# `NePhi`: Neural Deformation Fields for Approximately Diffeomorphic Medical Image Registration

**Lin Tian**
UNC Chapel Hill

**Soumyadip Sengupta**
UNC Chapel Hill

**Hastings Greer**
UNC Chapel Hill

**Raúl San José Estépar**
Harvard Medical School

**Marc Niethammer**
UNC Chapel Hill

## Abstract

This work proposes `NePhi`, a neural deformation model which results in approximately diffeomorphic transformations. In contrast to the predominant voxel-based approaches, `NePhi` represents deformations functionally which allows for memory-efficient training and inference. This is of particular importance for large volumetric registrations. Further, while medical image registration approaches representing transformation maps via multi-layer perceptrons have been proposed, `NePhi` facilitates both pairwise optimization-based registration *as well as* learning-based registration via predicted or optimized global and local latent codes. Lastly, as deformation regularity is a highly desirable property for most medical image registration tasks, `NePhi` makes use of gradient inverse consistency regularization which empirically results in approximately diffeomorphic transformations. We show the performance of `NePhi` on two 2D synthetic datasets as well as on real 3D lung registration. Our results show that `NePhi` can achieve similar accuracies as voxel-based representations in a single-resolution registration setting, while using less memory and allowing for faster instance-optimization.

## 1 Introduction

Given a moving image $I^A$ and a fixed image $I^B$, the goal of image registration is to find a spatial transformation, $\Phi$, which spatially transforms $I^A$ so that it corresponds to $I^B$, i.e., $I^A \circ \Phi^{-1} \approx I^B$. A mathematical formulation of registration usually involves the specification of a transformation model (i.e., how $\Phi$ is expressed, for example, parametrically or non-parametrically) and a similarity measure to assess if a spatially transformed moving image $I^A \circ \Phi^{-1}$ should be considered a good registration with respect to $I^B$. Determining $\Phi^{-1}$ requires solving an inverse problem of the form

$$\theta^* = \underset{\theta}{\mathrm{argmin}}\ \mathrm{Sim}(I^A \circ \Phi_\theta^{-1}, I^B) + \lambda \mathrm{Reg}(\Phi_\theta^{-1}), \quad \lambda \geq 0, \tag{1}$$

where $\mathrm{Sim}(\cdot, \cdot)$ denotes the similarity measure (for example mean squared error, normalized cross correlation, or mutual information), $\mathrm{Reg}(\cdot)$ denotes a regularizer, and $\theta$ are the parameters describing the transformation $\Phi_\theta^{-1}$ (these can for example be the parameters of an affine transformation, vector fields for a non-parametric specification of $\Phi_\theta^{-1}$, or coefficients of a multi-layer perceptron (MLP)).

Classic registration algorithms simply solve for $\theta$ by numerical optimization for a pair of images. More recently, deep neural networks (DNN) have been used to predict these registration parameters from $I^A$ and $I^B$. I.e., a DNN takes the high-dimensional images as inputs and outputs the transformation parameters: $\theta = \mathrm{DNN}(I^A, I^B)$. This works well in general. However, it also poses a formidable

challenge with respect to memory consumption. In both the classical and deep learning registration methods, the transformation map $\Phi$ is represented on a voxel grid. Hence, memory consumption scales with the number of voxels which can quickly become prohibitive for 3D convolutional DNNs for large 3D image volumes and with large network feature dimensions. While more memory efficient approaches have been proposed [9] these in general make the network design more complex.

Hence, it is desirable to explore alternative parameterizations of the transformations to reduce memory requirements and to open up the possibility to use more structurally agnostic network architectures, specifically, fully-connected networks. Intuitively, such alternative parameterizations should be possible, because deformation spaces are expected to be much lower-dimensional than the dimensionality of typical deformation parameterizations, e.g., when estimating a displacement field for a $100^3$ 3D image one estimates a discretized displacement field with 3 million parameters. Recently, there has being an increasing interest in neural representations, namely the implicitly defined, continuous, differentiable signal representations parameterized by neural networks [35]. These representations have been used in various computer vision and computer graphics tasks, where is was shown that one neural network has the capacity to embed the information contained in one 3D scene and that a 3D scene can be represented by the weights of the neural network. In such an approach, the memory use is fixed after training, even if the resolution of the 3D scene changes. Our goal in this work is explore such implicit representations for image registration.

The contributions of our work are as follows:

- We propose a neural deformation field, `NePhi`, which is based on an implicit representation and which empirically results in approximately diffeomorphic transformations (i.e., transformations that are smooth, have a smooth inverse, and are bijective) by using gradient inverse consistency regularization. Such diffeomorphic transformations are highly desirable for many medical image registration tasks to be able to move between image spaces.

- In contrast to existing neural deformations for image registration, `NePhi` is not only suitable for optimization-based registration but generalizes via predicted local and global latent codes in a learning-based registration framework. Inspired by work on deformable shape representations we use a CNN encoder to predict these latent codes thereby obtaining the first generalizable, learning-based registration approach.

- We demonstrate the performance of `NePhi` on synthetic data as well as on a real 3D lung registration task between inhale and exhale computed tomography (CT) images showing competitive registration performance.

## 2 Related work

### 2.1 Medical Image Registration

The typical goal of medical image registration is to find the spatial transformation between a pair of images. Traditionally, this task is formulated as an optimization problem [1, 37, 34, 14, 15], aiming at identifying the optimal transformation parameters that minimize the dissimilarity between the a warped moving image and a fixed image while adhering to certain regularity constraints. As there is generally no closed-form solution to these optimization problems they are solved using iterative optimization algorithms (typically a form of gradient descent), resulting in significant computation time. To improve capture range multi-resolution approaches have also been proposed [2, 21]. To reduce registration run-time, learning-based methods [47, 33, 6, 3, 24, 23, 20, 11, 40] have been proposed. The pioneering works [47, 6] demonstrated how convolutional neural networks (CNN) can be used for learning-based medical image registration. Subsequent work has added multi-step and multi-resolution capabilities to learning-based registration networks thereby achieving competitive results with optimization-based multi-resolution registration methods but with significant less run-time. Instance-optimization, where the registration networks are finetuned for a given image pair, have further increased registration performance [40]. Despite the state-of-the-art performance and nearly instantaneous estimation during inference provided by the learning-based multi-resolution registration methods, they consume significant memory during training, because the used CNNs predict dense deformation fields. Further, these approaches are relatively slow when using instance optimization as large neural networks are finetuned. Running registrations at high spatial resolutions for 3D image volumes is therefore not easily possible with these learning-based methods due to

memory constraints. *Our motivation is therefore to develop a neural deformation model which has a much more flexible memory profile (as a tradeoff between runtime and memory consumption) and allows for faster instance optimization.*

