

Meta Deformation Network: Meta Functionals for Shape Correspondence

Daohan Lu^{1,4} and Yi Fang^{1,2,3}

¹ Multimedia and Visual Computing Lab, New York University, New York, United States.

² Tandon School of Engineering, New York University, New York, United States.

³ Department of Electrical and Computer Engineering, New York University, Abu Dhabi, United Arab Emirates.

⁴ College of Arts and Science, New York University, New York, United States.

Abstract. We present a new technique named *Meta Deformation Network* for 3D shape matching via deformation, in which a deep neural network maps a reference shape onto the parameters of a second neural network whose task is to give the correspondence between a learned template and query shape via deformation. We categorize the second neural network as a meta-function, or a function generated by another function, as its parameters are dynamically given by the first network on a per-input basis. This leads to a straightforward overall architecture and faster execution speeds, without loss in the quality of the deformation of the template. We show in our experiments that Meta Deformation Network leads to improvements on the MPI-FAUST Inter Challenge over designs that utilized a conventional decoder design that has non-dynamic parameters.

Keywords: dynamic neural networks, shape registration, shape correspondence

1 Introduction

In the pairwise 3D shape correspondence problem, a program is given two query shapes and is asked to compute all pairs of corresponding points between the two shapes. Recent works [7],[5] yielded improved results on pairwise 3D shape correspondence by aligning a template point cloud with individual query shapes, which gives a correspondence relationship between the template and each query shape. After the correspondence between the template and each query shape was obtained, this information was used to infer the correspondence between the two query shapes. This approach usually employed an encoder-decoder design with the encoder generating a feature embedding \mathbf{E} that captures the characteristics of the query shape \mathbf{S}_q . \mathbf{E} is then fed into the decoder along with the of points on the template shape as input in order to output a set of points representing the deformation of the template into the query shape. Every point in the template is said to be in correspondence with the point on the query shape that is nearest

to its deformation. Originally named in [12], we refer to this decoder scheme as *Latent Vector Concatenation* (LVC), since the input to the decoder is the concatenation of individual 3-D coordinates and the latent vector \mathbf{E} .

Although LVC is widely used in recent works to represent or deform 3-D shapes, [14],[7],[8],[19],[5], in this paper, we investigate the possibility of an alternative decoder structure and compare it against LVC on the task of computing correspondence for human 3-D scans. Specifically, we evaluate an alternative decoder design scheme that uses only the template point cloud as input but has dynamic parameters that are predicted by a neural network from \mathbf{E} and also outputs the deformed template points. We call this architecture *Meta Deformation Network* because the deformation process is carried out by a neural network whose parameters are not independent but generated by another neural network. The decoder could be thought of as a second-order function that is defined or generated by another function. This formulation leads to a speedup in training and testing, and the results on the MPI-FAUST Inter correspondence challenge show that the meta decoder yields improved correspondence accuracy over a similar design [5] that employs an LVC decoder.

2 Related Works

2.1 Dynamic Neural Networks

Dynamic Neural Networks are a design in which one network is trained to predict the parameters of another network in test-time in order to give the second network higher adaptability. The possibility of using dynamic networks on 2-D and more recently 3-D data has been researched by several authors. [17], [9], [1], [3] applied dynamic convolutional kernels to 2-D images to input-aware single-image super-resolution, short-range weather prediction, one-shot character classification and tracking, few-shot video frame prediction, respectively. These tasks demanded highly adaptive behaviors from the neural networks, and they show improvements when using dynamic neural networks over conventional methods that used static-parameter convolutional kernels. [10] was one of the first works to apply the idea of dynamically generated parameters to the 3-D domain, where it achieved state-of-the-art results on single-image 3D reconstruction on ShapeNet-core [4] compared to conventional methods in [20], [6], [16]. We are motivated by these works to further investigate the potential of dynamic neural networks. We build on previous studies by examining the effectiveness of dynamic neural networks for 3-D shape correspondence in which both the input shapes and output correspondence are of a 3-dimensional nature.

2.2 3-D Shape Correspondence

The registration of and correspondence among 3-D structures with non-rigid geometrical variations is a problem that has been extensively studied and one in which machine-learned based methods have had great success. MPI-FAUST [2] is

a widely used benchmark for finding correspondence between pairs of 3-D scans of humans with realistic noises and artifacts as a result of scanning real-world humans. There are various approaches to the task. [10] used a Siamese structure to learn compact descriptors of input shapes and using functional maps [13] to solve for inter-shape correspondences. [21] fitting query meshes to a pre-designed part-based model (“stitched puppet”) and aligning the models for the two query scans to get correspondence. [7] took a simpler approach, by choosing a set of points to be a template and getting inter-shape correspondences by computing the query shapes’ individual correspondences to the common template, effectively factoring the problem of computing correspondence between two unknown shapes into two problems each involving the correspondence between a known template and an unknown shape. This method proved to give the best correspondence accuracy on FAUST and is the inspiration for the training procedure of our method, the Meta Deformation Network. [7] obtains correspondence between a template and a query shape by using a multi-layer perceptron g to deform the points on templates into the shape of the query and obtaining correspondence by Euclidean proximity between the deformed points of the template and actual points on the query shape. However, g holds fixed parameters and takes a latent vector concatenated with template points as input. By contrast, our Meta Deformation Network has a decoder g that holds dynamic parameters, which adds more flexibility and adaptability of the deformation with respect to different query shapes.

