

# DispVoxNets: Non-Rigid Point Set Alignment with Supervised Learning Proxies\*

Soshi Shimada<sup>1,2</sup> Vladislav Golyanik<sup>3</sup> Edgar Tretschk<sup>3</sup> Didier Stricker<sup>1,2</sup> Christian Theobalt<sup>3</sup>

<sup>1</sup>University of Kaiserslautern

<sup>2</sup>DFKI

<sup>3</sup>MPI for Informatics, SIC

## Abstract

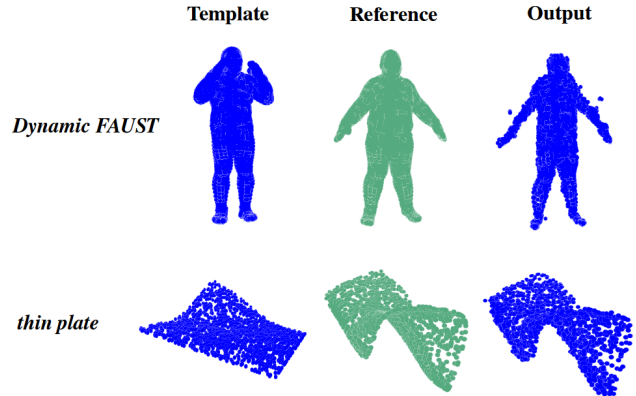
We introduce a supervised-learning framework for non-rigid point set alignment of a new kind — Displacements on Voxels Networks (*DispVoxNets*) — which abstracts away from the point set representation and regresses 3D displacement fields on regularly sampled proxy 3D voxel grids. Thanks to recently released collections of deformable objects with known intra-state correspondences, *DispVoxNets* learn a deformation model and further priors (e.g., weak point topology preservation) for different object categories such as cloths, human bodies and faces. *DispVoxNets* cope with large deformations, noise and clustered outliers more robustly than the state-of-the-art. At test time, our approach runs orders of magnitude faster than previous techniques. All properties of *DispVoxNets* are ascertained numerically and qualitatively in extensive experiments and comparisons to several previous methods.

## 1. Introduction

Point sets are raw shape representations which can implicitly encode surfaces and volumetric structures with inhomogeneous sampling densities. Many 3D vision techniques generate point sets which need to be subsequently aligned for various tasks such as shape recognition, appearance transfer and shape completion, among others.

The objective of non-rigid point set registration (NRPSR) is the recovery of a general displacement field aligning *template* and *reference* point sets, as well as correspondences between those. In contrast to rigid or affine alignment, where all template points transform according to a single shared transformation, in NRPSR, every point of the template has an individual transformation. Nonetheless, real structures do not evolve arbitrarily and often preserve the point topology.

\*supported by the ERC Consolidator Grant 4DReply (770784) and the BMBF project VIDETE (01IW18002).

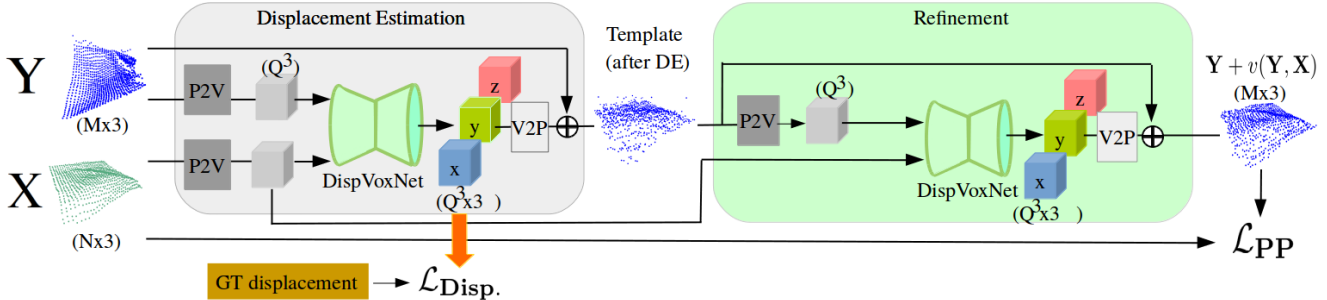


**Figure 1:** Alignment results of human body scans [5] and thin plate [17] with our DispVoxNets. In both cases, the template and reference differ by large non-linear deformations (articulated motion in the case of human body scans). To the best of our belief, DispVoxNet is the first non-rigid point set alignment approach which learns object-specific deformation models purely from data and does not rely on engineered priors.

### 1.1. Motivation and Contributions

On the one hand, existing general-purpose NRPSR techniques struggle to align point clouds differing by large non-linear deformations or articulations (e.g., significantly different facial expressions or body poses) and cause overregularisation, flattening and structure distortions [3, 9, 41, 31]. On the other hand, specialised methods exploit additional engineered (often class-specific) priors to align articulated and highly varying structures [40, 44, 13, 55, 57, 18]. In contrast, we are interested in a general-purpose method supporting large deformations and articulations (such as those shown in Fig. 1), which is robust to noise and clustered outliers and which can adapt to various object classes.

