring integration. To achieve this, we propose a new neural  
069 network architecture, which uses a frame-aware UNet Ronneberger et al. (2015) to encode two con-  
070 secutive images with frame information and a series of transformer blocks to obtain time-continuous  
071 intra-frame motion fields. Benefitting from the time-continuous property, we further propose a time-  
072 continuous Lagrangian motion constraint to achieve global temporally-continuous motion tracking,  
073 as shown in Figure 1. To assess the effectiveness of our method, we conduct extensive experiments  
074 on the public ACDC dataset. Our results show our CineMorph is superior to the previous state-of-  
075 the-art models.

076 To sum up, our contributions can be summarized as the following:

- 077 • We introduce a novel unsupervised learning method for tracking cardiac motion in cine  
078 MRI images, which integrates a frame-aware UNet architecture with Transformer blocks  
079 to generate time-continuous Lagrangian motion fields.
- 080 • We propose a time-continuous Lagrangian motion constraint to ensure temporal continuity  
081 and diffeomorphism with semigroup regularization.
- 082 • We provide extensive experiments on the ACDC dataset, which demonstrate the superior  
083 performance of CineMorph over recent state-of-the-art methods.

## 085 2 RELATED WORK

087 **Optical Flow-Based Methods.** Optical flow (OF) is a widely used technique in video sequences  
088 to track objects by estimating the motion of objects between consecutive frames Brox & Malik  
089 (2010); Zhang et al. (2021); Xu et al. (2022); Shi et al. (2023); Saxena et al. (2024). OF can pro-  
090 vide dense motion vectors for every pixel in the image, enabling detailed motion analysis across the  
091 entire frame. OF-based methods estimate cardiac motion field based on several basic assumptions  
092 regarding image appearance and motion strength, such as brightness consistency and small motion  
093 between the fixed and moving frames Carranza-Herrezuelo et al. (2010); Wang et al. (2019). How-  
094 ever, these assumptions are not always valid in cardiac image sequences due to lighting changes,  
095 noise, or large displacements of the myocardium. Another challenge is that most OF-based methods  
096 require supervised learning, which is nearly impractical for medical images.

097 **Image Registration-Based Methods.** Image registration-based methods aim to find a transforma-  
098 tion directly to obtain a dense displacement field that describes motion. Conventional non-rigid  
099 registration approaches, such as parametric B-Splines Rueckert et al. (1999), are formulated as it-  
100 erative optimization procedures that maximize a similarity criterion between the fixed and moving  
101 images to determine the optimal transformation. Shi *et al.* developed a spatial and temporal reg-  
102 istration approach that utilizes free-form deformations to estimate motion within the myocardium  
103 using a spatially-varying, weighted similarity measure Shi et al. (2012). Some studies have also uti-  
104 lized or extended this method to estimate cardiac motion for both untagged and tagged MR images  
105 Chandrashekara et al. (2004); De Craene et al. (2012). However, these methods are often associated  
106 with high computational costs and long execution times.

107 In recent years, there has been a surge of interest in applying deep learning to medical image registra-  
tion and motion tracking Dalca et al. (2019); Niethammer et al. (2019); Chen et al. (2023). Compared

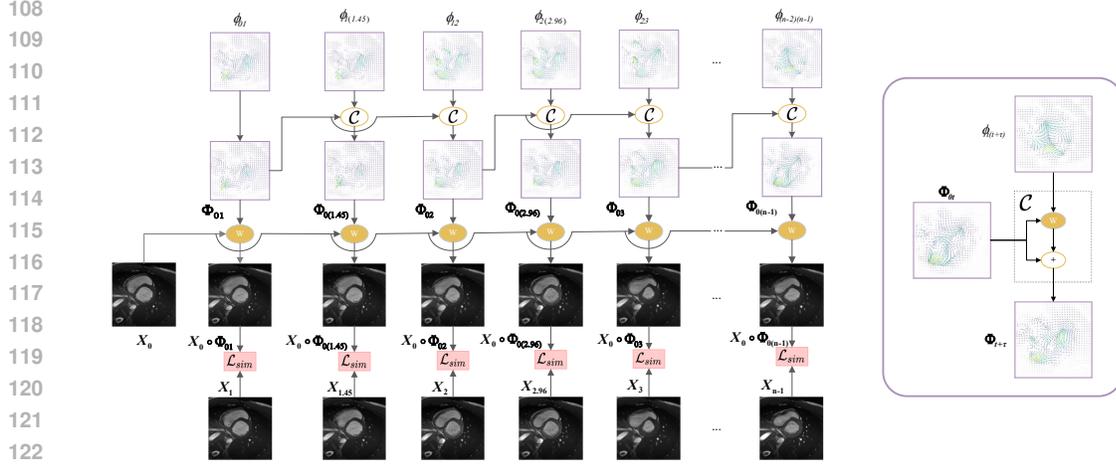


Figure 1: An overview of transforming motion field  $\phi$  to Lagrangian motion field  $\Phi$  with a composition layer  $\mathcal{C}$ . As shown in the figure, our method can achieve temporally continuous motion tracking by estimating time-continuous Lagrangian motion fields. “w” means “warp”.

to traditional iterative methods, deep learning-based approaches are faster and more accurate. In the context of motion tracking, the tracking problem is typically decomposed into pairwise registration processes to directly or indirectly generate Lagrangian motion fields Fechter & Baltas (2020); Yu et al. (2020) using registration networks Balakrishnan et al. (2019); Wu et al. (2022); Joshi & Hong (2023); Wang et al. (2024). Ye *et al.* proposed a bi-directional diffeomorphic registration network to estimate the inter-frame motion fields between consecutive image pairs and recombine them into Lagrangian motion fields through a differentiable composition layer Ye et al. (2021; 2023). Lu *et al.* proposed to model the temporal relations of cardiac cine MRI images through a bidirectional recurrent neural network to obtain the Lagrangian motion field between the reference image and other images Lu et al. (2023).

### 3 METHOD

#### 3.1 PRELIMINARIES

##### 3.1.1 MOTION DECOMPOSITION AND RECOMPOSITION

Cardiac cine MRI images capture a complete cardiac cycle, which comprises two phases: diastole and systole. Typically, the cine sequence starts at the end of diastole (ED), reaches peak contraction at the end of systole (ES), and then relaxes back to the ED phase. For a point  $m$  in a cine image that moves from position  $x_0$  at time  $t_0$ , we need to track its motion trajectory  $x_t$ . In an  $N$ -frame cine MRI image sequence, we only have the finite positions  $x_n$  ( $n = 0, 1, \dots, N - 1$ ) of  $m$ . Over the time interval  $\Delta t = t_{n-1} - t_{n-2}$ , the displacement can be represented as a vector  $\phi_{(n-2)(n-1)}$ , also called inter-frame motion field. A sequence of such inter-frame motions  $\{\phi_{t(t+1)}\}_{t=0}^{n-2}$  is composed to the Lagrangian motion field  $\Phi_{0(n-1)}$  Wang et al. (2019). Based on  $\Phi_{0(n-1)}$ , we can shift the point  $m$  from position  $x_0$  to  $x_{n-1}$ . For motion tracking, given the first frame at time  $t_0$  as the reference frame, our goal is to derive the Lagrangian motion field  $\Phi_{0(n-1)}$  between the reference frame and any subsequent frame at time  $t_{n-1}$ . Direct estimation of the Lagrangian motion field may lead to considerable motion errors due to large heart motion and intensity differences between temporally distant frames during the cardiac cycle. To address this, following Ye et al. (2023), we adopt the motion decomposition and recombination principle, which first estimates the inter-frame motions  $\{\phi_{t(t+1)}\}_{t=0}^{n-2}$  and then recomposes them to the Lagrangian motion field  $\Phi_{0(n-1)}$ .