3 Network Design

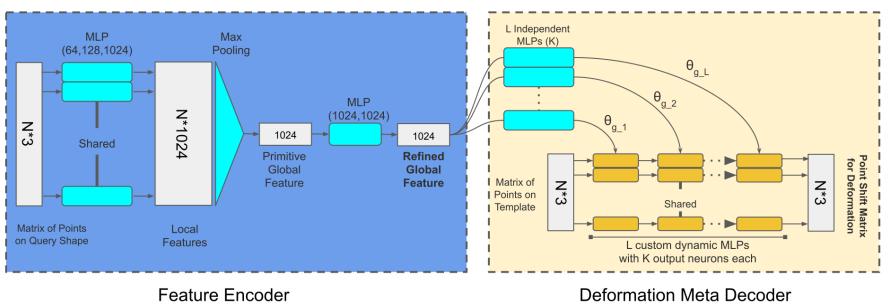


Fig. 1. Architecture Overview of the Meta Deformation Network. Blue rounded rectangles represent multi-layer perceptrons (MLPs) with numbers in parentheses specifying the number of neurons for all layers, yellow rounded rectangles represent dynamic MLPs, and white rectangles indicate feature sizes at different stages in the pipeline

3.1 Architecture Overview

A graphical description of the Meta Deformation Network is shown in fig. 1. We used a simplified version of the PointNet [15] encoder following [7], which we will denote as f , to output a feature vector \mathbf{E} of size 1024 from a set of points from the query shape \mathbf{S}_q . L independent MLPs, each denoted h_i , then \mathbf{E} takes as input outputs $\hat{\theta}_{g_i}$, the predicted optimal parameters for layer $i \in \{1, \dots, L\}$ of the meta decoder g . If layer i of g has K_{in} input neurons and K_{out} output neurons, then $\hat{\theta}_{g_i}$ is a vector of size $(K_{in} * K_{out} + K_{out} + K_{out})$ since we define the forward-pass for each layer of g as

$$g_i : x_{i+1} = (W_i x_i) * s_i + b_i \quad (1)$$

where $*$ denotes element-wise multiplication, and x_{out} is the pre-activation output. Note that we place a scaling term s in addition to the typical feed-forward calculation of an MLP with the intuition that it facilitates the learning of h (reasoning behind this intuition is given in 3.5).

With $\hat{\theta}_i$ computed, the meta decoder takes as input a point $\mathbf{p}_t \in \mathbf{S}_t$ outputs a 3-D residual vector $\Delta(\mathbf{p}_t)$ that will shift \mathbf{p}_t to the location of the corresponding point on the query shape \mathbf{S}_q . Note that the input layer to g is only a vector of 3 elements, which leads to a major speed up over a decoder that uses the LVC input, which take in a vector of size 1024+3 if using an \mathbf{E} . The difference in time is greater as the resolution of the template increases because the deformation computation is repeated N times to get the respective $\Delta(\mathbf{p}_t)$ for all $\mathbf{p}_t \in \mathbf{S}_t$.

In all experiments, we pick g 's structure $L = 6$ and $K_{out} = 64$ for every layer of g except the last, which has $K_{out} = 3$ since it outputs a 3 dimensional translation vector $\Delta\mathbf{p}_t$ for each input point \mathbf{p}_t .

3.2 Encoder

We used a simplified variant of the PointNet feature encoder [15]. The encoder, denoted f , applies a shared 3-layered MLP of output sizes (64, 128, 1024) to each point $\mathbf{p}_q \in \mathbf{S}_q$ to obtain a local feature of size $(N \times 1024)$ where N is the number of points chosen from the template and is equal to 6890 in all experiments. It then applies max-pooling along the first dimension to obtain a primitive global feature vector of size (1024), which is refined by a 2-layered MLP of output sizes (1024, 1024) to become the final global feature embedding \mathbf{E} . Formally,

$$\mathbf{E} = f_{\theta_f}(\mathbf{S}_q) \quad (2)$$

We pick $N = 6980$ for all of our experiments.

3.3 Parameter Predictor for Meta Decoder

Having defined operations performed by each layer of g , we define a set of 1-layered MLPs (or more accurately, SLPs), which we will denote as h_i , that maps

the global feature embedding \mathbf{E} to $\hat{\theta}_i$. Formally,

$$\begin{aligned}\hat{\theta}_{g_i} &= h_{i,\theta_{h_i}}(\mathbf{E}) = (W_i; s_i; b_i) \\ \hat{\theta}_g &= \{\hat{\theta}_{g_i} : i \in \{1, \dots, L\}\}\end{aligned}\tag{3}$$

So h_i has an input size of 1024 and an output size that matches the number of parameters needed for the i^{th} layer of g (see the calculation above eq. 1 in 3.1).

