

3D-CODED : 3D Correspondences by Deep Deformation

Thibault Groueix¹, Matthew Fisher², Vladimir G. Kim², Bryan C. Russell²,
and Mathieu Aubry¹

¹ LIGM (UMR 8049), École des Ponts, UPE

² Adobe Research

<http://imagine.enpc.fr/~groueixt/3D-CODED/>

Abstract. We present a new deep learning approach for matching deformable shapes by introducing *Shape Deformation Networks* which jointly encode 3D shapes and correspondences. This is achieved by factoring the surface representation into (i) a template, that parameterizes the surface, and (ii) a learnt global feature vector that parameterizes the transformation of the template into the input surface. By predicting this feature for a new shape, we implicitly predict correspondences between this shape and the template. We show that these correspondences can be improved by an additional step which improves the shape feature by minimizing the Chamfer distance between the input and transformed template. We demonstrate that our simple approach improves on state-of-the-art results on the difficult FAUST-inter challenge, with an average correspondence error of 2.88cm. We show, on the TOSCA dataset, that our method is robust to many types of perturbations, and generalizes to non-human shapes. This robustness allows it to perform well on real unclean, meshes from the the SCAPE dataset.

Keywords: 3D deep learning, computational geometry, shape matching

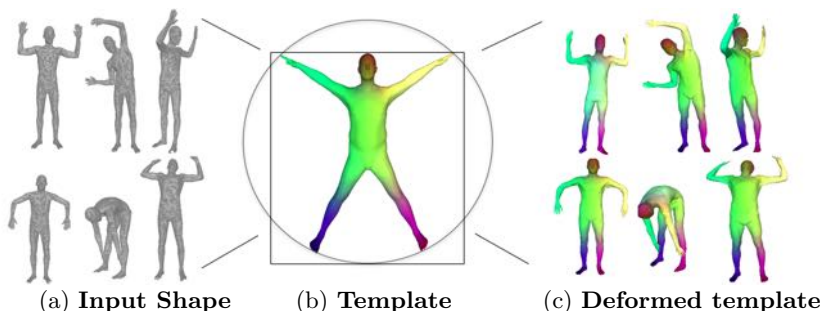


Fig. 1: Our approach predicts shape correspondences by learning a consistent mesh parameterization with a shared template. Colors show correspondences.

1 Introduction

There is a growing demand for techniques that make use of the large amount of 3D content generated by modern sensor technology. An essential task is to establish reliable 3D shape correspondences between scans from raw sensor data or between scans and a template 3D shape. This process is challenging due to low sensor resolution and high sensor noise, especially for articulated shapes, such as humans or animals, that exhibit significant non-rigid deformations and shape variations.

Traditional approaches to estimating shape correspondences for articulated objects typically rely on intrinsic surface analysis either optimizing for an isometric map or leveraging intrinsic point descriptors [43]. To improve correspondence quality, these methods have been extended to take advantage of category-specific data priors [9]. Effective human-specific templates and registration techniques have been developed over the last decade [49], but these methods require significant effort and domain-specific knowledge to design the parametric deformable template, create an objective function that ensures alignment of salient regions and is not prone to being stuck in local minima, and develop an optimization strategy that effectively combines a global search for a good heuristic initialization and a local refinement procedure.

In this work, we propose *Shape Deformation Networks*, a comprehensive, all-in-one solution to template-driven shape matching. A Shape Deformation Network learns to deform a template shape to align with an input observed shape. Given two input shapes, we align the template to both inputs and obtain the final map between the inputs by reading off the correspondences from the template.

We train our Shape Deformation Network as part of an encoder-decoder architecture, which jointly learns an encoder network that takes a target shape as input and generates a global feature representation, and a decoder Shape Deformation Network that takes as input the global feature and deform the template into the target shape. At test time, we improve our template-input shape alignment by optimizing locally the Chamfer distance between target and generated shape over the global feature representation which is passed in as input to the Shape Deformation Network. Critical to the success of our Shape Deformation Network is the ability to learn to deform a template shape to targets with varied appearances and articulation. We achieve this ability by training our network on a very large corpus of shapes.

In contrast to previous work [49], our method does not require a manually designed deformable template; the deformation parameters and degrees of freedom are implicitly learned by the encoder. Furthermore, while our network can take advantage of known correspondences between the template and the example shapes, which are typically available when they have been generated using some parametric model [6,45], we show it can also be trained without correspondence supervision. This ability allows the network to learn from a large collection of shapes lacking explicit correspondences.

We demonstrate that with sufficient training data this simple approach achieves state-of-the-art results and outperforms techniques that require complex multi-term objective functions instead of the simple reconstruction loss used by our method.

2 Related work

Registration of non-rigid geometries with pose and shape variations is a long standing problem with extensive prior work. We first provide a brief overview of generic correspondence techniques. We then focus on category specific and template matching methods developed for human bodies, which are more closely related to our approach. Finally, we present an overview of deep learning approaches that have been developed for shape matching and more generally for working with 3D data.

Generic shape matching. To estimate correspondence between articulated objects, it is common to assume that their intrinsic structure (e.g., geodesic distances) remains relatively consistent across all poses [29]. Finding point-to-point correspondences that minimize metric distortion is a non-convex optimization problem, referred to as generalized multi-dimensional scaling [11]. This optimization is typically sensitive to an initial guess [10], and thus existing techniques rely on local feature point descriptors such as HKS [43] and WKS [5], and use hierarchical optimization strategies [37,14]. Some relaxations of this problem have been proposed such as: formulating it as Markov random field and using linear programming relaxation [13], optimizing for soft correspondence [40,21,41], restricting correspondence space to conformal maps [23,22], heat kernel maps [32], and aligning functional bases [33].

While these techniques are powerful generic tools, some common categories, such as humans, can benefit from a plethora of existing data [6] to leverage stronger class-specific priors.

Template-based shape matching. A natural way to leverage class-specific knowledge is through the explicit use of a shape model. While such template-based techniques provide the best correspondence results they require a careful parameterization of the template, which took more than a decade of research to reach the current level of maturity [11,23,25,49]. For all of these techniques, fitting this representation to an input 3D shape requires also designing an objective function that is typically non-convex and involves multiple terms to guide the optimization to the right global minima. In contrast, our method only relies on a single template 3D mesh and surface reconstruction loss. It leverages a neural network to learn how to parameterize the human body while optimizing for the best reconstruction.

Deep learning for shape matching. Another way to leverage priors and training data is to learn better point-wise shape descriptors using human models with ground truth correspondence. Several neural network based methods have recently been developed to this end to analyze meshes [36,27,7,31] or depth maps [46]. One can further improve these results by leveraging global context, for example, by estimating an inter-surface functional map [24]. These methods still rely on hand-crafted point-wise descriptors [44] as input and use neural networks to improve results. The resulting functional maps only align basis functions and additional optimization is required to extract consistent point-to-point correspondences [33]. One would also need to optimize for template deformation to use these matching techniques for surface reconstruction. In contrast our method does not rely on hand-crafted features (it only takes point coordinates as input) and implicitly learns a human body representation. It also directly outputs a template deformation.