## 2.2 Neural Deformation Models

Recently, there has been increasing interest in using functional representations of deformation fields (parameterized via DNNs) for natural images in dynamic scenes [8, 18, 18, 28, 30, 41, 27, 17], for dynamic objects [16, 38, 43, 7, 25], and for animatable humans [46, 48, 49, 5, 10, 19, 29, 32]. The goal of these approaches is to reconstruct the underlying 3D scene, object, or human with respect to the motion contained in the given images. There is so far limited work on using such functional representations for medical image registration [12, 44, 36, 50, 45]. Existing approaches either use DNN-parameterized function to directly represent a displacement vector field (DVF) [50, 44] or to represent velocity or momentum fields [12, 36, 45] to capture large deformations. Velocity or momentum field approaches indirectly parameterize a transformation map which is recovered from the velocity or momentum fields by numerical integration. Such an indirect parameterization can assure diffemorphic transformations. but requires costly numerical integrations. *Most importantly, all existing registration work has only focused on* optimization-based *registration, which is slow compared to learning-based registration which predicts transformations at test time. In contrast, our proposed* NePhi *approach is suitable for optimization-based and learning-based registration.*

## 3 Background

Given two images $I^A : \Omega^A \to \mathbb{R}$ and $I^B : \Omega^B \to \mathbb{R}$, where $\Omega^{A,B} \in \mathbb{R}^d$ indicate the domain of the images, our goal is to find a mapping $\varphi^{AB} : \Omega^A \to \Omega^B$ such that $I^A \circ \varphi^{AB} \approx I^{B1}$. This continuous transformation map, $\varphi^{AB}$, can be parameterized as a function with parameters $\theta$, denoted as $\Phi_\theta^{AB}$, where the parameters depend on the function class to be parameterized. E.g., these could be parameters for an affine transformations, a discretized vector-field on a voxel grid (from which the map itself can be recovered by interpolation), or the parameters of a neural network. As discussed in Sec. 2.1, there are two general types of approaches for image registration.

**Optimization-based registration.** Given a pair of images $I^A$ and $I^B$, one can solve for the transformation map $\Phi_\theta^{AB}$ by optimizing the parameters $\theta$ as follows:

$$\theta^* = \arg\min_\theta \; \mathcal{L}_{\text{sim}} \left( I^A \circ \Phi_\theta^{AB}, I^B \right) + \lambda \mathcal{L}_{\text{reg}}(\Phi_\theta^{AB}), \tag{2}$$

where $\mathcal{L}_{\text{sim}}(\cdot, \cdot)$ is the *similarity measure*, $\mathcal{L}_{\text{reg}}(\cdot)$ is a *regularizer* with $\lambda \geq 0$.

**Learning-based registration.** Instead of optimizing for a single pair of images $I^A$ and $I^B$, we can train a convolutional neural network (CNN) $f_\theta(I_i^A, I_i^B) = \Phi_i^{AB}$ that outputs a *discretized* $\varphi$. The neural network is trained over a set of image pairs $I = \{(I_i^A, I_i^B)\}_{i=1}^N$ by solving

$$\theta^* = \arg\min_\theta \frac{1}{N} \sum_{i=1}^N \mathcal{L}_{\text{sim}} \left( I_i^A \circ f_\theta(I_i^A, I_i^B), I_i^B \right) + \lambda \mathcal{L}_{\text{reg}}(\Phi_i^{AB}). \tag{3}$$

**Diffeomorphic transformations.** In medical image registration it is frequently desirable to obtain diffeomorphic transformation maps, i.e., transformation maps $\varphi$ that are smooth, bijective and have a smooth inverse $\varphi^{-1}$. Such transformations are desirable because they allow moving between the image spaces without worrying about folds in the transformation, e.g., to go between an image and an atlas space. The standard approach to obtain diffeomorphic transformations is by parameterizing them by velocity fields. A transformation map is then obtained by costly numerical integration. Instead, we will use a regularizer (see Sec. 4) that directly asks for invertibility and empirically result in smooth and hence approximately diffeomorphic transformations.

## 4 Methods

### 4.1 Diffeomorphisms for Neural Deformation Fields

We adopt the approach to obtain *approximately* diffeomorphic neural transformation maps from [40] where gradient inverse consistency is used as a regularizer to push the network to output an approx-

---

[1]A perfect match can in general not be achieved due to image differences or image noise.

imate diffeomorphism. Specifically, given a forward transformation map $\Phi^{AB} = \text{Id} + f_\theta(I^A, I^B)$ and the backward map $\Phi^{BA} = \text{Id} + f_\theta(I^A, I^B)$, where Id denotes the identity map, regularity can be encouraged by minimizing the loss

$$\mathcal{L}_{reg} = \left\| \nabla \left[ \Phi_\theta^{AB} \circ \Phi_\theta^{BA} \right] - \mathbf{I} \right\|_F^2 , \tag{4}$$

where $I$ is the identity matrix and $\nabla$ denotes the Jacobian. To use this loss with the neural representation of NePhi, we use two MLPs, denoted $\varphi_{\theta_1}$ and $\varphi_{\theta_2}$, to represent the continuous forward transformation $\Phi_\theta^{AB}$ and backward transformation $\Phi_\theta^{BA}$, respectively. Note that this is a more flexible design thatn the one from the original GradICON work [40] where only one CNN is used for both directions. To make these MLPs generalizable we condition them using the *same* latent code $z$. With such a design, we can then express the gradient inverse consistency loss (to encourage approximate diffeomorphisms) as

$$\mathcal{L}_{reg}(\theta_1, \theta_2, z) = E_{x \sim q(x)} \left\| \nabla \left[ \varphi_{\theta_1}(\varphi_{\theta_2}(x; z); z) - x \right] \right\|_F^2 , \tag{5}$$

where $x$ is the point coordinate and we use a uniform distribution $q(x)$ across the image domain.