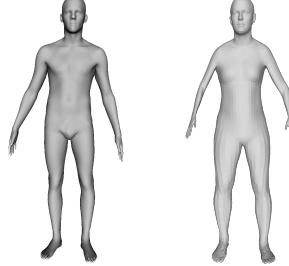


Fig. 2. Learnable Template. From left to right: initialization, template learned after 40 epochs of training

3.4 Learnable Template

Following [5], instead of choosing a template prior to training and making it permanently fixed, we initialize our template with a subsampled scan (so that $N = 6890$) from the MPI-FAUST training dataset and learns a translation vector v_i for each original point on the template $\mathbf{q}_{\mathbf{t}_i}$. So $\mathbf{p}_{\mathbf{t}_i} = \mathbf{q}_{\mathbf{t}_i} + v_i, i \in \{1, \dots, N\}$. This allows the network to learn an template that results in the best deformation correspondence during training. The chosen template initialization is depicted in fig 2. In test-time, we hold the template fixed by fixing the learned translation vectors $\{v_i\}$. Note that in test-time, if a high-resolution template is desired, we will no longer be able to use the learned translation vectors $\{v_i\}_{i=1}^N$ because each v_i is learned specifically for $\mathbf{p}_{\mathbf{t}_i}$. We will not have a dense set of translation vectors for every point in $\{p'_{t_j}\}_{j=1}^{N'}$ if we wish to use a high-resolution version (i.e. $N' > N$) of the initial template in test-time.

3.5 Meta Decoder

We formulate each layer of g as a customized version of a fully-connected layer (from eq. 1),

$$g_i : x_{i+1} = (W_i x_i) * s_i + b_i\tag{1}$$

The full network g outputs a residual vector for every input point $\mathbf{p}_t \in \mathbf{S}_t$ which will take it to the location of the corresponding point $\mathbf{p}_q \in \mathbf{S}_q$. Formally,

$$\begin{aligned}\Delta(\mathbf{p}_t) &= g_{\hat{\theta}_g}(\mathbf{p}_t) \\ \hat{\mathbf{p}}_q &= \mathbf{p}_t + \Delta(\mathbf{p}_t)\end{aligned}\tag{4}$$

Optimizing correspondence among query shapes is equivalent to minimizing the distance between $\hat{\mathbf{p}}_q$ and \mathbf{p}_q for every query shape. Thus, the formulations above leads to a straightforward supervised training loss with just a single term when training for correspondence.

3.6 Training and Losses

Since $\hat{\theta}_g$ is predicted, we only need to update the parameters of f and h in train-time. Let \mathbf{S}_q be the set of all query shapes in the training dataset, this gives the following supervised optimization problem, which is the same equation used in [7]:

$$\min_{\theta_f, \theta_h} \sum_{\mathbf{S}_q \in \mathbf{S}_q} \sum_{i=1}^N \|g_{\hat{\theta}_g}(\mathbf{p}_{t,i}) - \mathbf{p}_{q,i}\|^2 \tag{5}$$

$$\text{where } \hat{\theta}_g = h_{\theta_h}(\mathbf{E}), \mathbf{E} = f_{\theta_f}(\mathbf{S}_q)$$

Note that in the supervised case we assume having knowledge of the ground-truth location of the point $\mathbf{p}_{q,i} \in \mathbf{S}_q$ in correspondence with $\mathbf{p}_{t,i} \in \mathbf{S}_t$ for all $i \in 1, \dots, N$. However, we do not need any explicit information on the correspondence among the query shapes themselves as it is implicitly given by the correspondence with a common template.

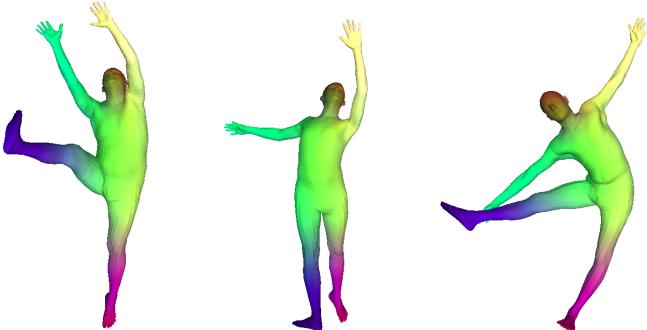


Fig. 3. Query Shapes from the training dataset by [7]

We adopt augmented SURREAL dataset by created T. Groueix et al. [7] using SMPL [11] and SURREAL [18], containing 229,984 meshes of synthetic

humans in various poses for training and 200 for validation. Samples of the training data are shown in fig. 3. We train the model using ADAM with at a learning rate of $2 * 10^{-5}$ for 25 epochs and then $2 * 10^{-5}$ for 15 epochs. This takes around 14 hours on a machine with a GTX 1080 Ti and a 6-core Intel CPU.