Deep Learning for 3D data. Following the success of deep learning approaches for image analysis, many techniques have been developed for processing 3D data, going beyond local descriptor learning to improve classification, segmentation, and reconstruction tasks. Existing networks operate on various shape representations, such as volumetric grids [17,47], point clouds [34,35,16], geometry images [39,38], seamlessly parameterized surfaces [26], by aligning a shape to a grid via distance-preserving maps [15], by folding a surface [48] or by predicting chart representations [18]. We build on these works in several ways. First, we process the point clouds representing the input shapes using an architecture similar to [34]. Second, similar to [38], we learn a surface representation. However, we do not explicitly encode correspondences in the output of a convolution network, but implicitly learn them by optimizing for parameters of the generation network as we optimize for reconstruction losses.

3 Method

Our goal is, given a reference shape \mathcal{S}_r and a target shape \mathcal{S}_t , to return a set of point correspondences \mathcal{C} between the shapes. We do so using two key ideas. First, we learn to predict a transformation between the shapes instead of directly learning the correspondences. This transformation, from 3D to 3D can indeed be represented by a neural network more easily than the association between variable and large number of points. The second idea is to learn transformations only from one template \mathcal{A} to any shape. Indeed, the large variety of possible poses of humans makes considering all pairs of possible poses intractable during training. We instead decouple the correspondence problem into finding two sets of correspondences to a common template shape. We can then form our final correspondences between the input shapes via indexing through the template shape. An added benefit is during training we simply need to vary the pose for a single shape and use the known correspondences to the template shape as the supervisory signal.

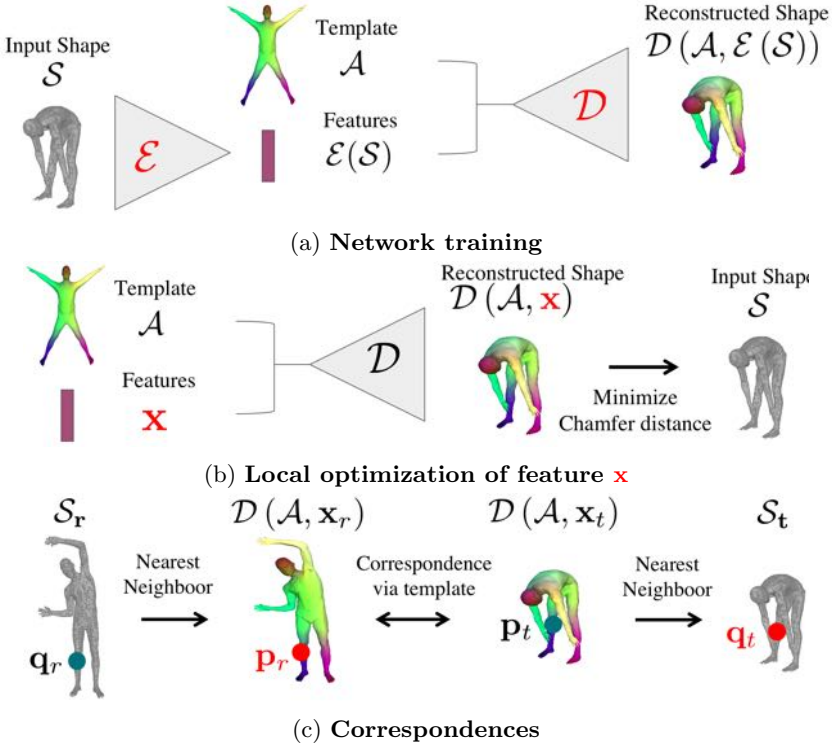


Fig. 2: **Method overview.** (a) A feed-forward pass in our autoencoder encodes input point cloud \mathcal{S} to latent code $\mathcal{E}(\mathcal{S})$ and reconstruct \mathcal{S} using $\mathcal{E}(\mathcal{S})$ to deform the template \mathcal{A} . (b) We refine the reconstruction $\mathcal{D}(\mathcal{A}, \mathcal{E}(\mathcal{S}))$ by performing a regression step over the latent variable \mathbf{x} , minimizing the Chamfer distance between $\mathcal{D}(\mathcal{A}, \mathbf{x})$ and \mathcal{S} . (c) Finally, given two point clouds \mathcal{S}_r and \mathcal{S}_t , to match a point \mathbf{q}_r on \mathcal{S}_r to a point \mathbf{q}_t on \mathcal{S}_t , we look for the nearest neighbor \mathbf{p}_r of \mathbf{q}_r in $\mathcal{D}(\mathcal{A}, \mathbf{x}_r)$, which is by design in correspondence with \mathbf{p}_t ; and look for the nearest neighbor \mathbf{q}_t of \mathbf{p}_t on \mathcal{S}_t . Red indicates what is being optimised.

Our approach has three main steps which are visualized figure 2. First, a feed-forward pass through our encoder network generates an initial global shape descriptor (Section 3.1). Second, we use gradient descent through our decoder Shape Deformation Network to refine this shape descriptor to improve the reconstruction quality (Section 3.2). We can then use the template to match points between any two input shapes (Section 3.3).

3.1 Learning 3D shape reconstruction by template deformation

To put an input shape \mathcal{S} in correspondence with a template \mathcal{A} , our first goal is to design a neural network that will take \mathcal{S} as input and predict transformation

parameters. We do so by training an encoder-decoder architecture. The encoder \mathcal{E}_ϕ defined by its parameters ϕ takes as input 3D points, and is a simplified version of the network presented in [34]. It applies to each input 3D point coordinate a multi-layer perceptron with hidden feature size of 64, 128 and 1024, then maxpooling over the resulting features over all points followed by a linear layer, leading to feature of size 1024 $\mathcal{E}_\phi(\mathcal{S})$. This feature, together with the 3D coordinates of a point on the template $\mathbf{p} \in \mathcal{A}$, are taken as input to the decoder \mathcal{D}_θ with parameters θ , which is trained to predict the position \mathbf{q} of the corresponding point in the input shape. This decoder Shape Deformation Network is a multi-layer perceptron with hidden layers of size 1024, 512, 254 and 128, followed by a hyperbolic tangent. This architecture maps any points from the template domain to the reconstructed surface. By sampling the template more or less densely, we can generate an arbitrary number of output points by sequentially applying the decoder over sampled template points.

This encoder-decoder architecture is trained end-to-end. We assume that we are given as input a training set of N shapes $\{\mathcal{S}^{(i)}\}_{i=1}^N$ with each shape having a set of P vertices $\{\mathbf{q}_j\}_{j=1}^P$. We consider two training scenarios: one where the correspondences between the template and the training shapes are known (supervised case) and one where they are unknown (unsupervised case). Supervision is typically available if the training shapes are generated by deforming a parametrized template, but real object scans are typically obtained without correspondences.

Supervised loss. In the supervised case, we assume that for each point \mathbf{q}_j on a training shape we know the correspondence $\mathbf{p}_j \leftrightarrow \mathbf{q}_j$ to a point $\mathbf{p}_j \in \mathcal{A}$ on the template \mathcal{A} . Given these training correspondences, we learn the encoder \mathcal{E}_ϕ and decoder \mathcal{D}_θ by simply optimizing the following reconstruction losses,

$$\mathcal{L}^{\text{sup}}(\theta, \phi) = \sum_{i=1}^N \sum_{j=1}^P |\mathcal{D}_\theta(\mathbf{p}_j; \mathcal{E}_\phi(\mathcal{S}^{(i)})) - \mathbf{q}_j^{(i)}|^2 \quad (1)$$

where the sums are over all P vertices of all N example shapes.