##### 3.1.2 DFFEOMORPHIC REGISTRATION FOR INTER-FRAME MOTION FIELD

For inter-frame motion field, deformable registration seeks for a vector field  $\phi_{t(t+1)} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ , which warps the moving image  $X_t$  at frame  $t$  smoothly towards the fixed image  $X_{t+1}$  at frame

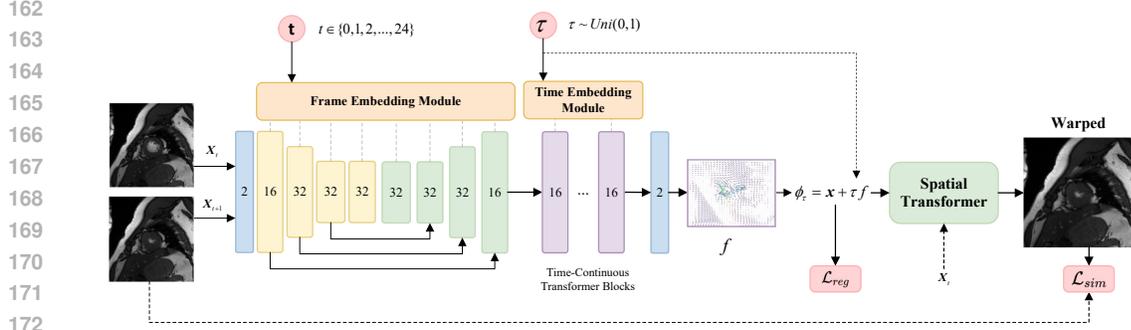


Figure 2: An overview of our proposed network. As illustrated in the figure, our frame-aware UNet is independent of time  $\tau$ . Therefore, when calculating the semigroup loss function, we only need to perform one forward propagation, reducing the training overhead.

$t + 1$ . The deformation field  $\phi$  is generally considered to be the flow map solution of the following ordinary differential equation (ODE) Beg et al. (2005); Chen et al. (2022); Joshi & Hong (2023); Wang et al. (2024):

$$\begin{cases} \frac{d\phi_\tau}{dt} = \mathbf{v}(\phi_\tau) = \mathbf{v} \circ \phi_\tau \\ \phi_0(\mathbf{x}) = \mathbf{x}, \end{cases} \quad (1)$$

where  $\tau \in [0, 1]$ ,  $\mathbf{x}$  is a spatial location,  $\circ$  is a composition operator,  $\mathbf{v}$  is a stationary velocity field and  $\phi_0$  is an identity transformation. The utility of Equation (1) is that its solution is guaranteed to be a diffeomorphism:

$$\phi_{1/2^T} = \mathbf{x} + \frac{\mathbf{v}(\mathbf{x})}{2^T}. \quad (2)$$

$\phi_1$  can be obtained by using the scaling and squaring integration scheme with the recurrence  $\phi_{1/2^i} = \phi_{1/2^{i+1}} \circ \phi_{1/2^{i+1}}$ , which can be expressed as:

$$\phi_{1/2^{T-1}} = \phi_{1/2^T} \circ \phi_{1/2^T} \Rightarrow \dots \Rightarrow \phi_1 = \phi_{1/2} \circ \phi_{1/2}. \quad (3)$$

A necessary and sufficient condition of  $\phi$  as the flow map solution of Equation (1) is that it satisfies the semigroup property, i.e., for any time steps  $\xi$  and  $\zeta$  it holds Biagi & Bonfiglioli (2019)

$$\phi_\xi \circ \phi_\zeta = \phi_\zeta \circ \phi_\xi = \phi_{\xi+\zeta}. \quad (4)$$

Assuming that  $\xi = -\zeta$ , we have  $\phi_\xi \circ \phi_{-\xi} = \phi_0$  to guarantee the bijectivity of the deformation field  $\phi$ . Meanwhile, if the deformation  $\phi$  satisfies Equation (4), then  $\phi$  is a diffeomorphism at any time  $\xi \in [-1, 1]$  Matinkia & Ray (2024).

## 3.2 PROPOSED METHOD

We propose an unsupervised deep learning method, dubbed as CineMorph, to learn a set of time-continuous motion fields  $\{\phi_{t(t+\tau)}\}_{t=0}^{n-2}$ , which are recomposed to the time-continuous Lagrangian motion fields. As shown in Figure 2, CineMorph consists of a frame-aware UNet and multiple time-continuous Transformer blocks. We decouple the frame  $t$  and time  $\tau$ , allowing us to perform only a single forward propagation calculation with UNet when calculating the semigroup loss function. This reduces the computational cost and enhances the flexibility of the overall framework, as UNet can be substituted with other more sophisticated models.

### 3.2.1 FRAME-AWARE UNET

Considering the differences in myocardium motion across different frames, we introduce a frame-aware UNet that better models the motion features of the image pairs using a frame embedding module. The frame-aware UNet takes an image pair and frame  $t$  as input and maps them to a motion feature. Formally, let  $\mathbf{X}_t$  and  $\mathbf{X}_{t+1}$  be a pair of 2D images with the same shape of  $H \times W$  and let  $\mathbf{Z} \in \mathbb{R}^{H \times W \times C}$  be the motion feature encoded by the frame-aware UNet  $\psi$ :

$$\mathbf{Z} = \psi(\mathbf{X}_t, \mathbf{X}_{t+1}, t; \boldsymbol{\theta}_1), \quad (5)$$

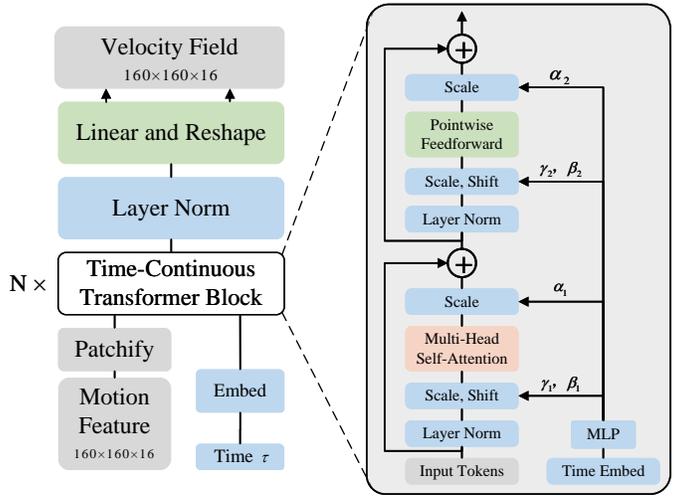


Figure 3: The architecture of time-continuous Transformer blocks. Left: The motion feature is decomposed into patches and processed by several transformer blocks. Right: Details of the time-continuous Transformer block.

where  $\theta_1$  represents the model parameters and  $C$  represents the number of channels. The frame  $t$  is encoded to an embedding vector of dimension  $\mathbb{R}^d$  using a sinusoidal positional embedding  $PE$  Vaswani (2017), followed by a multi-layer perceptron (MLP):

$$\mathbf{W}_2 \sigma(\mathbf{W}_1 PE(t)), \tag{6}$$

where  $\mathbf{W}_1 \in \mathbb{R}^{d \times d}$  and  $\mathbf{W}_2 \in \mathbb{R}^{d \times C}$  are learnable weights and  $\sigma$  is the SiLU activation function. The embedding vector is added to UNet.