4 Inferencing Correspondences

4.1 Optimizing Query Shape Rotation

As in [7], we use a procedure during testing to improve the deformation before the optimization of the latent vector \mathbf{E} . Before calculating the correspondence, we rotate each query shape \mathbf{S}_q by selecting different pitch and yaw angles so that the deformation of the template given the rotated \mathbf{S}_q has the smallest Chamfer distance to the rotated \mathbf{S}_q . (When returning the correspondence output, we apply the inverse rotation matrices to recover the pre-rotation predicted locations).

$$\mathbf{S}_q \leftarrow \arg \min_{\alpha, \beta} \mathcal{L}^{\text{CD}}(h((R_{z,\beta} R_{y,\alpha}) \mathbf{S}_q), (R_{z,\beta} R_{y,\alpha}) \mathbf{S}_q) \quad (6)$$

where $R_{z,\beta}$ stands for the counterclockwise rotation matrix around the z axis by angle β (yaw), $R_{y,\alpha}$ stands for the counterclockwise rotation matrix around the y axis (pitch) by angle α , and $(R_{z,\beta} R_{y,\alpha}) \mathbf{S}_q$ denotes the set of points after applying the pitch and yaw rotation matrices to every point in \mathbf{S}_q . In practice, we try every combination of (α, β) where $\alpha \in [-\pi/2, \pi/2]$, $\beta \in [-\pi/4, \pi/4]$. The intervals are discrete with strides $\pi/100$ and $\pi/50$ respectively for α and β . This leads to a total of $100 * 25$ tries during which we remember the pair (α^*, β^*) that gives the best deformation and replace \mathbf{S}_q with $(R_{z,\beta^*} R_{y,\alpha^*}) \mathbf{S}_q$.

4.2 Latent Vector Optimization

As 6.1 illustrates, using $f(\mathbf{S}_q)$ directly as the feature embedding \mathbf{E} produces suboptimal deformations. To mitigate this problem, during testing only, we use $f(\mathbf{S}_q)$ as an initialization for \mathbf{E} and optimizes over \mathbf{E} for 3000 iterations using ADAM at a learning rate of $5 * 10^{-5}$ to implicitly find the $\hat{\theta}_g$ that minimizes the Chamfer distance between the deformation $\{g_{\hat{\theta}_g}(\mathbf{p}_t) : \mathbf{p}_t \in \mathbf{S}_t\}$ and the query shape $\{\mathbf{p}_q : \mathbf{p}_q \in \mathbf{S}_q\}$. We do not perform this step in training as we reduces training speed and we want to train the encoder to predict feature embeddings directly. The Chamfer loss given \mathbf{E} and \mathbf{S}_q is defined as:

$$\mathcal{L}^{\text{CD}}(\mathbf{E}, \mathbf{S}_q) = \sum_{\mathbf{p}_t \in \mathbf{S}_t} \min_{\mathbf{p}_q \in \mathbf{S}_q} \|\mathbf{p}_q - g_h(\mathbf{E})(\mathbf{p}_t)\|^2 + \sum_{\mathbf{p}_q \in \mathbf{S}_q} \min_{\mathbf{p}_t \in \mathbf{S}_t} \|\mathbf{p}_q - g_h(\mathbf{E})(\mathbf{p}_t)\|^2 \quad (7)$$

Algorithm 1: Algorithm for finding 3D shape correspondences [7]

Input : Query shape 1 \mathbf{S}_{q1} and query shape 2 \mathbf{S}_{q2}
Output: Set of 3D point correspondences \mathcal{C}

- 1 #Regression steps over latent code to find best deformation into \mathbf{S}_{q1} and \mathbf{S}_{q2}
- 2 $\mathbf{E}_1 \leftarrow \arg \min_{\mathbf{E}} \mathcal{L}^{CD}(\mathbf{E}, \mathbf{S}_{q1})$ # \mathbf{E} initialized with $f(\mathbf{S}_{q1})$
- 3 $\mathbf{E}_2 \leftarrow \arg \min_{\mathbf{E}} \mathcal{L}^{CD}(\mathbf{E}, \mathbf{S}_{q2})$ # \mathbf{E} initialized with $f(\mathbf{S}_{q1})$
- 4 $(\hat{\theta}_g)_1 \leftarrow h(\mathbf{E}_1)$
- 5 $(\hat{\theta}_g)_2 \leftarrow h(\mathbf{E}_2)$
- 6 $\mathcal{C} \leftarrow \emptyset$
- 7 # Matching of $\mathbf{p}_{q1} \in \mathbf{S}_{q1}$ to $\mathbf{p}_{q2} \in \mathbf{S}_{q2}$
- 8 **foreach** $\mathbf{p}_{q1} \in \mathbf{S}_{q1}$ **do**
- 9 $\hat{\mathbf{p}}_t \leftarrow \arg \min_{\mathbf{p}_t \in \mathbf{S}_t} \|g_{(\hat{\theta}_g)_1}(\mathbf{p}_t) - \mathbf{p}_{q1}\|^2$
- 10 $\hat{\mathbf{p}}_{q2} \leftarrow \arg \min_{\mathbf{p}_{q2} \in \mathbf{S}_{q2}} \|g_{(\hat{\theta}_g)_2}(\hat{\mathbf{p}}_t) - \mathbf{p}_{q2}\|^2$
- 11 $\mathcal{C} \leftarrow \mathcal{C} \cup \{(\mathbf{p}_{q1}, \hat{\mathbf{p}}_{q2})\}$
- 12 **end**
- 13 **return** \mathcal{C}