Unsupervised loss. In the case where correspondences between the exemplar shapes and the template are not available, we also optimize the reconstructions, but also regularize the deformations toward isometries. For reconstruction, we use the Chamfer distance \mathcal{L}^{CD} between the inputs \mathcal{S}_i and reconstructed point clouds $\mathcal{D}_\theta(\mathcal{A}; \mathcal{E}_\phi(\mathcal{S}^{(i)}))$. For regularization, we use two different terms. The first term \mathcal{L}^{Lap} encourages the Laplacian operator defined on the template and the deformed template to be the same (which is the case for isometric deformations of the surface). The second term $\mathcal{L}^{\text{edges}}$ encourages the ratio between edges length in the template and its deformed version to be close to 1. More details on these different losses are given in supplementary material. The final loss we optimize is:

$$\mathcal{L}^{\text{unsup}} = \mathcal{L}^{\text{CD}} + \lambda_{\text{Lap}} \mathcal{L}^{\text{Lap}} + \lambda_{\text{edges}} \mathcal{L}^{\text{edges}} \quad (2)$$

where λ_{Lap} and λ_{edges} control the influence of regularizations against the data term \mathcal{L}^{CD} . They are both set to $5 \cdot 10^{-3}$ in our experiments.

We optimize the loss using the Adam solver, with a learning rate of 10^{-3} for 25 epochs then 10^{-4} for 2 epochs, batches of 32 shapes, and 6890 points per shape.

One interesting aspect of our approach is that it learns jointly a parameterization of the input shapes via the decoder and to predict the parameters $\mathcal{E}_\phi(\mathcal{S})$ for this parameterization via the encoder. However, the predicted parameters $\mathcal{E}_\phi(\mathcal{S})$ for an input shape \mathcal{S} are not necessarily optimal, because of the limited power of the encoder. Optimizing these parameters turns out to be important for the final results, and is the focus of the second step of our pipeline.

3.2 Optimizing shape reconstruction

We now assume that we are given a shape \mathcal{S} as well as learned weights for the encoder \mathcal{E}_ϕ and decoder \mathcal{D}_θ networks. To find correspondences between the template shape and the input shape, we will use a nearest neighbor search to find correspondences between that input shape and its reconstruction. For this step to work, we need the reconstruction to be accurate. The reconstruction given by the parameters $\mathcal{E}_\phi(\mathcal{S})$ is only approximate and can be improved. Since we do not know correspondences between the input and the generated shape, we cannot minimize the loss given in equation (1), which requires correspondences. Instead, we minimize with respect to the global feature \mathbf{x} the Chamfer distance between the reconstructed shape and the input:

$$\mathcal{L}^{CD}(\mathbf{x}; \mathcal{S}) = \sum_{\mathbf{p} \in \mathcal{A}} \min_{\mathbf{q} \in \mathcal{S}} |\mathcal{D}_\theta(\mathbf{p}; \mathbf{x}) - \mathbf{q}|^2 + \sum_{\mathbf{q} \in \mathcal{S}} \min_{\mathbf{p} \in \mathcal{A}} |\mathcal{D}_\theta(\mathbf{p}; \mathbf{x}) - \mathbf{q}|^2. \quad (3)$$

Starting from the parameters predicted by our first step $\mathbf{x} = \mathcal{E}_\phi(\mathcal{S})$, we optimize this loss using the Adam solver for 3,000 iterations with learning rate $5 * 10^{-4}$. Note that the good initialization given by our first step is key since Equation(3) corresponds to a highly non-convex problem, as shown in Figure 6.

3.3 Finding 3D shape correspondences

To recover correspondences between two 3D shapes \mathcal{S}_r and \mathcal{S}_t , we first compute the parameters to deform the template to these shapes, \mathbf{x}_r and \mathbf{x}_t , using the two steps outlined in section 3.1 and 3.2. Next, given a 3D point \mathbf{q}_r on the reference shape \mathcal{S}_r , we first find the point \mathbf{p} on the template \mathcal{A} such that its transformation with parameters \mathbf{x}_r , $\mathcal{D}_\theta(\mathbf{p}; \mathbf{x}_r)$ is closest to \mathbf{q}_r . Finally we find the 3D point \mathbf{q}_t on the target shape \mathcal{S}_t that is the closest to the transformation of \mathbf{p} with parameters \mathbf{x}_t , $\mathcal{D}_\theta(\mathbf{p}; \mathbf{x}_t)$. Our algorithm is summarized in Algorithm 1 and illustrated in Figure 2.

Algorithm 1: Algorithm for finding 3D shape correspondences

Input : Reference shape \mathcal{S}_r and target shape \mathcal{S}_t

Output: Set of 3D point correspondences \mathcal{C}

```

1 #Regression steps over latent code to find best reconstruction of  $\mathcal{S}_r$  and  $\mathcal{S}_t$ 
2  $\mathbf{x}_r \leftarrow \arg \min_{\mathbf{x}} \mathcal{L}^{\text{CD}}(\mathbf{x}; \mathcal{S}_r)$  #detailed in equation (3)
3  $\mathbf{x}_t \leftarrow \arg \min_{\mathbf{x}} \mathcal{L}^{\text{CD}}(\mathbf{x}; \mathcal{S}_t)$  #detailed in equation (3)
4  $\mathcal{C} \leftarrow \emptyset$ 
5 # Matching of  $\mathbf{q}_r \in \mathcal{S}_r$  to  $\mathbf{q}_t \in \mathcal{S}_t$ 
6 foreach  $\mathbf{q}_r \in \mathcal{S}_r$  do
7    $\mathbf{p} \leftarrow \arg \min_{\mathbf{p}' \in \mathcal{A}} |\mathcal{D}_\theta(\mathbf{p}'; \mathbf{x}_r) - \mathbf{q}_r|^2$ 
8    $\mathbf{q}_t \leftarrow \arg \min_{\mathbf{q}' \in \mathcal{S}_t} |\mathcal{D}_\theta(\mathbf{p}; \mathbf{x}_t) - \mathbf{q}'|^2$ 
9    $\mathcal{C} \leftarrow \mathcal{C} \cup \{(\mathbf{q}_r, \mathbf{q}_t)\}$ 
10 end
11 return  $\mathcal{C}$ 

```

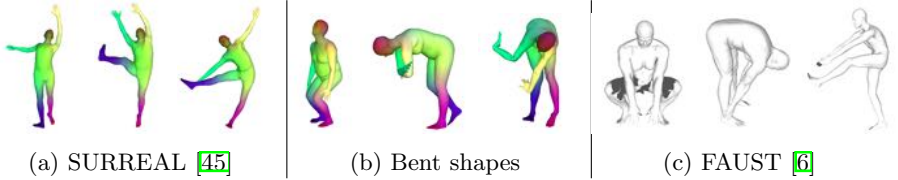


Fig. 3: Examples of the different datasets used in the paper.

4 Results

4.1 Datasets

Synthetic training data. To train our algorithm, we require a large set of shapes. We thus rely on synthetic data for training our model.

For human shapes, we use SMPL [6], a state-of-the-art generative model for synthetic humans. To obtain realistic human body shape and poses from the SMPL model, we sampled 2.10^5 parameters estimated in the SURREAL dataset [45]. One limitation of the SURREAL dataset is it does not include any humans bent over. Without adapted training data, our algorithm generalized poorly to these poses. To overcome this limitation, we generated an extension of the dataset. We first manually estimated 7 key-joint parameters (among 23 joints in the SMPL skeletons) to generate bent humans. We then sampled randomly the 7 parameters around these values, and used parameters from the SURREAL dataset for the other pose and body shape parameters. Note that not all meshes generated with this strategy are realistic as shown in figure 3. They however allow us to better cover the space of possible poses, and we added $3 \cdot 10^4$ shapes generated with this method to our dataset. Our final dataset thus has $2.3 \cdot 10^5$ human meshes with a large variety of realistic poses and body shapes.