### 3.2.2 TIME-CONTINUOUS TRANSFORMER BLOCK

Inspired by Scalable Diffusion Transformers Peebles & Xie (2023), we propose to learn the time-continuous motion field  $\phi$  using time-continuous transformer blocks. As illustrated in Figure 3, the time-continuous transformer block has a similar architecture to other transformer blocks Vaswani (2017). The key difference is that time  $\tau$  is utilized as additional conditional information to regress the scale and shift parameters  $\gamma$  and  $\beta$ , as well as the dimension-wise scaling parameters  $\alpha$ , through an MLP layer. The MLP is initialized to output the zero-vector for all  $\alpha$ , effectively setting the entire transformer block as the identity function. This ensures that the model focuses on learning inter-frame motion fields at the beginning of training. As training progresses, it gradually shifts to learning intra-frame motion fields.

**Patchify and Unpatchify.** The motion feature  $\mathbf{Z}$  has a high spatial resolution ( $160 \times 160$  in our experiment), significantly increasing the computational cost of the transformer blocks. Following Peebles & Xie (2023), we introduce a ‘‘patchify’’ layer as the first layer, which converts the motion feature  $\mathbf{Z}$  into a sequence of tokens, each of dimension  $d$ , using a convolutional layer with kernel size  $k$ . After the final transformer block, we apply a final layer norm and linearly decode each sequence of image tokens. Finally, we rearrange the decoded tokens into their original spatial layout to obtain the predicted velocity field.

Different from Equation (2), we follow Matinkia & Ray (2024) to model the motion field  $\phi$ . Specifically, we construct a sequence of transformer blocks to map the motion feature  $\mathbf{Z}$  to the motion field  $\phi$ :

$$\phi_\tau(\mathbf{x}, \mathbf{Z}; \theta_2) = \mathbf{x} + \tau f(\mathbf{x}, \mathbf{Z}, \tau; \theta_2), \forall \tau \in [-1, 1], \tag{7}$$

where  $f$  is a sequence of transformer blocks with learnable parameters  $\theta_2$ , which receives the motion feature  $\mathbf{Z}$  rather than the pair of images. When  $\tau = 0$ , we have  $\phi_0(\mathbf{x}) = \mathbf{x}$ , hence satisfying the initial condition of the ODE 1. Additionally, to ensure that  $\phi$  is a valid flow map, we enforce the

model to satisfy the semigroup property stated in Equation (4). We achieve this by setting  $\xi = \tau$  and  $\varsigma = \tau - 1$ , which can be expressed as:

$$\phi_\tau \circ \phi_{\tau-1} = \phi_{\tau-1} \circ \phi_\tau = \phi_{2\tau-1}, \forall \tau \in [0, 1]. \quad (8)$$

By randomly sampling  $\tau$ , we can obtain the motion field  $\phi_\tau$  at any time  $\tau$ , thus achieving the prediction of a continuous motion field.

### 3.2.3 INTER-FRAME AND INTRA-FRAME MOTION CONSTRAINTS

According to the bijectivity of the motion field, warping  $\mathbf{X}_t$  up to time  $\tau$  using the motion field  $\phi_\tau$  must be equivalent to warping  $\mathbf{X}_{t+1}$  up to time  $1 - \tau$  using the inverse motion field  $\phi_{\tau-1}$  due to the continuity of the trajectory of  $\phi$ . Hence we can define a time-continuous similarity loss:

$$\mathcal{L}_{sim}(\tau) = MSE(\phi_\tau[\mathbf{X}_t], \phi_{\tau-1}[\mathbf{X}_{t+1}]) = \|\phi_\tau[\mathbf{X}_t], \phi_{\tau-1}[\mathbf{X}_{t+1}]\|_2^2, \quad (9)$$

where  $MSE$  is the mean squared error and  $\phi_\tau[\mathbf{X}_t]$  represents warping  $\mathbf{X}_t$  with  $\phi_\tau$  using a spatial transformer network Jaderberg et al. (2015).  $\mathcal{L}_{sim}$  measures inter-frame motion similarity when  $\tau = 0$  or  $\tau = 1$ , and intra-frame motion similarity when  $0 < \tau < 1$ . The MSE loss is more suitable than normalized local cross-correlation (NCC) for image pairs that have similar intensity distributions and local contrast, such as cardiac Cine-MRI images Joshi & Hong (2023). Hence, we use the MSE loss in our experiments.

Using Equation 8, we impose the semigroup constraint on the motion field to ensure that  $\phi$  is invertible and a diffeomorphism at all time steps:

$$\mathcal{L}_{reg}(\tau) = \|\phi_{2\tau-1} - \phi_\tau \circ \phi_{\tau-1}\|_2 + \|\phi_{2\tau-1} - \phi_{\tau-1} \circ \phi_\tau\|_2, \forall \tau \in [0, 1]. \quad (10)$$

We use an explicit smoothness to the motion field to ensure reasonable deformation by penalizing its gradients:

$$\mathcal{L}_{smooth}(\phi) = \|\nabla\phi\|_2^2. \quad (11)$$

Therefore, the inter-frame and intra-frame motion constraints are:

$$\mathcal{L}_1 = \mathbb{E}_{\tau \sim Uni(0,1)}[\lambda_0 \mathcal{L}_{sim}(\tau) + \lambda_1 \mathcal{L}_{reg}(\tau) + \lambda_2 \mathcal{L}_{smooth}(\phi)], \quad (12)$$

where  $Uni(0, 1)$  is the uniform distribution on  $[0, 1]$ , and  $\lambda_0, \lambda_1$  and  $\lambda_2$  are the regularization factors.

### 3.2.4 TIME-CONTINUOUS LAGRANGIAN MOTION CONSTRAINTS

Benefiting from the prediction of the continuous motion fields  $\{\phi_{t(t+\tau)}\}_{t=0}^{n-2}$ , we can recompose them as time-continuous Lagrangian motion fields  $\{\Phi_{0(t+\tau)}\}_{t=0}^{n-2}$ , with  $\tau \in [0, 1]$ , by a differentiable composition layer  $\mathcal{C}$ , as shown in Figure 1. Formally, we formulate the time-continuous Lagrangian motion fields as:

$$\Phi_{0(t+\tau)} = \phi_{t(t+\tau)} \circ \Phi_{0t}, \forall \tau \in [0, 1], \quad (13)$$

where  $t = 0, 1, \dots, N - 1$ ,  $\Phi_{00} = \phi_{00}$ , and  $\Phi_{01} = \phi_{01}$ .