## 4.2 Hybrid Conditioning

To enable the generalizability of NePhi to new transformations (for new image pairs) we use a latent code, $z$, introduced above to condition $\varphi_{\theta_1}$ and $\varphi_{\theta_2}$. Specifically, $z$ is composed of a global latent code *vector* $z_g$ with dimension of $C$ and a local latent code *map* $z_l$ with shape $C \times H \times W$ for a 2D transformation and shape $C \times D \times H \times W$ for a 3D transformation. For any point $x$ we obtain a local latent code via bilinear interpolation of the local latent code from a pixel/voxel representation of latent codes. The interpolated local latent code at point $x$ is denoted as $z_l(x)$. For a given $x$, we then obtain the overall latent code, $z$ of dimension $2C$ via concatenation

$$z(x) = [z_g, z_l(x)], \tag{6}$$

where $[\cdot, \cdot]$ represents the concatenation of two vectors. During the learning of the latent codes in optimization-based or learning-based registration, we adopt the following loss from [26] to encourage a latent codes of low magnitude

$$\mathcal{L}_{prior}(z_g, z_l) = \frac{1}{\sigma^2} ||z_g||_2^2 + \frac{1}{|\Omega_{z_l}|} \sum_{i \in \Omega_{z_l}} \frac{1}{\sigma^2} ||z_l(i)||_2^2 \tag{7}$$

where $\Omega_{z_l}$ represents the set grid points where the $z_l$ are defined. $z_g$ and $z_l$ are randomly initialized from a zero-mean multivariate Gaussian $\mathcal{N}(0, \sigma^2)$.
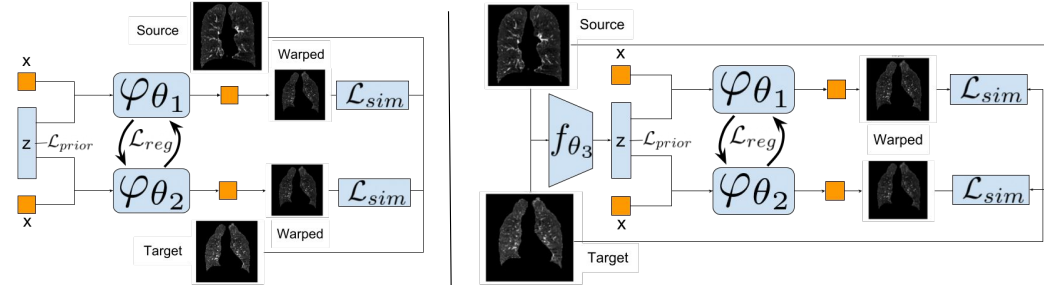
## 4.3 Learning the Latent Space of Deformations



Figure 1: NePhi in optimization-based registration and learning-based registration. x denotes the position of a point.

The proposed NePhi approach is a representation of a displacement vector field that can be used either in an optimization setting as in Eq. 2 or in a learning-based setting as in Eq. 3. Here, we describe how NePhi can be used for both settings. Fig. 1 illustrates the approach.

**Optimization-based `NePhi`.** Here, we optimize over the latent codes and parameters of $\varphi_{\theta_1}$ and $\varphi_{\theta_2}$ by minimizing the following loss

$$\{z, \theta_1, \theta_2\} = \arg\min \mathcal{L}_{\text{sim}} \left( I^A \circ \varphi_{\theta_1}(\cdot, z), I^B \right) + \mathcal{L}_{\text{sim}} \left( I^A, I^B \circ \varphi_{\theta_2}(\cdot, z) \right)$$
$$+ \lambda_1 \mathcal{L}_{\text{reg}}(\theta_1, \theta_2, z) + \lambda_2 \mathcal{L}_{prior}(z), \quad (8)$$

where $\mathcal{L}_{sim}$ is the similarity loss. In this work, we use normalized cross correlation as the similarity loss. As $\varphi$ is a continuous function that can represent transformations at any spatial location, we randomly sample three different sets of points within the image space to compute $\mathcal{L}_{\text{sim}} \left( I^A \circ \varphi_{\theta_1}, I^B \right)$, $\mathcal{L}_{\text{sim}} \left( I^A, I^B \circ \varphi_{\theta_2} \right)$ and $\mathcal{L}_{\text{reg}}$. The number of points contained in each set can be different. By sampling few points wwe can control the memory consumption.

**Learning-based `NePhi`.** We condition `NePhi` on latent codes because we want `NePhi` to provide equivalent functionality to the voxel representation used in learning-based registration. For a learning-based registration network, the CNN can directly predict the transformation given a pair of images during inference. For `NePhi`, we achieve a similar functionality by letting a neural network $f_{\theta_3}$ predict the latent codes $z_g$ and $z_l$ based on an input image pair. We use the same loss as in Eq. 8 but also optimize over the parameters of $f_{\theta_3}$ and `NePhi` as

$$\{\theta_3, \theta_1, \theta_2\} = \arg\min \mathcal{L}_{\text{sim}} \left( I^A \circ \varphi_{\theta_1}(\cdot, z), I^B \right) + \mathcal{L}_{\text{sim}} \left( I^A, I^B \circ \varphi_{\theta_2}(\cdot, z), I^B \right)$$
$$+ \lambda_1 \mathcal{L}_{\text{reg}}(\theta_1, \theta_2, z) + \lambda_2 \mathcal{L}_{prior}(z), z = f_{\theta_3}(I^A, I^B). \quad (9)$$

Similarly, we only need to sample a very small set of points to obtain regular transformations via `NePhi`. For more implementation details, refer to the Appendix.
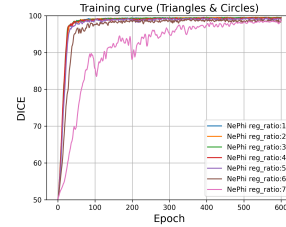
## 5 Experiments

### 5.1 Datasets

**Triangles and Circles (T&C)** is a synthetic 2D dataset [11] where images are either triangles or circles with either solid or hollow shapes. We generate a training set of 10000 pairs of images that of size $512 \times 512$ and a validation set with 1000 pairs of images with the same size for both solid shape mode and hollow mode.

**COPDGene** [31]. We use a subset of 999 inspiratory/expiratory lung CT pairs from the **COPDGene** study[2] [31] with provided lung segmentation masks for training. We resample the images to isotropic spacing of 2mm, leading to a volume of size $175 \times 175 \times 175$. CT intensities correspond to Hounsfield units. We clamp them within $[-1000, 0]$ and scale them linearly to $[0, 1]$. We then apply the lung segmentation mask to the images to extract the lung region of interest (ROI). We use 899 pairs for training and 100 pairs for validation.
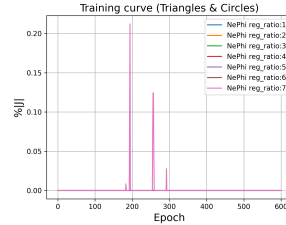
(a)

**DirLab**[4]. This dataset contains 10 pairs of inspiration/expiration lung CT images (from COPDGene, different from the 999 pairs above) with 300 anatomical landmarks per pair, manually identified by an expert in thoracic imaging. We process these images in the same way as all other **COPDGene** images. This dataset is only used to test a trained network.