4.3 Finding Pairwise Correspondences

We use the algorithm presented in [7] to optimize over \mathbf{E} and find the correspondence between a pair of 3-D shapes, shown here in algorithm 1. Note that although we only compute correspondence between two shapes for the FAUST-Inter challenge, the same algorithm can easily be extended to predict correspondence across multiple shapes with an execution time linear to the number of query shapes by matching each shape to the common template.

5 Analysis on the Characteristics of the "Meta" Decoder

5.1 Defining the Meta Decoder

In a Latent Vector Concatenation decoder construction, a fixed-parameter decoder g takes as input the coordinates of a template point $\mathbf{p}_t \in \mathbf{S}_t$ concatenated with the feature embedding \mathbf{E}_q representing the query shape \mathbf{S}_q and outputs the predicted coordinates $\hat{\mathbf{p}}_q$ of the corresponding point in the query shape. Since we predict pairwise correspondence based on the query shapes' individual correspondences to the common template, the quality of the template deformation was crucial to the correspondence accuracy. In LVC, because the decoder has fixed parameters and needs to accurately deform all points on the template into corresponding points on various query shapes when given only \mathbf{p}_t and \mathbf{E}_q as input, its structure is often complex in terms of the number of trainable parameters, which could lead to problems like lengthy training times and the overfitting. Denoting the feature encoder as f and the decoder as g , this leads

to the formulation:

$$\mathbf{E}_q = f_{\theta_f}(\mathbf{S}_q) \quad (8)$$

$$\hat{\mathbf{p}}_q = g_{\theta_g}(\mathbf{p}_t; \mathbf{E}_q) \quad (9)$$

Where θ_f and θ_g are jointly optimized during training.

Our approach simplifies the structure of the decoder by formulating it as a concise dynamic multi-layer perceptron, that is, a MLP whose parameters are dynamically and uniquely determined by another neural network for each input. Denoting the encoder f and the dynamic decoder g , we have:

$$\hat{\theta}_g = f_{\theta_f}(\mathbf{S}_q) \quad (10)$$

$$\hat{\mathbf{p}}_q = g_{\hat{\theta}_g}(\mathbf{p}_t) + \mathbf{p}_t \quad (11)$$

Because $\hat{\theta}_g$ is determined by f , we only need to optimize θ_f during training. Note that the decoder computes a location residual instead of the coordinates of the actual deformed point $\hat{\mathbf{p}}_q$ that corresponds to \mathbf{p}_t . We found that predicting the residual was an easier task to learn for the network.

5.2 Mathematical differences between LVC and Meta Decoder

Although they perform different tasks, the neural networks in [7], [15], and [14] have all employed g as an LVC decoder that is an MLP with non-dynamic parameters. E. Mitchell pointed out in [12] that this formulation is equivalent to having a MLP that has fixed input and weights but a variable bias. Consider a case in which g needs to deform a template into query shapes, then the LVC deformation decoder g would take in of $\mathbf{p}_t \in \mathbf{S}_t$ concatenated with a feature embedding representing the query shape $\mathbf{E} = f_{\theta_f}(\mathbf{S}_q)$. Such is the case in [7]. Let vector x be the concatenation of \mathbf{p}_t and \mathbf{E} , that is, $x = (\mathbf{p}_t; \mathbf{E})$, then the first layer of the LVC MLP decoder will perform the following operation:

$$a = Wx + b \quad (12)$$

where a is the pre-activation output, W is the fixed weights matrix, and b is the fixed bias. We can rewrite x as the sum of two vectors y and z where $y = (\mathbf{p}_{t,x}, \mathbf{p}_{t,y}, \mathbf{p}_{t,z}, 0, 0, \dots)$ and $z = (0, 0, 0, \mathbf{E})$

$$a = Wx + b = W(y + z) + b = Wy + Wz + b \quad (13)$$

Note that W, b, x are fixed with regards to \mathbf{S}_q and only z changes with different \mathbf{S}_q , so the only term that varies with the query shape is Wz . Thus, Wz can be seen as a dynamic bias predicted from query shapes for the first layer, which has fixed weights and input (points on the template were fixed in test-time). Using the concatenation of \mathbf{p}_t and \mathbf{E} as input is therefore equivalent to having

an MLP that has fixed input, weights, and biases, except that a variable bias term Wz is added to layer 1’s pre-activation output that depends on \mathbf{S}_q . The low adaptability of the LVC decoder prevents it from producing high-quality deformations with a concise structure.