With the Lagrangian motion field  $\Phi_{0(t+\tau)}$ , we can warp the reference frame image  $\mathbf{X}_0$  to any other time  $t + \tau$ :  $\Phi_{0(t+\tau)}[\mathbf{X}_0]$ . By measuring the similarity between  $\mathbf{X}_{t+\tau}$  and  $\Phi_{0(t+\tau)}[\mathbf{X}_0]$ , we form a time-continuous Lagrangian motion consistency constraint:

$$\mathcal{L}_{lag}(\tau) = \frac{1}{N-1} \sum_{t=0}^{N-2} \mathcal{L}_{sim}(\mathbf{X}_{t+\tau}, \Phi_{0(t+\tau)}[\mathbf{X}_0]), \quad (14)$$

where  $N$  is the total frame number of a cine image sequence.  $\tau$  follows a uniform distribution on  $[0, 1]$ . When  $\tau = 0$  or  $\tau = 1$ , we use the ground truths  $\mathbf{X}_t$  and  $\mathbf{X}_{t+1}$  as labels to compute the loss  $\mathcal{L}_{lag}$ . Otherwise, we use  $\mathbf{X}_{t+\tau} = \phi_\tau[\mathbf{X}_t]$  as a pseudo-label to compute the loss  $\mathcal{L}_{lag}$ . Note that  $\tau$  is independently sampled for each frame. Further, we also enforce the explicit smoothness of the Lagrangian motion field  $\Phi_{0(t+\tau)}$  by penalizing its gradients:

$$\mathcal{L}_{smooth}(\Phi) = \|\nabla\Phi\|_2^2. \quad (15)$$

The Lagrangian motion constraints are:

$$\mathcal{L}_2 = \mathbb{E}_{\tau \sim \text{Uni}(0,1)}[\lambda_3 \mathcal{L}_{lag}(\tau) + \lambda_4 \mathcal{L}_{smooth}(\Phi)], \quad (16)$$

where  $\lambda_3$  and  $\lambda_4$  are the regularization factors to balance the contribution of each loss term. To sum up, the complete loss function  $\mathcal{L}_{total}$  of our method is the sum of  $\mathcal{L}_1$  and  $\mathcal{L}_2$ :

$$\mathcal{L}_{total} = \mathcal{L}_1 + \mathcal{L}_2. \quad (17)$$

## 4 EXPERIMENTS

### 4.1 DATASET AND PRE-PROCESSING

We evaluated our method on the Automatic Cardiac Diagnostic Challenge (ACDC) dataset Bernard et al. (2018). ACDC is a public cine MR dataset that only consists of SAX view cine MR images from 150 subjects. Each scan includes 9 to 10 slices to cover the whole heart. In the original data split, there are 100 subjects in the training set, which includes segmentation mask annotations for the ED and ES frames, and another 50 subjects are in the testing set without any annotation masks. We rearranged and randomized the data based on subgroups, resulting in a revised configuration of 90 cases in the training set, 20 in the validation set, and 40 in the test set. We excluded slices located near the heart’s base or apex due to the absence of annotation masks. The modified data contains 921 two-dimensional sequences in the training set, 180 in the validation set, and 388 in the test set, respectively. For each sequence, the number of frames varies from 12 to 35, covering only the ED to ES phases. If a sequence contains more than 25 frames, we removed extra frames from the sequence, except for the beginning and ending ones. Sequences with fewer than 25 frames remain unchanged. We first extracted the region of interest from the images to cover the heart, then resampled them to the same in-plane spatial size  $160 \times 160$ . Each sequence is used as input to the model for tracking the cyclic cardiac motion. Each input is a 2D sequence with a spatial resolution of  $160 \times 160$  and a maximum of 25 frames. Following Ye et al. (2023), for each 2D image, we normalized the pixel values by first dividing them by 8 times the median intensity value of the image and then truncating the values to the range  $[0, 1]$ . Additionally, we performed data augmentation for each image with random rotation, translation, scaling, and Gaussian noise addition.

### 4.2 EVALUATION METRICS

We evaluated the motion tracking performance using the segmentation masks of the left ventricle (LV), myocardium wall (MYO), right ventricle (RV), and left atrium (LA). Since the mask annotations are available only on the ED and ES frames, we warped the mask from the ED frame to the ES frame using the estimated Lagrangian motion field. Here we used two metrics, the Dice score Dice (1945) and the 95% maximum Hausdorff distance (HD95) Huttenlocher et al. (1993). The Dice score evaluates the degree of overlap between the estimated ES mask and the ground truth ES mask, while the HD95 measures the similarity of the region contours. A higher Dice and lower HD95 scores indicate better overlap between the two segmentation masks, reflecting superior tracking performance.

### 4.3 BASELINE METHODS

We compared our proposed method with three state-of-the-art methods: VoxelMorph (VM) Balakrishnan et al. (2018); Dalca et al. (2019), DeepTag Ye et al. (2021), and SequenceMorph (SM) Ye et al. (2023). For VM and DeepTag, we used their public implementations and retrained them from scratch, following the optimal hyper-parameters suggested by the authors. Since the code has not been released for SM, we report the results directly from their paper. We compare our method with SM without Lagrangian motion refinement (SM woR) for fair comparisons. VM is based on direct Lagrangian motion tracking, whereas DeepTag, SM, and our method are based on Lagrangian motion recomposition.

### 4.4 IMPLEMENTATION DETAILS

Our method was implemented with PyTorch. The architecture of the frame-aware UNet is similar to that described in Matinkia & Ray (2024). Specifically, the encoder has 3 down-sampling layers