### 5.2 `NePhi` is More Efficient Than a Voxel Transformation

To demonstrate the memory efficiency of `NePhi` during training compared to a voxel-represented DVF, we conducted experiments on the the **T&C** solid dataset. We focused on the optimization-based registration setting and analyzed `NePhi`'s training curve and model statistics. We randomly generated a pair of images from the dataset, each with a shape of 512x512. The regularizer ratio, defined as the number of points where we computed $\mathcal{L}_{reg}$ over the volume of the

(b)

Figure 2: `NePhi` with varying regularizer ratio in **optimization-based registration**.

images used in the optimization, was varied according to $\frac{1}{4^n}$. By examining how this variation affected

---

registration accuracy and transformation regularity, we gained insights. The result is shown in Fig. 2a and Fig. 2b. It is evident from these figures that reducing the regularizer ratio does not significantly impact the regularity of the obtained transformations until reaching a ratio of approximately $0.0244\%$, which corresponds to approximately $64$ points ($2.44e^{-4} \times 512 \times 512$). Interestingly, since lowering the regularizer ratio does not affect regularity while keeping the similarity ratio the same, the registration accuracy remains unaffected. This experiment highlights that in the optimization-based registration setting, one can achieve satisfactory accuracy by evaluating regularity using a significantly reduced number of points, resulting in low memory usage.
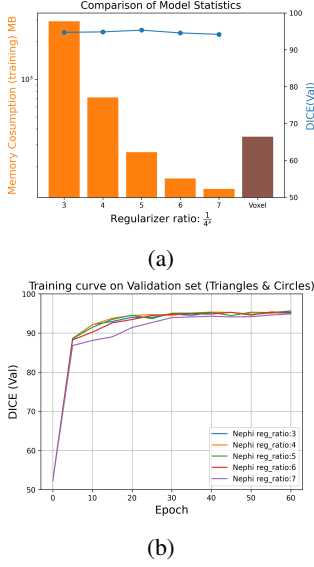


(a)



(b)

Figure 3: The peak memory consumption when training a **learning-based registration** with `NePhi` compared to voxel representation and the corresponding DICE curve of the validation set during training.
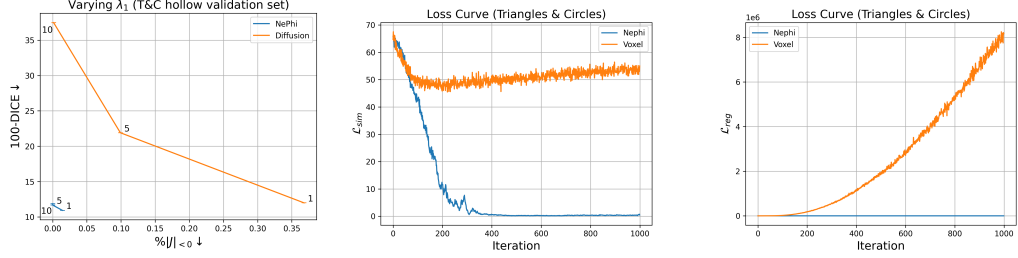
To study whether the same observation holds in learning-based registration, we train a CNN to predict `NePhi` while varying the regularizer ratio. The networks is trained on a **T&C** solid dataset that contains 10000 pairs of training images and 1000 pairs of validation images. The DICE on the validation set during the training is shown in Fig. 3b and the corresponding statistics of the network and training is shown in Fig. 3a. The metric $\%|J|_{<0}$, which indicates irregularity of the transformation, on the validation set is $0$ for all the trained models with varying regularizer ratio and for CNN+Voxel model (single-res. GradICON). And in Fig. 3a, we can see that CNN+`NePhi` consumes less memory than CNN+Voxel when achieving the same regularity of the predicted transformations.

### 5.3 `NePhi` Provides Better Regularity

Image registration is ill-posed as multiple transformation solutions can yield the same similarity loss. The choice of regularizer determines which set of transformations is preferred. Therefore, in this section, we examine how `NePhi` preserves regularity from various perspectives. Specifically, we investigate the following aspects: 1) Does `NePhi` exhibit a similar ability to balance between similarity and regularity as the voxel-based gradient inverse consistency (GradICON)? 2) Can `NePhi` maintain regularity in an optimization-based registration setting? 3) Does gradient inverse consistency promote better optimal solutions compared to diffusion regularizer in optimization-based registration? We will delve into the experiments and their motivations in detail in the subsequent discussion.

**Does `NePhi` exhibit a similar ability to balance between similarity and regularity as the voxel-based gradient inverse consistency (GradICON)?** For a pair of images that contain complex transformation in between, the more similar the registered images are, the more complex the transformation solution is, and easier to produce irregular deformations. Thus, the ability to sacrifice as little as possible of registration accuracy to fulfill transformation map regularity is important. In [40], the authors shows that gradient inverse consistency results in a better trade off between the similarity of the images and the regularity of the transformations in a learning-based registration setting with voxel representations. We show that such a property is preserved when applying gradient inverse consistency to neural transformation fields. We train learning-based registration networks with `NePhi` or with a neural transformation regularized by a diffusion regularizer. The experiment is conducted over the **T&C** hollow dataset. We show the DICE score and $\%|J|_{<0}$ of the trained networks over **T&C** the hollow validation set in Fig. 4a. We observe that the results of `NePhi` is clustered at the left bottom of Fig. 4a. When we loosen the regularizer factor $\lambda_1$, the regularity does not drop too much for `NePhi` compared to the transformations using the diffusion regularizer.

**Can `NePhi` maintain regularity in an optimization-based registration setting?** To pursue the best registration accuracy in practice, a post-processing step called instance optimization (IO) is commonly adopted after using a learning-based registration network for inference. Specifically, one conducts optimization either over the neural network or over the inferred transformation map for one pair of images for further accuracy. If one optimizes over the voxel-represented transformations with gradient inverse consistency in the optimization-based setting, the optimal solution is not regular (Fig. 4c). One needs to optimize the parameters of the overall CNN instead of the inferred transformation

6

(a) `NePhi` shows better trade-off between similarity and regularity than diffusion regularization applied to a MLP transform.

(b) $\mathcal{L}_{sim}$ of `NePhi` and voxel representation in the **optimization-based registration**.

(c) $\mathcal{L}_{reg}$ curve of `NePhi` and voxel representation in the **optimization-based registration**.

Figure 4: Experiments exploring the regularity properties of `NePhi`.
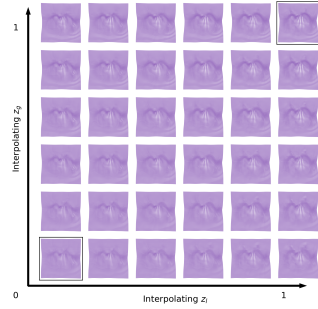
map in IO, at the same large memory cost during testing and training. Since `NePhi` contains MLPs which presumably have some implicit regularization, it is plausible to expect that `NePhi` can achieve regularity in the optimization-based registration setting. To validate this hypothesis, we compared the loss curves of `NePhi` and voxel-based transformation with gradient inverse consistency in the optimization-based registration setting. Fig. 4b and Fig. 4c show the results. Notably, the voxel-based representation converges to a much poorer similarity measure compared to `NePhi`.