It is desirable but challenging to combine all these properties into a single technique. We address this difficult problem with supervised learning on collections of deformable objects with known intra-state correspondences. Even though deep learning is broadly and successfully applied to various tasks in computer vision, its applications to NRPSR have not been demonstrated in the literature so far (see Sec. 2). One of the reasons is the varying cardinalities



**Figure 2:** Overview of our approach. The objective is to non-rigidly align a template  $Y$  to a reference  $X$ . In the displacement estimation stage, we first convert the point sets to a voxel representation (P2V). DispVoxNets then regress per-voxel displacements that we apply to  $Y$  (V2P), see Table 1 for more details on the network architecture. The first stage is trained in a supervised manner with ground truth displacements using the *displacement loss* ( $\mathcal{L}_{\text{Disp}}$ ). The subsequent stage refines the displacements with the unsupervised *point projection loss* ( $\mathcal{L}_{\text{PP}}$ ). Trilinear weights are applied to the displacements for interpolation and are used to compute weighted gradients in the backward pass, see Fig. 3 and our supplement for more details on the trilinear interpolation.

of the inputs, which poses challenges in the network architecture design. Another reason is that sufficiently comprehensive collections of deformable shapes with large deformations, suitable for the training have just recently become available [5, 2, 17, 36].

To take advantage of the latter, our core idea is to associate deformation priors with point displacements and predict feasible category-specific deformations between input samples on an abstraction layer. At its core, our framework contains geometric proxies — deep convolutional encoder-decoders operating on voxel grids — which learn a class-specific deformation model. We call the proposed proxy component *Displacements on Voxels Network* (DispVoxNet). Our architecture contains two identical DispVoxNets, *i.e.*, one for global displacements (trained in a supervised manner) and one for local refinement (trained in an unsupervised manner).

The proposed DispVoxNets abstract away from low-level properties of point clouds such as point sampling density and ordering. They realise a uniform and computationally feasible lower-dimensional parametrisation of deformations which are eventually transferable to the template in its original resolution and configuration. At the same time, DispVoxNets handle inputs of arbitrary sizes. To bridge a possible discrepancy in resolution between the 3D voxel grids and the point clouds, we maintain a point-to-voxel affinity table and apply a super-resolution approach. Due to all these properties, DispVoxNet enables the level of generalisability of our architecture which is essential for a general-purpose NRPSR approach.

A schematic overview of the proposed architecture with DispVoxNets is given in Fig. 2. Our general-purpose NRPSR method can be trained for arbitrary types of deformable objects. During inference, no further assumptions about the input point sets except of the object class are made. All class-specific priors including the weak topology preserving constraint are learned directly from the data.

Whereas some methods model noise distributions to enable robustness to noise [41], we augment the training datasets by adding uniform noises and removing points uniformly at random. We do not rely on parametric models, pre-defined templates, landmarks or known segmentations (see Sec. 3).

In our experiments, DispVoxNets consistently outperform other tested approaches in scenarios with large deformations, noise and missing data. In total, we perform a study on four object types and show that DispVoxNets can efficiently *learn* class-specific priors (see Sec. 4).

## 2. Related Work

**Methods with Global Regularisers.** When correspondences between points are given, an optimal rigid transformation between the point sets can be estimated in a closed form [27]. Iterative Closest Point (ICP) alternates between estimating the correspondences based on the nearest neighbour rule and local transformations until convergence [3, 7]. ICP is a simple and widely-used point set alignment technique, with multiple policies available to improve its convergence properties, runtime and robustness to noise [23, 49, 21, 12]. In practice, conditions for a successful alignment with ICP (an accurate initialisation and no disturbing effects such as noise) are often not satisfied. Extensions of ICP for the non-rigid case employ thin splines for topology regularisation [9]<sup>1</sup> or Markov random fields linked by a non-linear potential function [24].

Probabilistic approaches operate with multiply-linked point associations. Robust Point Matching (RPM) with a thin-plate splines (TPS) deformation model [9] uses a combination of soft-assign [14] and deterministic annealing for non-rigid alignment. As the transformation approaches the optimal solution, the correspondences become more and more certain. In [62], point set alignment is formulated as a graph matching problem which aims to maximise the num-

<sup>1</sup>introduced by Chui *et al.* [9] as a baseline ICP modification

ber of matched edges in the graphs. The otherwise NP-hard combinatorial graph matching problem is approximated as a constrained optimisation problem with continuous variables by relaxation labeling. In Gaussian Mixture Model (GMM) Registration (GMR) [31], NRPSR is interpreted as minimising the distance between two mixtures of Gaussians with a TPS mapping. The method was shown to be more tolerant to outliers and more statistically robust than TPS-RPM and non-rigid ICP. Myronenko and Song [41] interpret NRPSR as fitting a GMM (template points) to data (reference points) and regularise displacement fields using motion coherence theory. The resulting Coherent Point Drift (CPD) was shown to handle noisy and partially overlapping data with unprecedented accuracy. Zhou *et al.* [63] investigate the advantages of the Student’s-t Mixture Model over a GMM. The method of Ma *et al.* [39] alternates between correspondence estimation with shape context descriptors and transformation estimation by minimising the  $\ell_2$  distance between two densities. Their method demands deformations to lie in the reproducing kernel Hilbert space.