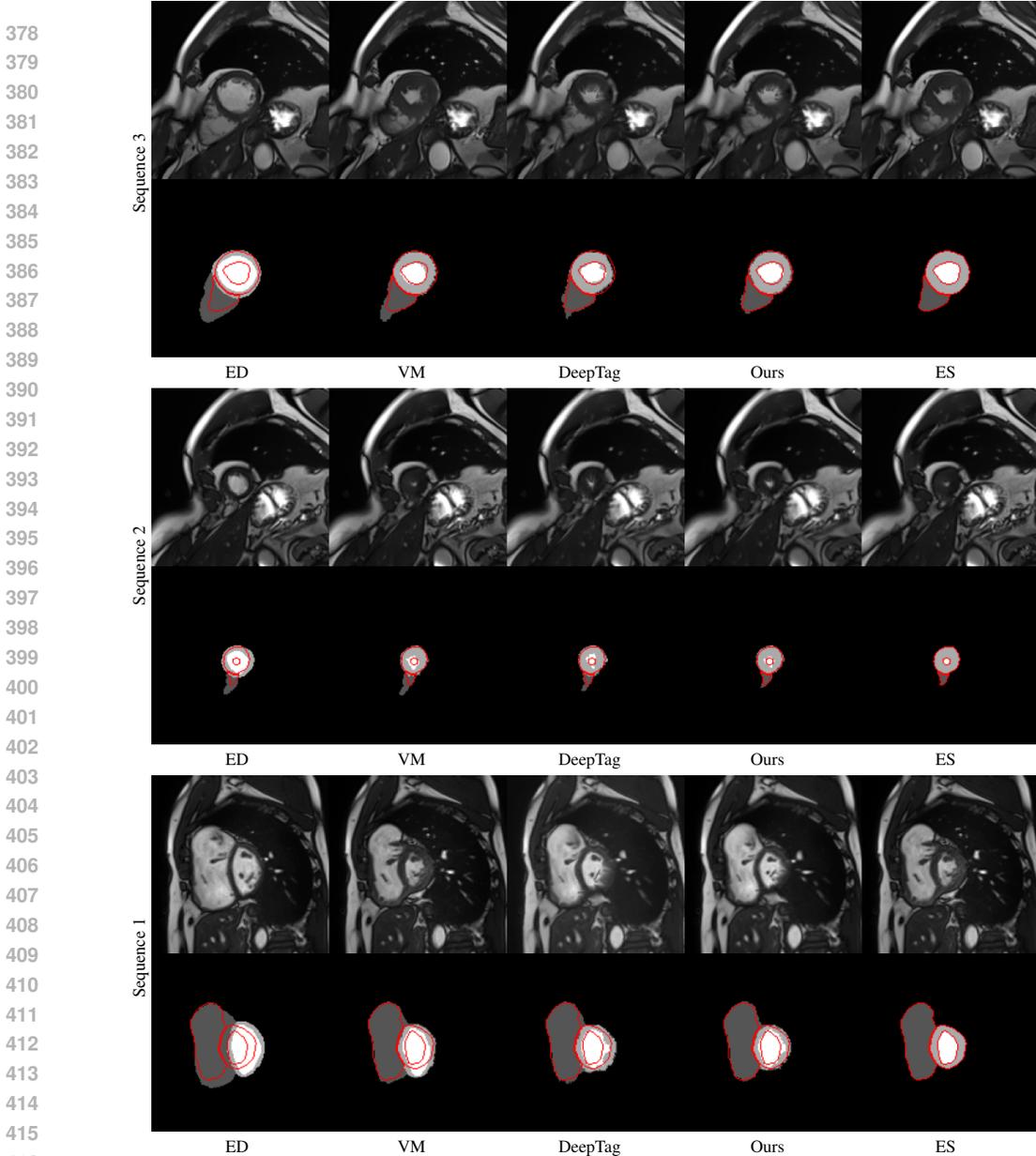


Figure 4: Motion tracking results on three cine MR image sequences (best viewed zoomed in). In each case, first row shows the images and second row shows the segmentation masks. Between ED and ES, we show the warped images by the estimated motion fields of different methods. Red contour shows the ground truth edge of LV, MYO and RV on the ES frame.

of dimensions 32, 32, and 32, and the decoder has 3 up-sampling layers with the same dimensions as the down-sampling layers. After the last up-sampling layer, we use a convolution layer to reduce the dimension to 16. All the activation functions for the layers are set to SiLU Hendrycks & Gimpel (2016) to provide more smoothness to the network. The number of time-continuous transformer blocks is set to 2. The time-embedding dimension is 64. The kernel size of the patchify layer is 8. We use the Adam optimizer with a  $1e^{-4}$  learning rate to train our model for 1000 epochs. The regularization factors are set to  $\lambda_0 = 100$ ,  $\lambda_1 = 5e^8$ ,  $\lambda_2 = 5$ ,  $\lambda_3 = 50$ , and  $\lambda_5 = 1$ , respectively.

#### 4.5 RESULTS

**Motion tracking performance.** Table 1 provides a comprehensive comparison of the motion tracking performance of our method against other baseline methods. All values are expressed as mean

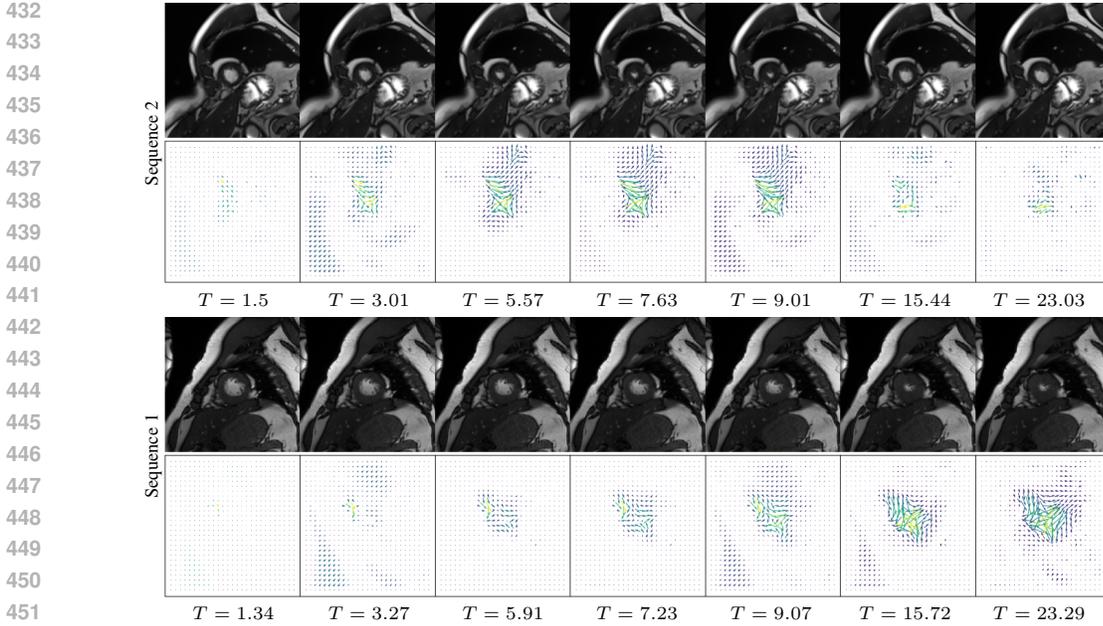


Figure 5: Motion tracking results on two cine MR image sequences (best viewed zoomed in). In each sequence, first row shows the warped images and second row shows the corresponding Lagrangian motion fields at different time  $T = t + \tau$ .

and standard deviation. Our implementation achieves similar motion tracking performance to that of Ye et al. (2023). As shown in Table 1, our method achieves the best performance regarding Dice and HD95 metrics. Compared to VM and DeepTag, our method consistently delivers better results for the LV, MYO, and RV regions. Compared to SM woR, our method shows significant performance improvements, except for the HD95 score in the MYO region. Figure 4 visualizes the warped images and motion tracking results of different methods from the ED phase to the ES phase on cine MR image sequences. The visualization shows that our method aligns more consistently with the ground truth of the ES mask. These results demonstrate the effectiveness of our method.

Table 1: Comparison of the performance of CineMorph with other methods. “woR” denotes “without Lagrangian motion refinement”. “\*” denotes that the results are reported in Ye et al. (2023).