**Does gradient inverse consistency result in better solutions compared to using a diffusion regularizer in optimization-based registration?** Since `NePhi` can provide regularity in an optimization-based setting, it immediately raises the question whether `NePhi` provides a better solution compared to other regularizers that also supports IO without optimizing the CNN. Thus, we compare `NePhi` to a neural transformation with the diffusion regularizer and to a neural transformation without any regularizer in an optimization-based registration setting. This experiment is conducted on one randomly generated pair of images from **T&C** and **T&C** hollow dataset. The loss curves are shown in Fig. 6. Fig. 6b and Fig. 6e show that `NePhi` is the only approach that results in zero foldings during optimization. Though, the metric $\%|J|_{<0}$ goes down to zero for the transformation with the diffusion regularizer. The irregularity of the transformation maps during the iterations might lead to a suboptimal solution, as the converged DICE curve of Diff.$(\lambda_1 = 1)$ converges to a lower DICE compared to `NePhi` in Fig. 6a and Fig. 6d. The visualization of the final registration in Fig. 6c and Fig. 6d also supports this conclusion.
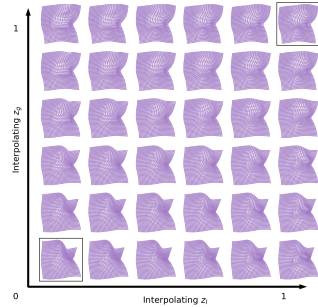
### 5.4 Interpolation of the Learned Latent Space

The expressiveness of the latent space of `NePhi` is critical, given that We want `NePhi` to be used in learning-based setting. we conduct the following two experiments to study whether `NePhi` can represent regularized transformations.

**Is `NePhi` generalizable?** In this experiment, we investigate whether a `NePhi` latent space, constructed through the optimization of a single `NePhi` model with a set of transformations, can be utilized to sample new regularized transformations. We use the **COPDGene**



(a) Interpolation of two latent codes from the *latent code book* on **COPDGene** dataset.



(b) Interpolation of two latent codes *predicted* by the CNN on **T&C** hollow dataset.

dataset in this experiment. Initially, we assign a latent code, comprising a $z_g$ and a $z_l$, to each image pair in the dataset. Consequently, a latent code table encompassing the entire dataset is generated. Subsequently, we perform joint optimization over the latent code table and the parameters of `NePhi` in an optimization-based registration setting. It is important to note that the parameters of `NePhi` are shared across all image pairs in the dataset. This experimental design aims to simulate the learning-based registration setting without a CNN encoder $f_{\theta_3}$. Upon completion of the optimization process, we linearly interpolate between two random latent codes from the latent code table. The
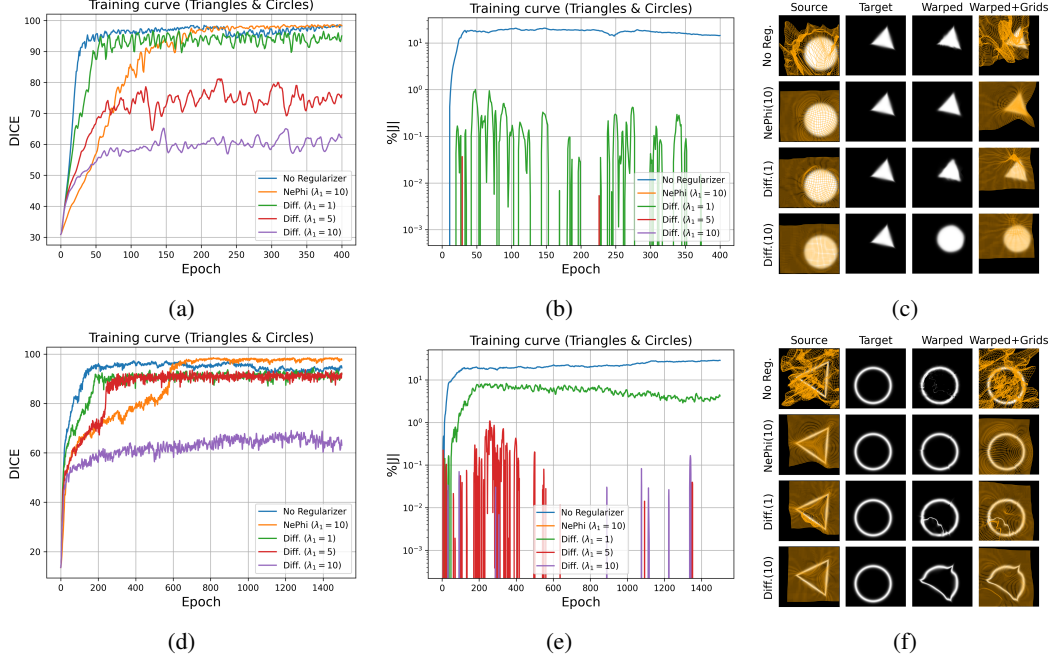
Figure 6: Comparisons between using diffusion regularizer and gradient inverse consistency (`NePhi`) on continuous representation of transformations in optimization-based registration. (a), (b) and (c) are the corresponding DICE curve, $\%|J|_{<0}$ curve and qualitative results of one pair of **C&T** solid images. (d), (e) and (f) are the results of the same experiment but on one pair of **C&T** hollow images.

resulting transformations, populated by the generated latent codes, are illustrated in Fig. 5a. Notably, the bottom left and top right transformations correspond to image pairs from the training set, while the remaining transformations are synthesized through latent code interpolation. As observed, all the transformation maps exhibit smoothness, with deformations gradually transitioning from the bottom left to the top right. Quantitatively, the mean $\%|J|_{<0}$ of the synthesized transformations depicted in Fig. 5a is 0, indicating the absence of folds in the synthesized transformations.

**Will the latent space learned via $f_{\theta_3}$ preserve regularity?** Having established that `NePhi` possesses the capacity to learn a latent space that maintains deformation regularity without an encoder $f_{\theta_3}$, it is important to investigate whether the same holds true for the latent space learned through $f_{\theta_3}$ in the learning-based registration setting. To address this question, we train $f_{\theta_3}$ with `NePhi` on the **T&C** hollow dataset. Subsequently, we perform interpolation between the latent codes predicted by $f_{\theta_3}$ using two randomly selected image pairs from the training set. From Fig. 5b we observe a smooth transition of transformations from the bottom left to the top right. Additionally, it is noteworthy that as the interpolation ratio of $z_g$ changes from 0 to 1, the global deformation becomes increasingly similar to the transformation in the top right. This finding suggests that $z_g$ and $z_l$ capture deformations at different levels, which aligns with our intended design.