In the case of the meta decoder, we set all parameters of g to be dynamic. With this approach, the input x is simply a 3-D vector ($\mathbf{p}_{t,x}, \mathbf{p}_{t,y}, \mathbf{p}_{t,z}$), and the first layer of the decoder MLP performs the operation

$$a = (Wx) * s + b \quad (14)$$

where W, s, b are given by $h_1(f(\mathbf{S}_q))$. The process is similar for subsequent layers. As a result, all parameters of g are dynamic a function of the query shape \mathbf{S}_q . Defining g ’s parameters to be a function of \mathbf{S}_q offers the decoder more flexibility in computing the optimal translation that deforms the template into diverse query shapes. We also include an element-wise scaling factor s in the formulation of g to facilitate the learning of h , which maps \mathbf{E} to θ_g . To see this, suppose that a particular output neuron of g needs to be changed in order to compute the optimal deformation. Multiplying s_1 by λ has the same effect on the output as scaling all elements in the first row of W by λ , but is an easier change to predict for h since it requires predicting a single number instead of multiple numbers to capture this change. Overall, the meta decoder has fully adaptive parameters to different query shapes. By contrast, the LVC decoder has only a variable bias in the output of the first layer and fixed parameters everywhere else. The enhanced adaptability of g with respect to query shapes allows it to produce better deformations of the template with fewer parameters and at higher speeds.

6 Experiment Results

We test the Meta Deformation Network on the MPI-FAUST Inter and Intra Subject Challenges. The ”inter” challenge contains 40 pairs of 3D scans of real people with each pair consisting of two different people at different poses. In test-time, to get the optimal deformation, we use $f(\mathbf{S}_q)$ as an initialization for the feature embedding \mathbf{E} , and optimizes over \mathbf{E} for 3000 iterations using ADAM at a learning rate of $5 * 10^{-5}$ to implicitly find the $\hat{\theta}_g$ that minimizes the Chamfer distance between the deformation $\{g_{\hat{\theta}_g}(\mathbf{p}_t) : \mathbf{p}_t \in \mathbf{S}_t\}$ and the query shape $\{\mathbf{p}_q : \mathbf{p}_q \in \mathbf{S}_q\}$.

6.1 Unoptimized Deformations

We show the deformation of the template into the query shapes before the optimization of \mathbf{E} , and compare Meta Deformation Network against the non-meta deformation network developed by T. Deprelle et al. in [5]. The figures are included in fig. 4.

Both approaches give a reasonable unoptimized deformation of a template into the query shape, though both suffer from varying amounts of distortions

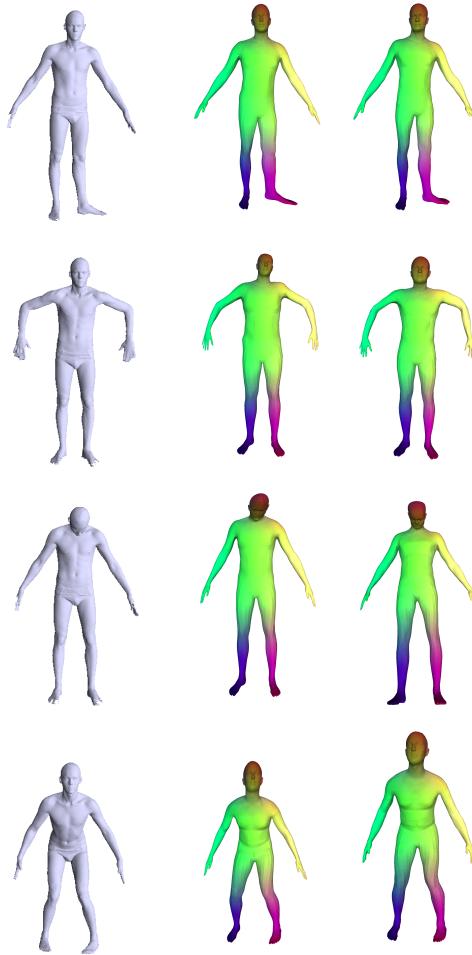


Fig. 4. Comparison of Meta Deformation Network and 3D-CODED with Learned templates [5] in unoptimized deformation quality. Left: Query Shape, middle: Meta Deformation Networks, right: 3D-CODED with Learned Template]

and imprecision from the actual query shape. The optimization step depicted in 6.2 over \mathbf{E} improves the quality of the deformation.

6.2 Feature Embedding Optimization

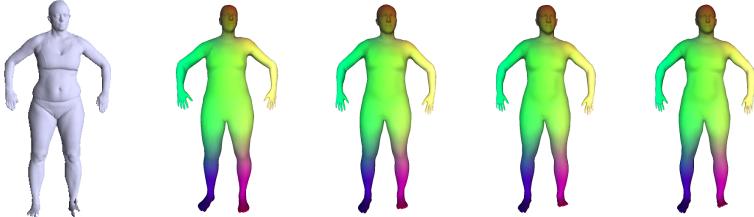


Fig. 5. The deformation of the template into the query shape after optimizing \mathbf{E} for different numbers of iterations. From left to right: ground-truth query shape, unoptimized deformation, deformation optimized for 200 iterations, 1,000 iterations, and 3,000 iterations. Note the little difference in quality between deformations optimized for 1,000 and 3,000 iterations

We also test how sensitive the quality of the template deformation is to different numbers of optimization steps for \mathbf{E} . We qualitatively demonstrate that the deformation continues to improve in quality as \mathbf{E} is optimized for more iterations, but the marginal gain diminishes as the number of iterations increases, with the difference between 1,000 and 3,000 iterations hardly noticeable (See fig. 5).