Method	Dice $\uparrow$				HD95(mm) $\downarrow$			
	LV	MYO	RV	avg	LV	MYO	RV	avg
VM*	0.824 $\pm$ 0.156	0.793 $\pm$ 0.105	0.785 $\pm$ 0.175	0.801 $\pm$ 0.021	3.752 $\pm$ 3.607	3.071 $\pm$ 2.399	7.037 $\pm$ 6.679	4.620 $\pm$ 2.121
VM (our impl.)	0.827 $\pm$ 0.170	0.797 $\pm$ 0.110	0.765 $\pm$ 0.208	0.798 $\pm$ 0.166	3.657 $\pm$ 2.508	3.418 $\pm$ 1.822	5.484 $\pm$ 3.731	4.099 $\pm$ 2.867
DeepTag*	0.825 $\pm$ 0.146	0.793 $\pm$ 0.094	0.803 $\pm$ 0.159	0.807 $\pm$ 0.016	3.632 $\pm$ 3.048	2.924 $\pm$ 1.819	6.066 $\pm$ 6.448	4.208 $\pm$ 1.648
DeepTag (our impl.)	0.838 $\pm$ 0.147	0.796 $\pm$ 0.093	0.794 $\pm$ 0.169	0.810 $\pm$ 0.139	3.698 $\pm$ 2.339	3.501 $\pm$ 1.672	4.664 $\pm$ 3.324	3.907 $\pm$ 2.523
SM woR*	0.833 $\pm$ 0.146	0.802 $\pm$ 0.094	0.808 $\pm$ 0.158	0.815 $\pm$ 0.017	3.367 $\pm$ 2.935	2.787 $\pm$ 1.808	5.804 $\pm$ 6.372	4.016 $\pm$ 1.652
Ours	<b>0.860 <math>\pm</math> 0.137</b>	<b>0.826 <math>\pm</math> 0.084</b>	<b>0.821 <math>\pm</math> 0.152</b>	<b>0.836 <math>\pm</math> 0.127</b>	<b>3.073 <math>\pm</math> 2.072</b>	<b>3.050 <math>\pm</math> 1.549</b>	<b>4.081 <math>\pm</math> 3.273</b>	<b>3.356 <math>\pm</math> 2.384</b>

**Visualization of the time-continuous Lagrangian motion field.** Benefiting from the prediction of the time-continuous Lagrangian motion fields, our method, compared to other tracking methods, can predict not only trajectories across frames but also intra-frame trajectories. By estimating the intra-frame motion field, our approach makes the motion field smoother, thereby improving tracking performance. In Figure 5, we visualize the warped images and corresponding Lagrangian motion fields at different time.

#### 4.6 ABLATION STUDY

**Effects of time-continuous transformer blocks.** To investigate the impact of time-continuous transformer blocks on model performance, we train our model with varying numbers of blocks. Considering when the number of the transformer block is 0, the semigroup property is not used to constrain our model. In this case, we change the input of the frame embedding module to the time  $\tau$  sampled from  $Uni(0, 1)$ . The results are reported in Table 2. We find that the transformer block yields considerable performance improvement, indicating the transformer block is critical to

improving motion tracking performance. Again, we observe that across different configurations, similar average Dice and HD95 scores are obtained by increasing the number of blocks, indicating that our method is insensitive to the number of transformer blocks. However, further increasing the number of blocks will increase the computational cost. Therefore, in our experiments, we set the number of blocks to 2 by default.

Table 2: Results of our method with varying numbers of transformer blocks.

Number	Dice $\uparrow$				HD95(mm) $\downarrow$			
	LV	MYO	RV	avg	LV	MYO	RV	avg
0	0.848 $\pm$ 0.145	0.818 $\pm$ 0.092	0.815 $\pm$ 0.158	0.828 $\pm$ 0.134	3.331 $\pm$ 2.130	3.147 $\pm$ 1.574	4.212 $\pm$ 3.284	3.520 $\pm$ 2.409
1	0.859 $\pm$ 0.140	0.828 $\pm$ 0.085	0.817 $\pm$ 0.157	0.836 $\pm$ 0.130	3.060 $\pm$ 2.066	2.987 $\pm$ 1.542	4.146 $\pm$ 3.291	3.348 $\pm$ 2.397
2	0.860 $\pm$ 0.137	0.826 $\pm$ 0.084	0.821 $\pm$ 0.152	0.836 $\pm$ 0.127	3.073 $\pm$ 2.072	3.050 $\pm$ 1.549	4.081 $\pm$ 3.273	3.356 $\pm$ 2.384
3	0.858 $\pm$ 0.141	0.828 $\pm$ 0.083	0.817 $\pm$ 0.159	0.835 $\pm$ 0.131	3.069 $\pm$ 2.127	3.040 $\pm$ 1.562	4.130 $\pm$ 3.304	3.366 $\pm$ 2.421

**Importance of time-continuous Lagrangian motion constraint.** To evaluate the effectiveness of our proposed time-continuous Lagrangian motion constraint (TCLMC), we train our model with and without TCLMC. Table 3 shows the motion tracking performance. Our proposed TCLMC significantly improves the model’s performance in terms of Dice and HD95 scores. TCLMC helps the model learn more continuous motion fields and reduces the drift error accumulating over time, resulting in better motion estimation on a series of frames.

Table 3: Ablation study on the time-continuous Lagrangian motion constraint (TCLMC).

TCLMC	Dice $\uparrow$				HD95(mm) $\downarrow$			
	LV	MYO	RV	avg	LV	MYO	RV	avg
$\times$	0.848 $\pm$ 0.150	0.815 $\pm$ 0.094	0.811 $\pm$ 0.160	0.826 $\pm$ 0.137	3.270 $\pm$ 2.135	3.117 $\pm$ 1.575	4.209 $\pm$ 3.326	3.487 $\pm$ 2.431
$\checkmark$	<b>0.860 <math>\pm</math> 0.137</b>	<b>0.826 <math>\pm</math> 0.084</b>	<b>0.821 <math>\pm</math> 0.152</b>	<b>0.836 <math>\pm</math> 0.127</b>	<b>3.073 <math>\pm</math> 2.072</b>	<b>3.050 <math>\pm</math> 1.549</b>	<b>4.081 <math>\pm</math> 3.273</b>	<b>3.356 <math>\pm</math> 2.384</b>

**Effects of frame embedding module.** Here we study the effects of the frame embedding module with three embedding ways. First, we remove the frame embedding module to explore its importance for motion tracking (Model A). Second, we replace the frame  $t$  with the time  $\tau$  as the input to the frame embedding module (Model B). Third, we maintain the frame embedding module (Model C). The results are shown in Table 4. We find that Model B is superior to Model A, demonstrating the effectiveness of the frame embedding module. Models A and B achieve similar Dice and HD95 scores, indicating that the frame  $t$  and the time  $\tau$  are interchangeable. However, the advantage of using the frame  $t$  as the input is that when implementing the semigroup property, only one forward propagation of the UNet is required, whereas using the time  $\tau$  requires three propagation, significantly reducing computation costs. Additionally, we can use more complex models, such as TransMorph Chen et al. (2022), to train our model for more accurate motion tracking.

Table 4: Ablation study on the frame embedding module. A: without the frame embedding module. B: Replacing the frame  $t$  with the time  $\tau$  as the input to the frame embedding module. C: with the frame embedding module.