## 5.5 Registration Performance on the DIRLab Benchmark

Throughout the preceding sections, we have examined `NePhi`'s properties from various perspectives, showing its suitability for both optimization-based and learning-based registration. In this section, we proceed to evaluate `NePhi` quantitatively on inspiration/expiration lung registration.

In the context of optimization-based registration, we initially evaluate `NePhi`'s performance and compare it with state-of-the-art (SOTA) multi-resolution optimization-based registration methods PTVReg[3] [42] and RRN [13]. It is essential to note that our primary focus in this work lies in the neural representation aspect rather than the registration framework itself. Consequently, we solely test `NePhi` within the single-resolution registration framework. However, as discussed in Sec. 2.1, multi-resolution frameworks are commonly employed to enhance the registration accuracy of pioneering single-resolution registration methods. Therefore, it is expected that `NePhi`'s registration

---

[3]The results are obtained from PTVReg's official github repository. Thus, it is slightly different from the value reported in PTVReg's paper.

performance may be inferior to the SOTA methods due to the differences in framework design. Table 1 presents the results of the performance comparison. While `NePhi` yields satisfactory results, it does not attain the performance level of the SOTA methods listed in Table 1. Nonetheless, it is crucial to emphasize that `NePhi` exhibits a significantly lower runtime of 150 seconds in contrast to the SOTA methods. We attribute the relatively lower performance of `NePhi` to the single-resolution registration framework we have currently employed, as well as the chosen sampling strategy, which can be improved on in future work.

| | Method | mTRE(mm)$\downarrow$ | $\%|J|_{<0}\downarrow$ | Time(s) |
|---|---|---|---|---|
| | Initial | 23.36 | — | — |
| M | PTVReg [42] | 0.84 | 0.60 | 442 |
| | RRN [13] | 0.83 | — | — |
| S | NePhi | 1.73 | **3.4e-5** | **150** |

Table 1: Performance on **DirLab** in optimization-based registration. M and S denotes multi-resolution framework and single-resolution framework, respectively.

For learning-based registration, we train a CNN with `NePhi` on the **COPDGene** dataset and evaluate the trained network on **DirLab**. As mentioned earlier, our main focus is on the neural representation aspect, thus we test `NePhi` using a straightforward CNN to encode the latent vevtors within a single-resolution registration framework. To ensure a fair comparison, we compare the results of CNN+`NePhi` with those of single-resolution GradICON, which utilizes a voxel representation. From Table 2 we observe that the registration network trained on `NePhi` performs similarly to single-resolution GradICON in terms of registration accuracy. However, the CNN+`NePhi` approach exhibits better regularity, uses fewer parameters, and has lower peak memory use during training compared to single-resolution GradICON.

Although it is expected that CNN+`NePhi` in the single-resolution registration setting would perform worse than multi-resolution GradICON, the performance gap can be largely closed through IO. Additionally, supported by the exploration in Sec. 5.3, we conduct IO without $f_{\theta_3}$, leading to less memory use and faster runtime compared to GradICON with IO.

| | Method | mTRE$\downarrow$ (mm) | $\%|J|_{<0}\downarrow$ | Rep. (M) | Enc. (M) | Time(Inference) (ms) | Peak Mem.(Train) (MB) |
|---|---|---|---|---|---|---|---|
| | Initial | 23.36 | — | — | — | — | — |
| M | LapIRN | 4.24 | 1.1e-2 | 0 | 0.92 | 235 | |
| | GradICON | 1.93 | 2.6e-4 | 0 | 70.68 | 161 | 13482 |
| S | GradICON | 5.41 | 1.4e-4 | 0 | 17.67 | 202 | 6394 |
| | GradICON+IO | 2.10 | 3.5e-4 | 0 | 17.67 | 53200 | 6394 |
| | NePhi | 5.45 | **0.0** | 0.28 | 3.07 | 461 | **2648** |
| | NePhi+IO | 1.88 | 0.0 | 0.28 | 0 | 7566 | 2496 |

Table 2: Performance of the registration network on **DirLab**. M and S denotes multi-resolution framework and single-resolution framework, respectively. Rep. is the parameter number of the representation and Enc. is the parameter of the encoder network. The inference time for IO means the time to finish the optimization iterations.

## 6   Conclusion

In this study, we introduced `NePhi`, an approximately diffeomorphic neural deformation representation. Our findings demonstrate that `NePhi` is effective in both an optimization-based and, most importantly, in a learning-based registration setting. Furthermore, we showed that `NePhi` can achieve comparable registration performance to voxel transformations within a single-resolution learning-based framework. However, there are certain *limitations* to our work. We acknowledge that we did not yet investigate the performance of `NePhi` within a multi-resolution framework, as discussed in Sec. 5.5. We believe that alternative advanced structures, such as patch-based registration frameworks or single-resolution transformer encoders, may be better suited for `NePhi` compared

to the multi-resolution approaches employed in existing methods. We take pride in being the first study to show the ability to predict neural deformations using a CNN encoder, while attaining similar results to voxel transformations within a basic single-resolution framework. We hope that our work is a starting point to inspire further research.

## 7 Acknowledgement

## References

[1] John Ashburner. A fast diffeomorphic image registration algorithm. *NeuroImage*, 38(1):95–113, 2007.

[2] Brian Avants, Nick Tustison, and Gang Song. Advanced normalization tools (ants). *Insight J*, 1–35, 11 2008.

[3] Guha Balakrishnan, Amy Zhao, Mert R. Sabuncu, John V. Guttag, and Adrian V. Dalca. Voxel-morph: A learning framework for deformable medical image registration. *IEEE Transactions on Medical Imaging*, 38(8):1788–1800, 2019.

[4] Richard Castillo, Edward Castillo, David Fuentes, Moiz Ahmad, Abbie M Wood, Michelle S Ludwig, and Thomas Guerrero. A reference dataset for deformable image registration spatial accuracy evaluation using the COPDgene study archive. *Phys. Med. Biol.*, 58(9):2861, 2013.

[5] Xu Chen, Yufeng Zheng, Michael J Black, Otmar Hilliges, and Andreas Geiger. Snarf: Differentiable forward skinning for animating non-rigid neural implicit shapes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11594–11604, 2021.

[6] Adrian Dalca, Guha Balakrishnan, John Guttag, and Mert Sabuncu. Unsupervised learning of probabilistic diffeomorphic registration for images and surfaces. *Medical Image Analysis*, 57:226–236, 2019.

[7] Shivam Duggal and Deepak Pathak. Topologically-aware deformation fields for single-view 3d reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1536–1546, 2022.

[8] Chen Gao, Ayush Saraf, Johannes Kopf, and Jia-Bin Huang. Dynamic view synthesis from dynamic monocular video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5712–5721, 2021.

[9] Deepak Ghimire, Dayoung Kil, and Seong-heum Kim. A survey on efficient convolutional neural networks and hardware acceleration. *Electronics*, 11(6):945, 2022.