The imprecision of the unoptimized deformation suggests that our ShapeNet encoder f is unable to produce an optimal feature embedding \mathbf{E} that fully captures the characteristics of the query shape; the unoptimized deformation has a more generic structure lacking details of the person’s body shape, which the post-optimization deformation is able to capture. The relatively weak encoding power of our simplified ShapeNet encoder, can be alleviated by optimizing over \mathbf{E} given $f(\mathbf{S}_q)$ as its initialization. Also note that although the query shape has incomplete data around the person’s feet, our decoder has learned to reconstruct the missing parts from the input.

6.3 Quantitative Assessment

Though the Meta Deformation Network did not surpass the state-of-the-art performance, it showed an improvement over 3D-CODED w/. Learned 3D Translation Template, an approach that is comparable to ours in that it shares the same simplified PointNet encoder design and also applies a learnable translation matrix to the points of the base template, with the difference being that it its LVC decoder has fixed parameters. This shows quantitatively that having the

Method	FAUST-Inter Mean Error (cm)
Deep functional maps [13]	4.82
Stitched Puppet [21]	3.12
3D-CODED w/. Learned 3D Translation Template [5]	3.05
3D-CODED w/. Learned 3D Deformed Template [5]	2.76
Meta Deformation Network (Ours)	2.97

Table 1. Comparison of Meta Deformation Network to Network Architectures with Conventional Methods

meta decoder produces more accurate correspondences on FAUST-Inter than does a static-parameter decoder (while also having speedier execution).

We report results from training under a low-resolution template and hot-swapping in a different high-resolution template in inferencing (details for this procedure are given in 6.4). See the accuracy numbers are shown in table 1.

6.4 Transference onto Unseen Templates

Template used in inferencing	FAUST-Iter Mean Error (cm)
Learned Template (6890 Vertices)	3.07
Unseen Template (176K Vertices)	2.97

Table 2. Comparison of infernencing with a low-resolution learned template and a high-resolution new template.

In the training step, we have trained a feature encoder f and a parameter mapper h to enable g to predict deformations from a known template into unseen query shapes. Here, we also test how well the model performs when given a template and query shapes that are both unseen in training. We run an experimental setup where we find the rotation matrices and the feature embedding initialization for \mathbf{S}_q using the learned template \mathbf{S}_t but swaps in a different, higher-resolution template \mathbf{S}_t' for the optimization step and subsequent correspondence calculation. The learned template \mathbf{S}_t contains 6,890 vertices while the unseen template \mathbf{S}_t' contains slightly over 176,000 vertices. A visualization is given in fig. 6. The visual results is consistent with table 2, with little difference in error rate between using a learned template and a high-resolution unseen template.

Results from table 2 suggest that h is able to transfer the learned predictions for $\hat{\theta}_g$ when g operates on a different template. As a result, the meta decoder adapts wells operating on a template that is different from the one used in training and takes advantage of the higher resolution of the new template to produce more accurate correspondences. The high adaptability of the Meta Decoder Network suggests the possibility of variations of the model that use

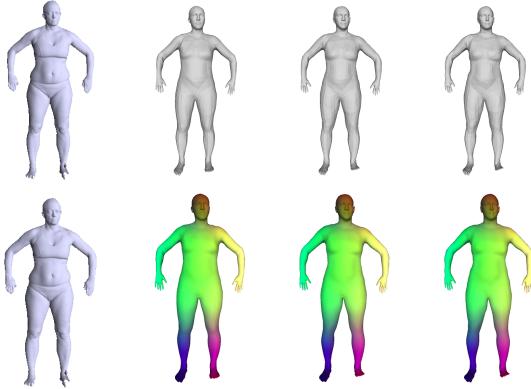


Fig. 6. Comparison of decoding from the high-resolution unseen template (top) and learned template (bottom) at 200, 1000, 3000 iterations. The image on the left of both rows is the query shape.

a dynamic template in addition to the dynamic decoder, which we desire to research in future works.

6.5 Training and Testing Speeds

Method	Time per training iteration(s)
3D-CODED w/. Learned 3D Translation Template [7]	0.1964
Meta Deformation Network (Ours)	0.1259

Table 3. Comparison of training speeds. Experiment ran with a GTX 1080 Ti and 6-core Intel CPU. Batch size was set to 32 on both runs.

Method	Time per pairwise corr. (s)
3D-CODED w/. Learned 3D Translation Template [7]	286.42
Meta Deformation Network (Ours)	273.66

Table 4. Comparison of testing speeds. Experiment ran with a GTX 1080 Ti and 6-core Intel CPU. Batch size was set to 1 on both runs.