Model	Dice $\uparrow$				HD95(mm) $\downarrow$			
	LV	MYO	RV	avg	LV	MYO	RV	avg
A	0.854 $\pm$ 0.143	0.825 $\pm$ 0.091	0.818 $\pm$ 0.157	0.833 $\pm$ 0.132	3.215 $\pm$ 2.165	3.079 $\pm$ 1.584	4.175 $\pm$ 3.323	3.444 $\pm$ 2.443
B	0.859 $\pm$ 0.139	0.827 $\pm$ 0.088	0.821 $\pm$ 0.151	0.836 $\pm$ 0.128	3.123 $\pm$ 2.088	3.064 $\pm$ 1.599	4.089 $\pm$ 3.291	3.381 $\pm$ 2.406
C	0.860 $\pm$ 0.137	0.826 $\pm$ 0.084	0.821 $\pm$ 0.152	0.836 $\pm$ 0.127	3.073 $\pm$ 2.072	3.050 $\pm$ 1.549	4.081 $\pm$ 3.273	3.356 $\pm$ 2.384

## 5 CONCLUSION

In this paper, we present a novel unsupervised learning method for generating time-continuous Lagrangian motion fields to improve cardiac motion tracking in cine MRI images. Our approach utilizes a frame-aware UNet to encode two consecutive images with frame information and employs a series of transformer blocks to derive time-continuous intra-frame motion fields. We train our model using semigroup regularization and time-continuous Lagrangian motion regularization to capture temporal continuity and ensure diffeomorphism. Extensive experiments on the public ACDC dataset demonstrate the effectiveness of our method.

## REFERENCES

- 540  
541  
542 Vincent Arsigny, Olivier Commowick, Xavier Pennec, and Nicholas Ayache. A log-euclidean frame-  
543 work for statistics on diffeomorphisms. In *Medical Image Computing and Computer-Assisted*  
544 *Intervention, October 1-6, 2006. Proceedings, Part 1 9*, pp. 924–931. Springer, 2006.
- 545 Guha Balakrishnan, Amy Zhao, Mert R Sabuncu, John Guttag, and Adrian V Dalca. An unsu-  
546 pervised learning model for deformable medical image registration. In *Proceedings of the IEEE*  
547 *Conference on Computer Vision and Pattern Recognition*, pp. 9252–9260, 2018.
- 548 Guha Balakrishnan, Amy Zhao, Mert R Sabuncu, John Guttag, and Adrian V Dalca. Voxelmorph:  
549 a learning framework for deformable medical image registration. *IEEE Transactions on Medical*  
550 *Imaging*, 38(8):1788–1800, 2019.
- 551 M Faisal Beg, Michael I Miller, Alain Trouvé, and Laurent Younes. Computing large deformation  
552 metric mappings via geodesic flows of diffeomorphisms. *International Journal of Computer*  
553 *Vision*, 61:139–157, 2005.
- 554 Ghalib A Bello, Timothy JW Dawes, Jinming Duan, Carlo Biffi, Antonio De Marvao, Luke SGE  
555 Howard, J Simon R Gibbs, Martin R Wilkins, Stuart A Cook, Daniel Rueckert, et al. Deep-  
556 learning cardiac motion analysis for human survival prediction. *Nature Machine Intelligence*, 1  
557 (2):95–104, 2019.
- 558 Olivier Bernard, Alain Lalonde, Clement Zotti, Frederick Cervenansky, Xin Yang, Pheng-Ann Heng,  
559 Irem Cetin, Karim Lekadir, Oscar Camara, Miguel Angel Gonzalez Ballester, et al. Deep learning  
560 techniques for automatic mri cardiac multi-structures segmentation and diagnosis: is the problem  
561 solved? *IEEE Transactions on Medical Imaging*, 37(11):2514–2525, 2018.
- 562 Stefano Biagi and Andrea Bonfiglioli. *An Introduction to the Geometrical Analysis of Vector Fields:*  
563 *with Applications to Maximum Principles and Lie Groups*. World Scientific, 2019.
- 564 Thomas Brox and Jitendra Malik. Large displacement optical flow: descriptor matching in varia-  
565 tional motion estimation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33  
566 (3):500–513, 2010.
- 567 Noemi Carranza-Herrezuelo, Ana Bajo, Filip Sroubek, Cristina Santamarta, Gabriel Cristóbal,  
568 Andrés Santos, and María J Ledesma-Carbayo. Motion estimation of tagged cardiac magnetic  
569 resonance images using variational techniques. *Computerized Medical Imaging and Graphics*, 34  
570 (6):514–522, 2010.
- 571 Raghavendra Chandrashekar, Raad H Mohiaddin, and Daniel Rueckert. Analysis of 3-d myocardial  
572 motion in tagged mr images using nonrigid image registration. *IEEE Transactions on Medical*  
573 *Imaging*, 23(10):1245–1250, 2004.
- 574 Junyu Chen, Eric C Frey, Yufan He, William P Segars, Ye Li, and Yong Du. Transmorph: Trans-  
575 former for unsupervised medical image registration. *Medical Image Analysis*, 82:102615, 2022.
- 576 Zeyuan Chen, Yuanjie Zheng, and James C Gee. Transmatch: A transformer-based multilevel dual-  
577 stream feature matching network for unsupervised deformable image registration. *IEEE Transac-*  
578 *tions on Medical Imaging*, 43(1):15–27, 2023.
- 579 Tiziana Ciarambino, Giovanni Menna, Gennaro Sansone, and Mauro Giordano. Cardiomyopathies:  
580 an overview. *International Journal of Molecular Sciences*, 22(14):7722, 2021.
- 581 Sean Coffey, Ross Roberts-Thomson, Alex Brown, Jonathan Carapetis, Mao Chen, Maurice  
582 Enriquez-Sarano, Liesl Zühlke, and Bernard D Prendergast. Global epidemiology of valvular  
583 heart disease. *Nature Reviews Cardiology*, 18(12):853–864, 2021.
- 584 Adrian V Dalca, Guha Balakrishnan, John Guttag, and Mert R Sabuncu. Unsupervised learning  
585 of probabilistic diffeomorphic registration for images and surfaces. *Medical Image Analysis*, 57:  
586 226–236, 2019.