[10] Philip-William Grassal, Malte Prinzler, Titus Leistner, Carsten Rother, Matthias Nießner, and Justus Thies. Neural head avatars from monocular rgb videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18653–18664, 2022.

[11] Hastings Greer, Roland Kwitt, François-Xavier Vialard, and Marc Niethammer. ICON: Learning regular maps through inverse consistency. In *ICCV*, 2021.

[12] Kun Han, Shanlin Sun, Xiangyi Yan, Chenyu You, Hao Tang, Junayed Naushad, Haoyu Ma, Deying Kong, and Xiaohui Xie. Diffeomorphic image registration with neural velocity field. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1869–1879, 2023.

[13] Xinzi He, Jia Guo, Xuzhe Zhang, Hanwen Bi, Sarah Gerard, David Kaczka, Amin Motahari, Eric Hoffman, Joseph Reinhardt, R Graham Barr, et al. Recursive refinement network for deformable lung registration between exhale and inhale ct scans. *arXiv preprint arXiv:2106.07608*, 2021.

[14] Mattias P. Heinrich, Mark Jenkinson, Michael Brady, and Julia A. Schnabel. Globally optimal deformable registration on a minimum spanning tree using dense displacement sampling. In *MICCAI*, volume 7512, pages 115–122, 2012.

[15] Mattias P. Heinrich, Bartlomiej W. Papiez, Julia A. Schnabel, and Heinz Handels. Non-parametric discrete registration with convex optimisation. In *Biomedical Image Registration - 6th International Workshop, WBIR*, volume 8545, pages 51–61, 2014.

[16] Jiahui Lei and Kostas Daniilidis. Cadex: Learning canonical deformation coordinate space for dynamic surface representation via neural homeomorphism. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6624–6634, 2022.

[17] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6498–6508, 2021.

[18] Jia-Wei Liu, Yan-Pei Cao, Weijia Mao, Wenqiao Zhang, David Junhao Zhang, Jussi Keppo, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Devrf: Fast deformable voxel radiance fields for dynamic scenes. *arXiv preprint arXiv:2205.15723*, 2022.

[19] Lingjie Liu, Marc Habermann, Viktor Rudnev, Kripasindhu Sarkar, Jiatao Gu, and Christian Theobalt. Neural actor: Neural free-view synthesis of human actors with pose control. *ACM Transactions on Graphics (TOG)*, 40(6):1–16, 2021.

[20] Risheng Liu, Zi Li, Xin Fan, Chenying Zhao, Hao Huang, and Zhongxuan Luo. Learning deformable image registration from optimization: Perspective, modules, bilevel training and beyond. *IEEE Transactions on Pattern Analysis Machine Intelligence*, 44(11):7688–7704, 2022.

[21] Frederik Maes, Dirk Vandermeulen, and Paul Suetens. Comparative evaluation of multiresolution optimization strategies for multimodality image registration by maximization of mutual information. *Medical Image Analysis*, 3(4):373–386, 1999.

[22] Ishit Mehta, Michaël Gharbi, Connelly Barnes, Eli Shechtman, Ravi Ramamoorthi, and Manmohan Chandraker. Modulated periodic activations for generalizable local functional representations. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14214–14223, 2021.

[23] Tony C. W. Mok and Albert C. S. Chung. Fast symmetric diffeomorphic image registration with convolutional neural networks. In *IEEE CVPR*, pages 4643–4652, 2020.

[24] Tony C. W. Mok and Albert C. S. Chung. Large deformation diffeomorphic image registration with laplacian pyramid networks. In *Medical Image Computing and Computer Assisted Intervention*, volume 12263, pages 211–221, 2020.

[25] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Occupancy flow: 4d reconstruction by learning particle dynamics. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5379–5389, 2019.

[26] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 165–174, 2019.

[27] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5865–5874, 2021.

[28] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M Seitz. Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields. *arXiv preprint arXiv:2106.13228*, 2021.

[29] Sida Peng, Junting Dong, Qianqian Wang, Shangzhan Zhang, Qing Shuai, Xiaowei Zhou, and Hujun Bao. Animatable neural radiance fields for modeling dynamic human bodies. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14314–14323, 2021.

[30] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10318–10327, 2021.

[31] Elizabeth A Regan, John E Hokanson, James R Murphy, Barry Make, David A Lynch, Terri H Beaty, Douglas Curran-Everett, Edwin K Silverman, and James D Crapo. Genetic epidemiology of COPD (COPDGene) study design. *COPD: J. Chronic Obstr. Pulm. Dis.*, 7(1):32–43, 2011.

[32] Ruizhi Shao, Hongwen Zhang, He Zhang, Mingjia Chen, Yan-Pei Cao, Tao Yu, and Yebin Liu. Doublefield: Bridging the neural surface and radiance fields for high-fidelity human reconstruction and rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15872–15882, 2022.

[33] Zhengyang Shen, Xu Han, Zhenlin Xu, and Marc Niethammer. Networks for joint affine and non-parametric image registration. In *IEEE CVPR*, pages 4224–4233, 2019.

[34] Hanna Siebert, Lasse Hansen, and Mattias P. Heinrich. Fast 3d registration with accurate optimisation and little learning for learn2reg 2021. In *Biomedical Image Registration, Domain Generalisation and Out-of-Distribution Analysis*, volume 13166, pages 174–179, 2021.

[35] Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. *Advances in Neural Information Processing Systems*, 33, 2020.

[36] Shanlin Sun, Kun Han, Deying Kong, Chenyu You, and Xiaohui Xie. Mirnf: medical image registration via neural fields. *arXiv preprint arXiv:2206.03111*, 2022.

[37] Wei Sun, Wiro J Niessen, and Stefan Klein. Free-form deformation using lower-order b-spline for nonrigid image registration. In *MICCAI*, pages 194–201, 2014.

[38] Ramana Subramanyam Sundararaman, Riccardo Marin, Emanuele Rodola, and Maks Ovsjanikov. Reduced representation of deformation fields for effective non-rigid shape matching. *Advances in Neural Information Processing Systems*, 35:10405–10420, 2022.

[39] Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *Advances in Neural Information Processing Systems*, 33:7537–7547, 2020.

[40] Lin Tian, Hastings Greer, François-Xavier Vialard, Roland Kwitt, Raúl San José Estépar, and Marc Niethammer. Gradicon: Approximate diffeomorphisms via gradient inverse consistency. *arXiv preprint arXiv:2206.05897*, 2022.

[41] Edgar Tretschk, Ayush Tewari, Vladislav Golyanik, Michael Zollhöfer, Christoph Lassner, and Christian Theobalt. Non-rigid neural radiance fields: Reconstruction and novel view synthesis of a dynamic scene from monocular video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12959–12970, 2021.