For animal shapes, we use the SMAL [51] model, which provides the equivalent of SMPL for several animals. Recent papers estimate model parameters from

images, but no large-scale parameter set is yet available. For training we thus generated models from SMAL with random parameters (drawn from a Gaussian distribution of *ad-hoc* variance 0.2). This approach works for the 5 categories available in SMAL. In SMALR [50], Zuffi et al. showed that the SMAL model could be generalized to other animals using only an image dataset as input, demonstrating it on 17 additional categories. Note that since the templates for two animals are in correspondences, our method can be used to get inter-category correspondences for animals. We qualitatively demonstrate this on hippopotamus/horses in the appendix [19].

Testing data. We evaluate our algorithm on the FAUST [6], TOSCA [12] and SCAPE [4] datasets.

The FAUST dataset consists of 100 training and 200 testing scans of approximately 170,000 vertices. They may include noise and have holes, typically missing part of the feet. In this paper, we never used the training set, except for a single baseline experiment, and we focus on the test set. Two challenges are available, focusing on intra- and inter-subject correspondences. The error is the average Euclidean distance between the estimated projection and the ground-truth projection. We evaluated our method through the online server and are the best public results on the 'inter' challenge at the time of submission³.

The SCAPE [4] dataset has two sets of 71 meshes : the first set consists of real scans with holes and occlusions and the second set are registered meshes aligned to the first set. The poses are different from both our training dataset and FAUST.

TOSCA is a dataset produced by deforming 3 template meshes (human, dog, and horse). Each mesh is deformed into multiple poses, and might have various additional perturbations such as random holes in the surface, local and global scale variations, noise in vertex positions, varying sampling density, and changes in topology.

Shape normalization. To be processed and reconstructed by our network, the training and testing shapes must be normalized in a similar way. Since the vertical direction is usually known, we used synthetic shapes with approximately the same vertical axis. We also kept a fixed orientation around this vertical axis, and at test time selected the one out of 50 different orientations which leads to the smaller reconstruction error in term of Chamfer distance. Finally, we centered all meshes according to the center of their bounding box and, for the training data only, added a random translation in each direction sampled uniformly between -3cm and 3cm to increase robustness.

4.2 Experiments

In this part, we analyze the key components of our pipeline. More results are available in the appendix [19].

³ http://faust.is.tue.mpg.de/challenge/Inter-subject_challenge

Results on FAUST. The method presented above leads to the best results to date on the FAUST-inter dataset: 2.878 cm : **an improvement of 8% over state of the art**, 3.12cm for [49] and 4.82cm for [24]. Although it cannot take advantage of the fact that two meshes represent the same person, our method is also the second best performing (average error of 1.99 cm) on FAUST-intra challenge.

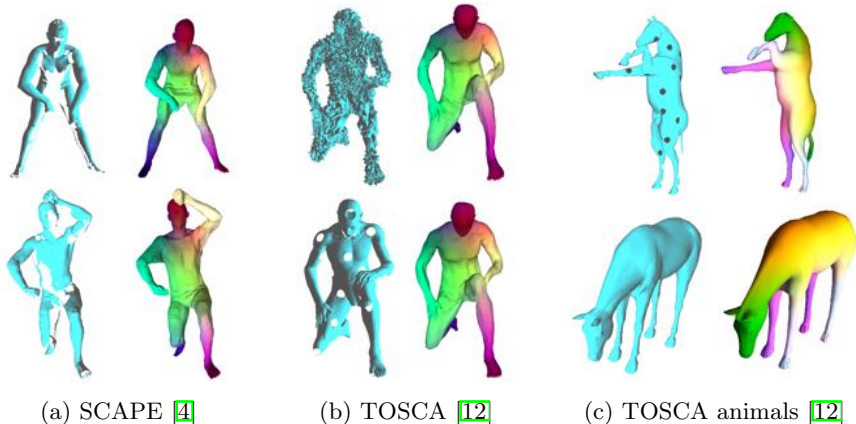


Fig. 4: **Other datasets.** Left images show the input, right images the reconstruction with colors showing correspondences. Our method works with real incomplete scans (a), strong synthetic perturbations (b), and on non-human shapes (c).

Results on SCAPE : real and partial data. The SCAPE dataset provides meshes aligned to real scans and includes poses different from our training dataset. When applying a network trained directly on our SMPL data, we obtain satisfying performance, namely 3.14cm average Euclidean error. Quantitative comparison of correspondence quality in terms of geodesic error are given in Fig [5]. We outperform all methods except for Deep Functional Maps [24]. SCAPE also allows evaluation on real partial scans. Quantitatively, the error on these partial meshes is 4.04cm, similar to the performance on the full meshes. Qualitative results are shown in Fig [4a].

Results on TOSCA : robustness to perturbations. The TOSCA dataset provides several versions of the same synthetic mesh with different perturbations. We found that our method, still trained only on SMPL or SMAL data, is robust to all perturbations (isometry, noise, shotnoise, holes, micro-holes, topology changes, and sampling), except scale, which can be trivially fixed by normalizing

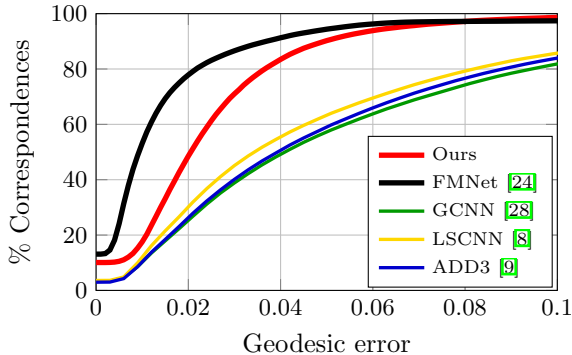


Fig. 5: Comparison with learning-based shape matching approaches on the SCAPE dataset. Our method is trained on synthetic data, FMNet was trained on FAUST data, and all other methods on SCAPE. We outperform all methods except FMNet even though our method was trained on a different dataset.

all meshes to have consistent surface area. Examples of representative qualitative results are shown Fig 4b and quantitative results are reported in appendix 19.

Reconstruction optimization. Because the nearest neighbors used in the matching step are sensitive to small errors in alignment, the second step of our pipeline which finds the optimal features for reconstruction, is crucial to obtain high quality results. This optimization however converges to a good optimum only if it is initialized with a reasonable reconstruction, as visualized in Figure 6. Since we optimize using Chamfer distance, and not correspondences, we also rely on the fact that the network was trained to generate humans in correspondence and we expect the optimized shape to still be meaningful.

Table 1 reports the associated quantitative results on FAUST-inter. We can see that: (i) optimizing the latent feature to minimize the Chamfer distance between input and output provides a strong boost; (ii) using a better (more uniform) sampling of the shapes when training our network provided a better initialization; (iii) using a high resolution sampling of the template ($\sim 200k$

Method	Faust error (cm)
Without regression	6.29
With regression	3.255
With regression + Regular Sampling	3.048
With regression + Regular Sampling + High-Res template	2.878

Table 1: **Importance of the reconstruction optimization step.** Optimizing the latent feature is key to our results. Regular point sampling for training and high resolution for the nearest neighbor step provide an additional boost.