- 594 Mathieu De Craene, Gemma Piella, Oscar Camara, Nicolas Duchateau, Etelvino Silva, Adelina  
595 Doltra, Jan D’hooge, Josep Brugada, Marta Sitges, and Alejandro F Frangi. Temporal diffeomor-  
596 phic free-form deformation: Application to motion and strain estimation from 3d echocardiogra-  
597 phy. *Medical Image Analysis*, 16(2):427–450, 2012.
- 598  
599 Lee R Dice. Measures of the amount of ecologic association between species. *Ecology*, 26(3):  
600 297–302, 1945.
- 601 Thor Edvardsen, Helge Skulstad, Svend Aakhus, Stig Urheim, and Halfdan Ihlen. Regional my-  
602 ocardial systolic function during acute myocardial ischemia assessed by strain doppler echocar-  
603 diography. *Journal of the American College of Cardiology*, 37(3):726–730, 2001.
- 604 Tobias Fechter and Dimos Baltas. One-shot learning for deformable medical image registration and  
605 periodic motion tracking. *IEEE Transactions on Medical Imaging*, 39(7):2506–2517, 2020.
- 606  
607 Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint*  
608 *arXiv:1606.08415*, 2016.
- 609  
610 Monica Hernandez, Matias N Bossa, and Salvador Olmos. Registration of anatomical images using  
611 geodesic paths of diffeomorphisms parameterized with stationary vector fields. In *2007 IEEE*  
612 *11th International Conference on Computer Vision*, pp. 1–8. IEEE, 2007.
- 613  
614 Daniel P Huttenlocher, Gregory A. Klanderman, and William J Rucklidge. Comparing images using  
615 the hausdorff distance. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 15(9):  
616 850–863, 1993.
- 617 Max Jaderberg, Karen Simonyan, Andrew Zisserman, et al. Spatial transformer networks. *Advances*  
618 *in Neural Information Processing Systems*, 28, 2015.
- 619  
620 Ankita Joshi and Yi Hong. R2net: Efficient and flexible diffeomorphic image registration using  
621 lipschitz continuous residual networks. *Medical Image Analysis*, 89:102917, 2023.
- 622  
623 Jiayi Lu, Renchao Jin, Manyang Wang, Enmin Song, and Guangzhi Ma. A bidirectional registration  
624 neural network for cardiac motion tracking using cine mri images. *Computers in Biology and*  
625 *Medicine*, 160:107001, 2023.
- 626  
627 Mohammadjavad Matinkia and Nilanjan Ray. Learning diffeomorphism for image registration with  
628 time-continuous networks using semigroup regularization. *arXiv preprint arXiv:2405.18684*,  
2024.
- 629  
630 Marc Niethammer, Roland Kwitt, and Francois-Xavier Vialard. Metric learning for image registra-  
631 tion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
632 pp. 8463–8472, 2019.
- 633  
634 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of*  
*the IEEE/CVF International Conference on Computer Vision*, pp. 4195–4205, 2023.
- 635  
636 Grant W Reed, Jeffrey E Rossi, and Christopher P Cannon. Acute myocardial infarction. *The*  
637 *Lancet*, 389(10065):197–210, 2017.
- 638  
639 Martin Reindl, Christina Tiller, Magdalena Holzknrecht, Ivan Lechner, Alexander Beck, David Plap-  
640 pert, Michelle Gorzala, Mathias Pamminger, Agnes Mayr, Gert Klug, et al. Prognostic implica-  
641 tions of global longitudinal strain by feature-tracking cardiac magnetic resonance in st-elevation  
myocardial infarction. *Circulation: Cardiovascular Imaging*, 12(11):e009404, 2019.
- 642  
643 Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomed-  
644 ical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention*,  
*Munich, Germany, October 5-9, 2015, proceedings, part III 18*, pp. 234–241. Springer, 2015.
- 645  
646 Daniel Rueckert, Luke I Sonoda, Carmel Hayes, Derek LG Hill, Martin O Leach, and David J  
647 Hawkes. Nonrigid registration using free-form deformations: application to breast mr images.  
*IEEE Transactions on Medical Imaging*, 18(8):712–721, 1999.

- 648 Saurabh Saxena, Charles Herrmann, Junhwa Hur, Abhishek Kar, Mohammad Norouzi, Deqing Sun,  
649 and David J Fleet. The surprising effectiveness of diffusion models for optical flow and monocular  
650 depth estimation. *Advances in Neural Information Processing Systems*, 36, 2024.
- 651 Wenzhe Shi, Xiahai Zhuang, Haiyan Wang, Simon Duckett, Duy VN Luong, Catalina Tobon-  
652 Gomez, KaiPin Tung, Philip J Edwards, Kawal S Rhode, Reza S Razavi, et al. A comprehensive  
653 cardiac motion estimation framework using both untagged and 3-d tagged mr images based on  
654 nonrigid registration. *IEEE Transactions on Medical Imaging*, 31(6):1263–1275, 2012.
- 655 Xiaoyu Shi, Zhaoyang Huang, Weikang Bian, Dasong Li, Manyuan Zhang, Ka Chun Cheung, Simon  
656 See, Hongwei Qin, Jifeng Dai, and Hongsheng Li. Videoflow: Exploiting temporal cues for  
657 multi-frame optical flow estimation. In *Proceedings of the IEEE/CVF International Conference  
658 on Computer Vision*, pp. 12469–12480, 2023.
- 660 Hisham Sliman, Ahmed Elnakib, G Beache, Adel Elmaghraby, and Ayman El-Baz. Assessment  
661 of myocardial function from cine cardiac mri using a novel 4d tracking approach. *Journal of  
662 Computer Science and Systems Biology*, 7:169–73, 2014.
- 663 A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.
- 664 Haiqiao Wang, Dong Ni, and Yi Wang. Recursive deformable pyramid network for unsupervised  
665 medical image registration. *IEEE Transactions on Medical Imaging*, 2024.
- 666 Liang Wang, Patrick Clarysse, Zhengjun Liu, Bin Gao, Wanyu Liu, Pierre Croisille, and Philippe  
667 Delachartre. A gradient-based optical-flow cardiac motion estimation method for cine and tagged  
668 mr images. *Medical Image Analysis*, 57:136–148, 2019.
- 671 Yu Wang, Changyu Sun, Sona Ghadimi, Daniel C Auger, Pierre Croisille, Magalie Viallon, Ken-  
672 neth Mangion, Colin Berry, Christopher M Haggerty, Linyuan Jing, et al. Strainnet: improved  
673 myocardial strain analysis of cine mri by deep learning from dense. *Radiology: Cardiothoracic  
674 Imaging*, 5(3):e220196, 2023.
- 675 Yifan Wu, Tom Z Jiahao, Jiancong Wang, Paul A Yushkevich, M Ani Hsieh, and James C Gee.  
676 Nodeo: A neural ordinary differential equation based optimization framework for deformable  
677 image registration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
678 Recognition*, pp. 20804–20813, 2022.
- 680 Haofei Xu, Jing Zhang, Jianfei Cai, Hamid Rezatofghi, and Dacheng Tao. Gmflow: Learning  
681 optical flow via global matching. In *Proceedings of the IEEE/CVF Conference on Computer  
682 Vision and Pattern Recognition*, pp. 8121–8130, 2022.
- 683 Meng Ye, Mikael Kanski, Dong Yang, Qi Chang, Zhennan Yan, Qiaoying Huang, Leon Axel, and  
684 Dimitris Metaxas. Deeptag: An unsupervised deep learning method for motion tracking on car-  
685 diac tagging magnetic resonance images. In *Proceedings of the IEEE/CVF Conference on Com-  
686 puter Vision and Pattern Recognition*, pp. 7261–7271, 2021.
- 687 Meng Ye, Dong Yang, Qiaoying Huang, Mikael Kanski, Leon Axel, and Dimitris N Metaxas. Se-  
688 quencemorph: A unified unsupervised learning framework for motion tracking on cardiac im-  
689 age sequences. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(8):10409–  
690 10426, 2023.
- 692 Hanchao Yu, Shanhu Sun, Haichao Yu, Xiao Chen, Honghui Shi, Thomas S Huang, and Terrence  
693 Chen. Foal: Fast online adaptive learning for cardiac motion estimation. In *Proceedings of the  
694 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4313–4323, 2020.
- 695 Feihu Zhang, Oliver J Woodford, Victor Adrian Prisacariu, and Philip HS Torr. Separable flow:  
696 Learning motion cost volumes for optical flow estimation. In *Proceedings of the IEEE/CVF  
697 International Conference on Computer Vision*, pp. 10807–10817, 2021.
- 698  
699  
700  
701