We compare the training and testing speeds of the Meta Deformation Network against those of 3D-CODED, which has the same encoder structure but uses an LVC decoder. We train both models on the same Extended SURREAL dataset created by [7].

We find that the Meta Deformation Network took 35.8% less time to carry out an iteration of training (table 3) with a batch size set of 32, yet it yields improved accuracy (table 1) on the FAUST Inter challenge than did the model from [5] that also had a translation-based learnable template.

For testing, we evaluate the time a model takes to compute the correspondence for one pair of 3D scans from FAUST Inter. In testing, our model took 4.45% less time to finish. Here the difference was less dramatic than in training, which we think was due to the under-utilization of the computer’s hardware in both methods due to the the latent vector optimization with a batch size of 1. See table 4.

7 Conclusion

Meta Deformation Network is the first implementation of a meta decoder used in a 3-D to 3-D task aim at solving correspondences. It has an encoder-meta-decoder design in which all parameters of the decoder are inferred from a given query shape during both training and testing to improve the deformation of the template. The meta decoder enjoys more adaptability when deforming the template into query shapes compared to the LVC decoder [5] with the same encoder. As suggested by testing the Meta Deformation Network on the FAUST-Inter challenge, the meta decoder achieves better correspondence accuracy with the added benefit of speedier training and inferencing.

References

1. Bertinetto, L., Henriques, J.F., Valmadre, J., Torr, P.H.S., Vedaldi, A.: Learning feed-forward one-shot learners (2016)
2. Bogo, F., Romero, J., Loper, M., Black, M.J.: FAUST: Dataset and evaluation for 3D mesh registration. In: Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR). IEEE, Piscataway, NJ, USA (Jun 2014)
3. Brabandere, B.D., Jia, X., Tuytelaars, T., Gool, L.V.: Dynamic filter networks (2016)
4. Chang, A.X., Funkhouser, T., Guibas, L., Hanrahan, P., Huang, Q., Li, Z., Savarese, S., Savva, M., Song, S., Su, H., Xiao, J., Yi, L., Yu, F.: Shapenet: An information-rich 3d model repository (2015)
5. Deprelle, T., Groueix, T., Fisher, M., Kim, V.G., Russell, B.C., Aubry, M.: Learning elementary structures for 3d shape generation and matching (2019)
6. Fan, H., Su, H., Guibas, L.: A point set generation network for 3d object reconstruction from a single image (2016)
7. Groueix, T., Fisher, M., Kim, V.G., Russell, B.C., Aubry, M.: 3d-coded : 3d correspondences by deep deformation (2018)
8. Groueix, T., Fisher, M., Kim, V.G., Russell, B.C., Aubry, M.: Atlasnet: A papier-mch approach to learning 3d surface generation (2018)
9. Klein, B., Wolf, L., Afek, Y.: A dynamic convolutional layer for short rangeweather prediction. In: 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 4840–4848 (June 2015). <https://doi.org/10.1109/CVPR.2015.7299117>
10. Littwin, G., Wolf, L.: Deep meta functionals for shape representation (2019)
11. Loper, M., Mahmood, N., Romero, J., Pons-Moll, G., Black, M.J.: SMPL: A skinned multi-person linear model. ACM Trans. Graphics (Proc. SIGGRAPH Asia) **34**(6), 248:1–248:16 (Oct 2015)
12. Mitchell, E., Engin, S., Isler, V., Lee, D.D.: Higher-order function networks for learning composable 3d object representations (2019)
13. Ovsjanikov, M., Ben-Chen, M., Solomon, J., Butscher, A., Guibas, L.: Functional maps: A flexible representation of maps between shapes. ACM Transactions on Graphics - TOG **31** (07 2012). <https://doi.org/10.1145/2185520.2185526>
14. Park, J.J., Florence, P., Straub, J., Newcombe, R., Lovegrove, S.: Deepsdf: Learning continuous signed distance functions for shape representation (2019)
15. Qi, C.R., Su, H., Mo, K., Guibas, L.J.: Pointnet: Deep learning on point sets for 3d classification and segmentation (2016)
16. Richter, S.R., Roth, S.: Matryoshka networks: Predicting 3d geometry via nested shape layers (2018)
17. Riegler, G., Schulter, S., Rther, M., Bischof, H.: Conditioned regression models for non-blind single image super-resolution. In: 2015 IEEE International Conference on Computer Vision (ICCV). pp. 522–530 (Dec 2015). <https://doi.org/10.1109/ICCV.2015.67>
18. Varol, G., Romero, J., Martin, X., Mahmood, N., Black, M.J., Laptev, I., Schmid, C.: Learning from synthetic humans. In: CVPR (2017)
19. Yang, Y., Feng, C., Shen, Y., Tian, D.: Foldingnet: Point cloud auto-encoder via deep grid deformation (2017)
20. Zeng, W., Karaoglu, S., Gevers, T.: Inferring point clouds from single monocular images by depth intermedation (2018)

21. Zuffi, S., Black, M.J.: The stitched puppet: A graphical model of 3d human shape and pose. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2015)