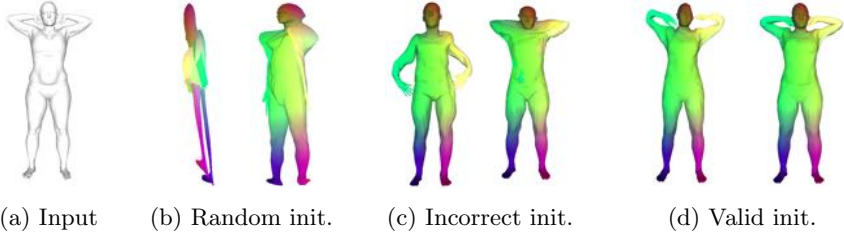


Fig. 6: **Reconstruction optimization.** The quality of the initialization (i.e. the first step of our algorithm) is crucial for the deformation optimization. For a given target shape (a) and for different initializations (left of (b), (c) and (d)) the figure shows the results of the optimization. If the initialization is random (b) or incorrect (c), the optimization converges to bad local minima. With a reasonable initialization (d) it converges to a shape very close to the target ((d), right).

vertices) for the nearest-neighbor step provide an additional small boost in performance.

Necessary amount of training data. Training on a large and representative dataset is also crucial for our method. To analyze the effect of training data, we ran our method without re-sampling FAUST points regularly and with a low resolution template for different training sets: FAUST training set, 2×10^5 SURREAL shapes, and 2.3×10^5 , 10^4 and 10^3 shapes from our augmented dataset. The quantitative results are reported Table 2 and qualitative results can be seen in Figure 7. The FAUST training set only include 10 different poses and is too small to train our network to generalize. Training on many synthetic shapes from the SURREAL dataset [45] helps overcome this generalization problem. However, if the synthetic dataset does not include any pose close to test poses (such as bent-over humans), the method will fail on these poses (4 test pairs of shapes out of 40). Augmenting the dataset as described in section 4.1 overcomes this

training data	Faust error (cm)
FAUST training set	18.22
non-augmented synthetic dataset 2×10^5 shapes	5.63
augmented synthetic data, 10^3 shapes	5.76
augmented synthetic data, 10^4 shapes	4.70
augmented synthetic data, 2.3×10^5 shapes	3.26

Table 2: **FAUST-inter results when training on different datasets.** Adding synthetic data reduce the error by a factor of 3, showing its importance. The difference in performance between the basic synthetic dataset and its augmented version is mostly due to failure on specific poses, as in Figure 3.

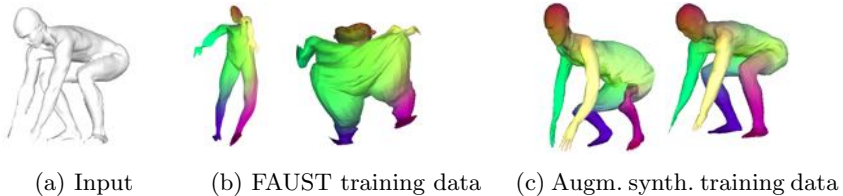


Fig. 7: **Importance of the training data.** For a given target shape (a) reconstructed shapes when the network is trained on FAUST training set (b) and on our augmented synthetic training set (c), before (left) and after (right) the optimization step.

limitation. As expected the performance decreases with the number of training shapes, respectively to 5.76cm and 4.70cm average error on FAUST-inter.

Unsupervised correspondences. We investigate whether our method could be trained without correspondence supervision. We started by simply using the reconstruction loss described in Equation (3). One could indeed expect that an optimal way to deform the template into training shapes would respect correspondences. However, we found that the resulting network did not respect correspondences between the template and the input shape, as visualized figure 8. However, these results improve with adequate regularization such as the one presented in Equation (2), encouraging regularity of the mapping between the template and the reconstruction. We trained such a network with the same training data as in the supervised case but **without any correspondence supervision** and obtained a 4.88cm of error on the FAUST-inter data, i.e. similar to Deep Functional Map [24] which had an error of 4.83 cm. This demonstrates that our method can be efficient even without correspondence supervision. Further details on regularization losses are given in the appendix 19.

Rotation invariance We handled rotation invariance by rotating the shape and selecting the orientation for which the reconstruction is optimal. As an alternative, we tried to learn a network directly invariant to rotations around the vertical axis. It turned out the performances were slightly worse on FAUST-

Loss	Faust error (cm)
Chamfer distance, eq. 3 (unsupervised)	8.727
Chamfer distance + Regularization, eq. 2 (unsupervised)	4.835
Correspondences, eq. 1 (supervised)	2.878

Table 3: Results with and without supervised correspondences. Adding regularization helps the network find a better local minimum in terms of correspondences.

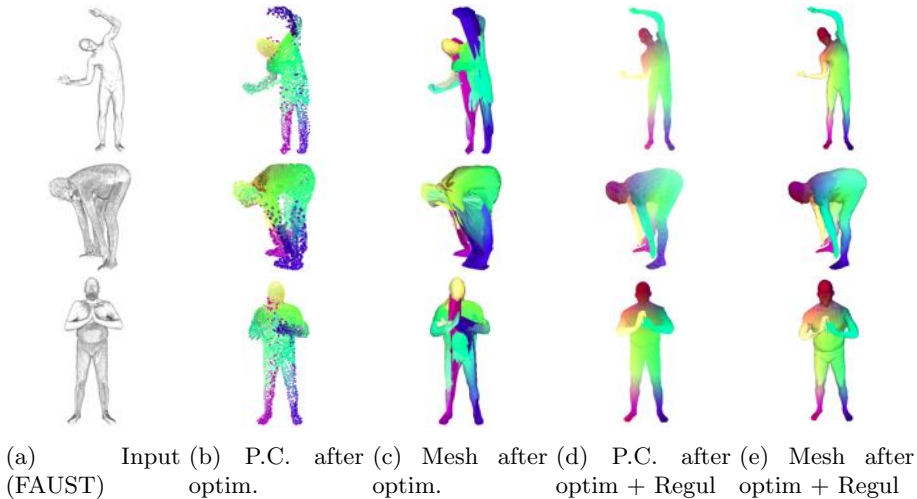


Fig. 8: **Unsupervised correspondences.** We visualize for different inputs (a), the point clouds (P.C.) predicted by our approach (b,d) and the corresponding meshes (c,e). Note that without regularization, because of the strong distortion, the meshes appear to barely match to the input, while the point clouds are reasonable. On the other hand surface regularization creates reasonable meshes.

inter (3.10cm), but still better than the state of the art. We believe this is due to the limited capacity of the network and should be tried with a larger network. However, interestingly, this rotation invariant network seems to have increased robustness and provided slightly better results on SCAPE.

5 Conclusion

We have demonstrated an encoder-decoder deep network architecture that can generate human shape correspondences competitive with state-of-the-art approaches and that uses only simple reconstruction and correspondence losses. Our key insight is to factor the problem into an encoder network that produces a global shape descriptor, and a decoder Shape Deformation Network that uses this global descriptor to map points on a template back to the original geometry. A straightforward regression step uses gradient descent through the Shape Deformation Network to significantly improve the final correspondence quality.

Acknowledgments. This work was partly supported by ANR project EnHerit ANR-17-CE23-0008, Labex Bézout, and gifts from Adobe to École des Ponts. We thank Gül Varol, Angjoo Kanazawa, and Renaud Marlet for fruitful discussions.