[42] Valery Vishnevskiy, Tobias Gass, Gabor Szekely, Christine Tanner, and Orcun Goksel. Isotropic total variation regularization of displacements in parametric image registration. *IEEE transactions on medical imaging*, 36(2):385–395, 2016.

[43] Ziyu Wang, Yu Deng, Jiaolong Yang, Jingyi Yu, and Xin Tong. Generative deformable radiance fields for disentangled image synthesis of topology-varying objects. *arXiv preprint arXiv:2209.04183*, 2022.

[44] Jelmer M Wolterink, Jesse C Zwienenberg, and Christoph Brune. Implicit neural representations for deformable image registration. In *International Conference on Medical Imaging with Deep Learning*, pages 1349–1359. PMLR, 2022.

[45] Nian Wu and Miaomiao Zhang. Neurepdiff: Neural operators to predict geodesics in deformation spaces. *arXiv preprint arXiv:2303.07115*, 2023.

[46] Hongyi Xu, Thiemo Alldieck, and Cristian Sminchisescu. H-nerf: Neural radiance fields for rendering and temporal reconstruction of humans in motion. *Advances in Neural Information Processing Systems*, 34:14955–14966, 2021.

[47] Xiao Yang, Roland Kwitt, Martin Styner, and Marc Niethammer. Quicksilver: Fast predictive image registration–a deep learning approach. *NeuroImage*, 158:378–396, 2017.

[48] Ruiqi Zhang and Jie Chen. Ndf: Neural deformable fields for dynamic human modelling. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXII*, pages 37–52. Springer, 2022.

[49] Yufeng Zheng, Victoria Fernández Abrevaya, Marcel C Bühler, Xu Chen, Michael J Black, and Otmar Hilliges. Im avatar: Implicit morphable head avatars from videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13545–13555, 2022.

[50] Jing Zou, Noémie Debroux, Lihao Liu, Jing Qin, Carola-Bibiane Schönlieb, and Angelica I Aviles-Rivero. Homeomorphic image registration via conformal-invariant hyperelastic regularisation. *arXiv preprint arXiv:2303.08113*, 2023.

# A Implementation Details

## A.1 Structure of NePhi

NePhi contains one forward transformation $\varphi_{\theta_1}$ and one backward transformation $\varphi_{\theta_2}$. The transformation maps are modeled adding a as displacement vector field to the identity transform, namely $\varphi_{\{\theta_1,\theta_2\}}(x, z(x)) = x + u_{\{\theta_1,\theta_2\}}(x, z(x))$. To be noted, our proposed NePhi does not restrict the type of implicit function we use for $u(x, z(x))$, as long as the input is kept as $(x, z(x))$ and the output is a vector of the same dimension as $x$. In this work, we use Modulated Periodic Activations from Mehta et al. [22] with MLPs as $u_\theta$. Different from standard MLPs which accept the concatenation of a positional embedding and a latent code as the input of the first MLP, Mehta et al. [22] proposed to construct the implicit function using two modules: a synthesis network and a modulation network. The synthesis network accepts the coordinates as the input and the modulation network accepts the latent code as the input. In each layer, the output of the modulation network is used as an amplifier over the output of the non-linear activation of the synthesis network. $u_{\theta_1}$ and $u_{\theta_1}$ share the same structure but have separate parameters $\theta_1$ and $\theta_2$. We use two hidden layers with 512 hidden feature dimensions with each followed by sin activation functions [35]. We use a Fourier feature mapping [39] to embed the coordinates.

## A.2 Structure of CNN encoder $f_{\theta_3}$

For the encoder, we use a convolutional neural network that accepts the concatenation of $I^A$ and $I^B$ as input and outputs the global latent code $z_g$. We branch from the intermediate feature maps followed by a convolutional layer to output the local latent code $z_l$. The position where we branch from the CNN is based on the shape of the intermediate feature maps. We choose to branch from the feature maps that are closest in shape to $16 \times 16$ for 2D and to $16 \times 16 \times 16$ for 3D feature maps. Fig. 7 shows the structure of $f_{\theta_3}$.
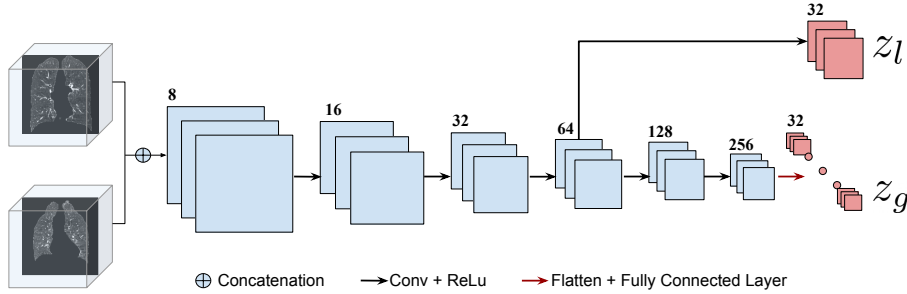


Figure 7: The structure of $f_{\theta_3}$ used in the experiments. The number of layers and the dimensions of the feature maps are kept the same for 2D and 3D registration. The position of the $z_l$ branch is adjusted according to the shape of the input image for each registration task to obtain a $z_l$ with a shape close to $16 \times 16$ or $16 \times 16 \times 16$ for 2D and 3D registration respecively.

## A.3 Similarity Loss $\mathcal{L}_{sim}$

We use Normalized Cross-Correlation (NCC) as the similarity measure and compute the similarity loss over a set of points drawn from a uniform distribution $q(x)$ across the image domain. The similarity loss is written as

$$\mathcal{L}_{sim}\left(I^A \circ \varphi_{\theta_1}(\cdot, z), I^B\right) = 1 - NCC(I^A \circ \varphi_{\theta_1}(x, z), I^B)^2, x \sim q(x). \tag{10}$$

Based on the difficulty of registration tasks, we sample different number of points when computing $\mathcal{L}_{sim}$. We sample 0.001, 0.01, 0.2 of the total number of pixel/voxel of the image for **T&C** soild, **T&C** hollow, and **COPDGene** dataset, respectively.

## A.4 Experimental Details

We set $\lambda_1 = 10$ and $\lambda_0 = 1e - 2$ for all experiments for learning-based registration and optimization-based registration. For both settings, we set the dimension of $z_g$ and $z_l$ to 32, respectively. For optimization-based registration, the shape of $z_l$ is $16 \times 16$ for 2D registration and $16 \times 16 \times 16$ for 3D registration regardless of the shape of the input images. For the learning-based framework, the shape of $z_l$ depends on the shape of the input images. We adjust the network according to Sec. A.2 to keep the shape of $z_l$ close to $16 \times 16$ for 2D registration and close to $16 \times 16 \times 16$ for 3D registration.