References

1. Allen, B., Curless, B., Popovic, Z.: Articulated body deformation from range scan data. SIGGRAPH (2002)
2. Allen, B., Curless, B., Popovic, Z.: The space of human body shapes: reconstruction and parameterization from range scans. SIGGRAPH (2003)
3. Allen, B., Curless, B., Popovic, Z.: Learning a correlated model of identity and pose-dependent body shape variation for real-time synthesis. Symposium on Computer Animation (2006)
4. Anguelov, D., Srinivasan, P., Koller, D., Thrun, S., Rodgers, J., Davis, J.: Scape: shape completion and animation of people. ACM transactions on graphics (TOG) **24**(3), 408–416 (2005)
5. Aubry, M., Schlickewei, U., Cremers, D.: The wave kernel signature: A quantum mechanical approach to shape analysis. IEEE International Conference on Computer Vision (ICCV) - Workshop on Dynamic Shape Capture and Analysis (4DMOD) (2011)
6. Bogo, F., Romero, J., Loper, M., Black, M.J.: FAUST: Dataset and evaluation for 3D mesh registration. Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) (Jun 2014)
7. Boscaini, D., Masci, J., Rodola, E., Bronstein, M.M.: Learning shape correspondence with anisotropic convolutional neural networks. NIPS (2016)
8. Boscaini, D., Masci, J., Melzi, S., Bronstein, M.M., Castellani, U., Vandergheynst, P.: Learning class-specific descriptors for deformable shapes using localized spectral convolutional networks. Computer Graphics Forum **34**(5), 13–23 (2015)
9. Boscaini, D., Masci, J., Rodolà, E., Bronstein, M.M., Cremers, D.: Anisotropic diffusion descriptors. Computer Graphics Forum **35**(2), 431–441 (2016)
10. Bronstein, A.M., Bronstein, M.M., Kimmel, R.: Efficient computation of isometry-invariant distances between surfaces. SIAM J. Scientific Computing (2006)
11. Bronstein, A.M., Bronstein, M.M., Kimmel, R.: Generalized multidimensional scaling: a framework for isometry-invariant partial surface matching. Proc. National Academy of Sciences (PNAS) (2006)
12. Bronstein, A.M., Bronstein, M.M., Kimmel, R.: Numerical geometry of non-rigid shapes. Springer Science & Business Media (2008)
13. Chen, Q., Koltun, V.: Robust nonrigid registration by convex optimization. International Conference on Computer Vision (ICCV) (2015)
14. D.Raviv, A.Dubrovina, R.Kimmel: Hierarchical framework for shape correspondence. Numerical Mathematics: Theory, Methods and Applications (2013)
15. Ezuz, D., Solomon, J., Kim, V.G., Ben-Chen, M.: Gwcnn: A metric alignment layer for deep shape analysis. SGP (2017)
16. Fan, H., Su, H., Guibas, L.: A point set generation network for 3D object reconstruction from a single image. Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017)
17. Girdhar, R., Fouhey, D., Rodriguez, M., Gupta, A.: Learning a predictable and generative vector representation for objects. Proceedings of European Conference on Computer Vision (ECCV) (2016)
18. Groueix, T., Fisher, M., Kim, V.G., Russell, B., Aubry, M.: AtlasNet: A Papier-Mâché Approach to Learning 3D Surface Generation. Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) (2018)
19. Groueix, T., Fisher, M., Kim, V.G., Russell, B., Aubry, M.: Supplementary material (appendix) for the paper <https://http://imagine.enpc.fr/~groueixt/3D-CODED/index.html> (2018)

20. Kanazawa, A., Tulsiani, S., Efros, A.A., Malik, J.: Learning category-specific mesh reconstruction from image collections. CoRR **abs/1803.07549** (2018)
21. Kim, V.G., Li, W., Mitra, N.J., DiVerdi, S., Funkhouser, T.: Exploring Collections of 3D Models using Fuzzy Correspondences. Transactions on Graphics (Proc. of SIGGRAPH) **31**(4) (2012)
22. Kim, V.G., Lipman, Y., Funkhouser, T.: Blended Intrinsic Maps. Transactions on Graphics (Proc. of SIGGRAPH) **30**(4) (2011)
23. Lipman, Y., Funkhouser, T.: Mobius voting for surface correspondence. ACM Transactions on Graphics (Proc. SIGGRAPH) **28**(3) (2009)
24. Litany, O., Remez, T., Rodola, E., Bronstein, A.M., Bronstein, M.M.: Deep functional maps: Structured prediction for dense shape correspondence. ICCV (2017)
25. Loper, M., Mahmood, N., Romero, J., Pons-Moll, G., Black, M.J.: Smpl: A skinned multi-person linear model. SIGGRAPH Asia (2015)
26. Maron, H., Galun, M., Aigerman, N., Trope, M., Dym, N., Yumer, E., Kim, V.G., Lipman, Y.: Convolutional neural networks on surfaces via seamless toric covers. SIGGRAPH (2017)
27. Masci, J., Boscaini, D., Bronstein, M.M., Vandergheynst, P.: Geodesic convolutional neural networks on riemannian manifolds. Proc. of the IEEE International Conference on Computer Vision (ICCV) Workshops pp. 37–45 (2015)
28. Masci, J., Boscaini, D., Bronstein, M.M., Vandergheynst, P.: Geodesic convolutional neural networks on riemannian manifolds. 3dRR (2015)
29. Mémoli, F., Sapiro, S.: A theoretical and computational framework for isometry invariant recognition of point cloud data. Foundations of Computational Mathematics (2005)
30. Meyer, M., Desbrun, M., Schr, P., Barr, A.: Discrete differential-geometry operators for triangulated 2-manifolds. Proceedings of Visualization and Mathematics **3** (11 2001)
31. Monti, F., Boscaini, D., Masci, J., Rodola, E., Svoboda, J., Bronstein, M.M.: Geometric deep learning on graphs and manifolds using mixture model cnns. CVPR (2017)
32. Ovsjanikov, M., Mériqot, Q., Mémoli, F., Guibas, L.: One point isometric matching with the heat kernel. Computer Graphics Forum (Proc. of SGP) (2010)
33. Ovsjanikov, M., Ben-Chen, M., Solomon, J., Butscher, A., Guibas, L.: Functional maps: A flexible representation of maps between shapes. ACM Trans. Graph. (2012)
34. Qi, C.R., Su, H., Mo, K., Guibas, L.J.: PointNet: Deep learning on point sets for 3D classification and segmentation. Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017)
35. Qi, C.R., Yi, L., Su, H., Guibas, L.J.: PointNet++: Deep hierarchical feature learning on point sets in a metric space. Advances in Neural Information Processing Systems (NIPS) (2017)
36. Rodola, E., Rota Bulo, S., Windheuser, T., Vestner, M., Cremers, D.: Dense non-rigid shape correspondence using random forests. CVPR (2014)
37. Sahillioğlu, Y., Yemez, Y.: Coarse-to-fine combinatorial matching for dense isometric shape correspondence. Computer Graphics Forum (2011)
38. Sinha, A., Unmesh, A., Huang, Q., Ramani, K.: Surfnet: Generating 3d shape surfaces using deep residual networks. Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017)
39. Sinha, A., Bai, J., Ramani, K.: Deep learning 3d shape surfaces using geometry images. Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)

40. Solomon, J., Nguyen, A., Butscher, A., Ben-Chen, M., Guibas, L.: Soft maps between surfaces. SGP (2012)
41. Solomon, J., Peyre, G., Kim, V.G., Sra, S.: Entropic metric alignment for correspondence problems. Transactions on Graphics (Proc. of SIGGRAPH) (2016)
42. Sorkine, O.: Differential representations for mesh processing. Comput. Graph. Forum **25**, 789–807 (12 2006)
43. Sun, J., Ovsjanikov, M., Guibas, L.: A concise and provably informative multi-scale signature-based on heat diffusion”. Computer Graphics Forum (Proc. of SGP) (2009)
44. Tombari, F., Salti, S., Stefano, L.D.: Unique signatures of histograms for local surface description. ECCV (2010)
45. Varol, G., Romero, J., Martin, X., Mahmood, N., Black, M.J., Laptev, I., Schmid, C.: Learning from synthetic humans. CVPR (2017)
46. Wei, L., Huang, Q., Ceylan, D., Vouga, E., Li, H.: Dense human body correspondences using convolutional networks. Computer Vision and Pattern Recognition (CVPR) (2016)
47. Wu, Z., Song, S., Khosla, A., Yu, F., Zhang, L., Tang, X., Xiao, J.: 3d shapenets: A deep representation for volumetric shapes. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition pp. 1912–1920 (2015)
48. Yang, Y., Feng, C., Shen, Y., Tian, D.: Foldingnet: Point cloud auto-encoder via deep grid deformation. CVPR (2018)
49. Zuffi, S., Black, M.J.: The stitched puppet: A graphical model of 3d human shape and pose. Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) (2015)
50. Zuffi, S., Kanazawa, A., Black, M.J.: Lions and tigers and bears: Capturing non-rigid, 3D, articulated shape from images. IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)
51. Zuffi, S., Kanazawa, A., Jacobs, D., Black, M.J.: 3D menagerie: Modeling the 3D shape and pose of animals. CVPR (2017)

6 Supplementary

6.1 Choice of template

The template is a critical element for our method. We experimented with three different templates: (i) a “FAUST” template associated with SMPL parameters fitted to a body in a neutral pose in the FAUST training set, (ii) a “zero” template corresponding to the “zero” shape of SMPL, and (iii) a “separated” template in which this “zero” shape is modified to have the legs better separated and the arms higher. Figure 9 shows the different templates, while table 4 shows quantitative results using the different templates. Interestingly, the best results were obtained with the more “natural” template, selected in the “FAUST” training dataset, rather than with the templates from simple SMPL parameters, where points from different body parts seem easier to separate.

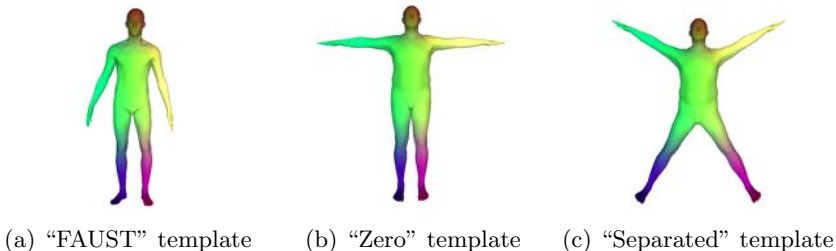


Fig.9: **Shapes for template study.** We evaluate three different template shapes used in our model.

template 0	Faust error (cm)
“FAUST” template	3.255
“Zero” template	3.385
“Separated” template	3.314

Table 4: **Comparison of different template shapes.** We compare different choices for the template shape shown in Figure 9. Notice that the neutral “FAUST” template performs best out of the three tested shapes.

6.2 Quantitative results for perturbations on TOSCA

We evaluate quantitatively the robustness of our method to perturbation on the TOSCA dataset. It consists of one horse shape with different added perturbations, namely noise, shotnoise, sampling, scale, local scale, topology, holes, microholes, and isometry. We report in 5 quantitative results for each perturbation

(with a gradual strength from 1 to 5) and show qualitative reconstruction with correspondences suggested by colors for each category with maximum strength in [10]. Surprisingly, adding noise can enhance the quantitative error.

Table 5: Quantitative results for perturbations on TOSCA for the horse category

Perturbation	Error (cm)	Perturbation	Error (cm)	Perturbation	Error (cm)
Noise	1 4.58	Scale	1 4.73	Holes	1 4.71
	2 3.87		2 4.78		2 4.71
	3 3.93		3 4.66		3 4.72
	4 3.67		4 4.62		4 4.69
	5 3.91		5 4.67		5 4.84
ShotNoise	1 4.66	Local scale	1 4.18	Microholes	1 4.71
	2 2.64		2 3.65		2 4.72
	3 3.03		3 3.62		3 4.82
	4 2.72		4 3.75		4 4.69
	5 3.00		5 3.56		5 3.53
Sampling	1 4.82	Topology	1 3.99	Isometry	1 4.72
	2 4.78		2 4.38		2 4.69
	3 4.61		3 4.37		3 4.79
	4 3.72		4 4.31		4 4.85
	5 9.93		5 7.53		5 4.74

6.3 Cross-category correspondances on animals

SMAL synthetic are in correspondences across categories. Hence the template for two different categories are in correspondences and our approach can be trivially extended to get correspondences for animals from different species. Qualitative evidence of this is show in Figure [11]

6.4 Regularization for the unsupervised case

We observe some distortion when the network is trained using the Chamfer distance alone. For example the left foot is propagated on left hand in Figure [8]. Even if this distortion is consistent across shapes, we hope that by regularizing the generator, the learned deformation on the template would respect the connectivity of the points of the templates. To achieve this, we tried two methods.

Ratio preservation Let (V, E) be the graph of the template and V^g the reconstructed vertices.

$$E_{ratio}(V^g) = \frac{1}{\#E} \cdot \sum_{i \sim j} \left\| \frac{V_i^g - V_j^g}{V_i - V_j} - 1 \right\|$$

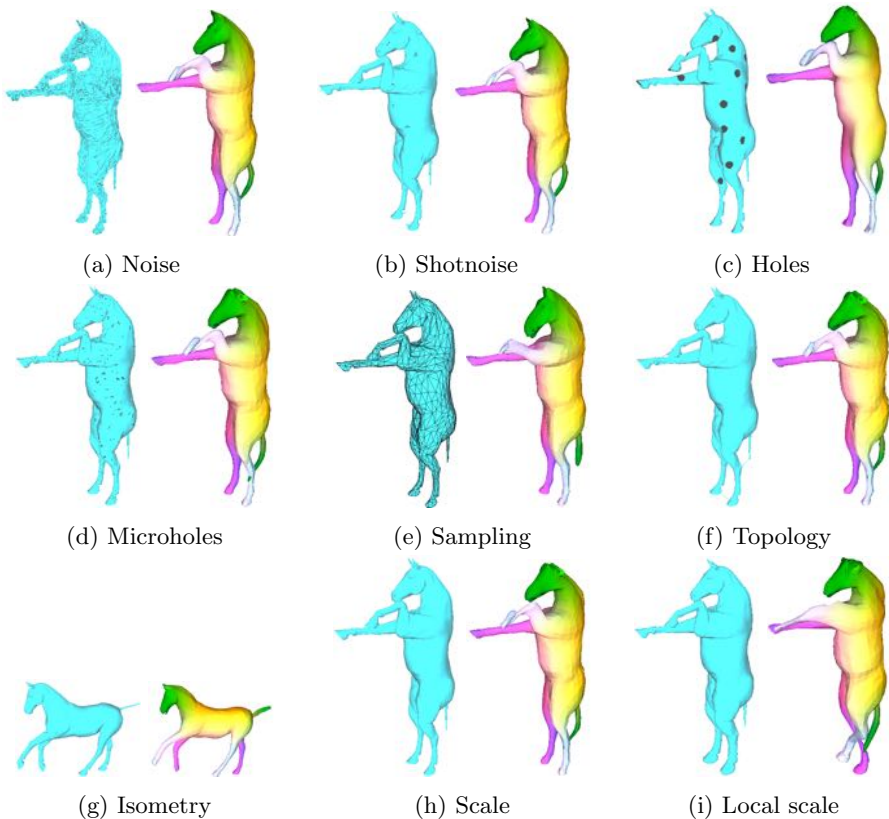


Fig. 10: Robustness to perturbations on TOSCA for the horse category. Correspondences are suggested by color. Notice the overall robustness to all perturbations, with small errors on the ears, tail or legs.

This enforces edges to keep the same length in the template and the generated mesh. We use $\lambda_{ratio} = 0.005$. For instance, if the length of an edge doubles the contribution to the loss is $\lambda_{ratio} \cdot 1 = 0.005$ which is equivalent (in terms of contribution to the loss function) to a error of placement of 7.1cm. In other words, in terms on loss for the network, it is equivalent to double an edge's length or to misplace a point by 3.2cm.

Laplacian regularization Similar to Kanazawa et. al. [20], we use the Laplacian regularization. The Laplacian matrix L is defined as :

$$L_{i,j} = \begin{cases} d_i & \text{if } i = j \\ -1 & \text{if } (i,j) \in E \\ 0 & \text{otherwise} \end{cases}$$

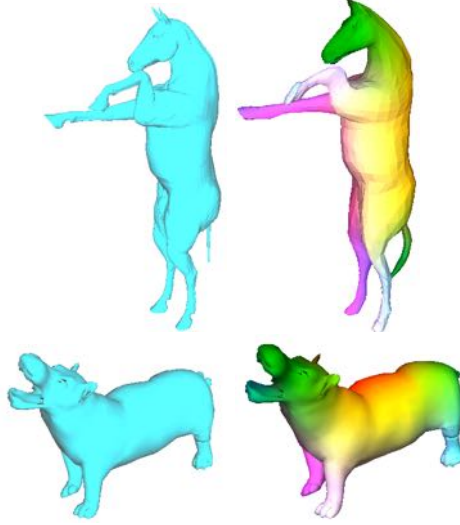


Fig. 11: **Inter-class correspondences on animals.** Correspondences are suggested by color.

$$LV_i = \sum_{i \sim j} V_i - V_j$$

$$LV_i = d_i \cdot (V_i - \frac{\sum_{i \sim j} V_j}{d_i})$$

This is an approximation of the following integral as explained in [42].

$$\lim_{\gamma \rightarrow 0} \frac{1}{|\gamma|} \int_{v \in \gamma} (v_i - v) dl(v) = -H(v_i) \cdot n_i$$

where:

- $H(v_i)$ is the mean curvature
- n_i is the surface normal

We follow [30] and use cotangent weights in the Laplacian to have better geometric discretization property.

$$L^c V_i = \frac{1}{\Omega_i} \sum_{i \sim j} \frac{1}{2} (\cot \alpha_{ij} + \cot \beta_{ij}) (V_i - V_j)$$

where :

- Ω_i is the size of the Voronoi cell of i
- α_{ij} and β_{ij} denote the two angles opposite of edge (i, j)

Our Laplacian loss is thus written :

$$E_{laplace}(V^g) = \mathbb{1}^t \cdot L^c \cdot (V^{template} - V^g)$$

We use $\lambda_{laplace} = 0.005$. In practice we notice that using Laplacian regularization constrain the network to keep sound surfaces. It may still suffer from error in symmetry and can still invert right and left, and front and back.

6.5 Asymetric Chamfer distance

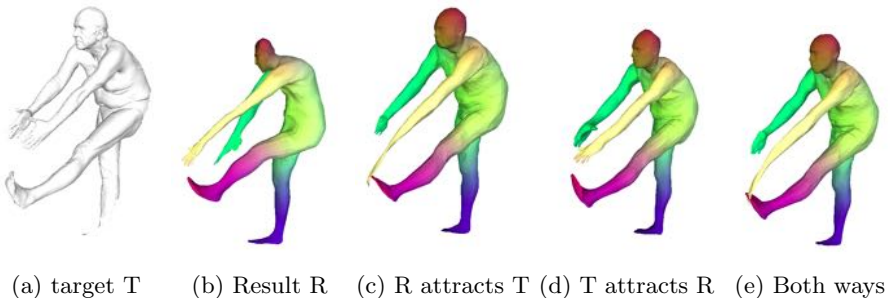


Fig. 12: **Asymmetric Chamfer loss in reconstruction optimization.** Given an input scan, with holes (a), our network outputs a reconstruction result (b), that can be improved by an optimization step. When the scan has holes, it is better to only consider a loss where the scan attracts the reconstruction (d), rather than using a loss where reconstruction attracts the scan (c), or the Chamfer distance where they attract each other (e).

Figure 12 illustrates that optimizing an asymmetric Chamfer distance can in some cases, especially when the 3D scans have holes, produce qualitatively better results.

6.6 Failure cases

Figure 13 shows the two main sources of error our algorithm faces. It can be an error in the nearest neighbor step in overlapping regions; here, a point is matched with the closest point in Euclidean distance but the match is very far in geodesic distance. This could be addressed by enforcing matches between the input mesh and its reconstruction in a way that takes into account the regularity of the matching. We leave this to future work.

The other source of error comes from failures in reconstruction: in such cases, the initial guess of the autoencoder is just too far away from the input, and the regression step fails.

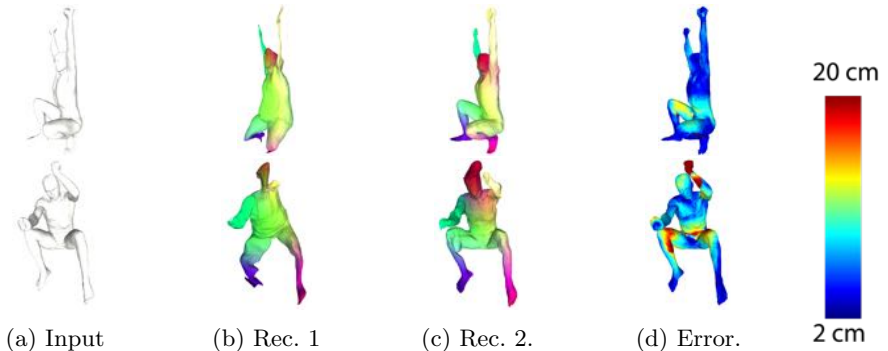


Fig. 13: **Error visualization** Given the input mesh (a), our autoencoder makes an initial reconstruction (b), optimized by a regression step (c). The average in centimeters over each vertex of (a), of the Euclidean distance between its projection and the ground truth, is reported (d). We use the jet colormap. Red vertices have an error higher than 10, blue ones lower than 2cm. The largest error are observed in places where the Euclidean distance is small, while the geodesic distance is high, such as touching skin (zoom in on the leg). In such region, the nearest neighbors step is match a vertex in mesh A in a distant (in terms of geodesic distance) vertex in mesh A’s reconstruction. High error can also come from bad reconstruction. See the head of the second example.