Recently, physics-based alignment approaches were discovered [11, 15, 1, 30, 19]. Deng *et al.* [11] minimise a distance metric between Schrödinger distance transforms performed on the point sets. Their method has shown an improved recall, *i.e.*, the portion of correctly recovered correspondences. Ali *et al.* [1] align point sets as systems of particles with masses deforming under simulated gravitational forces. Gravitational simulation combined with smoothed particle hydrodynamics regularisation place this approach among the most resilient to large amounts of uniform noise in the data, and, at the same time, most computationally expensive techniques. In contrast, the proposed approach executes in just a few seconds and is robust to large amounts of noise due to our training policy with noise augmentation.

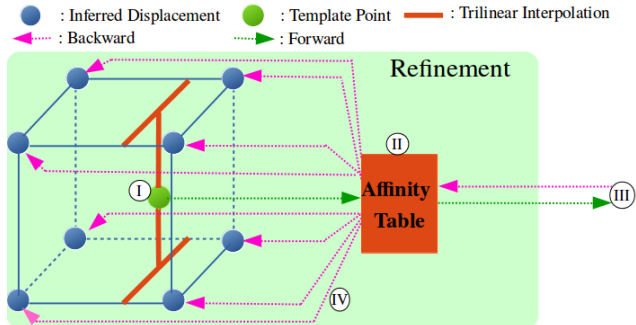
**Large Deformations and Articulations.** If point sets differ by large deformations and articulations, global topology regularisers of the methods discussed so far often overconstrain local deformations. Several extended versions of ICP address the case of articulated bodies with the primary applications to human hands and bodies [40, 44, 55]. Ge *et al.* [13] extend CPD with a local linear embedding which accounts for multiple non-coherent motions and local deformations. The method assumes a uniform sampling density and its accuracy promptly decays with an increasing level of noise. Some methods align articulated bodies with problem-specific segmented templates [16]. In contrast to all these techniques, our DispVoxNets can be trained for an arbitrary object class and are not restricted to a single template. Furthermore, our approach is resilient to sampling densities, large amounts of outliers and missing data. It grasps the intrinsic class-specific deformation model on multiple scales (global and localised deformations) directly

from data. A more in-depth overview and comparison of NRPSR methods can be found in [56, 64].

**Voxel-Based Methods.** Voxel-based methods have been an active research domain of 3D reconstruction over the past decades [51, 6, 4, 45, 33, 37, 50, 59]. Their core idea is to discretise the target volume of space in which the reconstruction is performed. With the renaissance of deep learning in the modern era [34, 20, 54, 29, 26], there have been multiple attempts to adapt voxel-based techniques to learning-based 3D reconstruction [8, 47]. These methods have been criticised for a high training complexity due to expensive 3D convolutions and discretisation artefacts due to a low grid resolution. In contrast to previous works, we use a voxel-based proxy to regress displacements instead of deformed shapes. Note that in many related tasks, deformations are parametrised by lower-resolution data structures such as deformation graphs [53, 42, 61]. To alleviate discretisation artefacts and enable superresolution of displacements, we apply point projection and trilinear interpolation.

**Learning Deformation Models.** Recently, the first supervised learning methods trained for a deformation model were proposed for monocular non-rigid 3D reconstruction [17, 46, 52, 58]. Their main idea is to train a deep neural network (DNN) for feasible deformation modes from collections of deforming objects with known correspondences between non-rigid states. Implicitly, a single shape at rest (a thin surface) is assumed which is deformed upon 2D observations. Next, several works include a free-form deformation component for monocular rigid 3D reconstruction with an object-class template [35, 28]. Hanocka *et al.* [25] align meshes in a voxelised representation with an unsupervised learning approach. They learn a shape-aware deformation prior from shape datasets and can handle incomplete data.

Our method is inspired by the concept of a learned deformation model. In NRPSR, both the reference and template can differ significantly from scenario to scenario, and we cannot assume a single template for all alignment problems. To account for different scenarios and inputs, we introduce a proxy voxel grid which abstracts away from the point cloud representation. We learn a deformation model for displacements instead of a space of feasible deformation modes for shapes. Thus, we are able to use the same data modality for training as in [17, 52, 58] and generalise to arbitrary point clouds for non-rigid alignment. Wang *et al.* [60] solve a related problem on a voxel grid: predicting object deformations under applied physical forces. Their network is trained in an adversarial manner with ground truth deformations conditioned upon the elastic properties of the material and applied forces. In contrast to 3D-PhysNet [60], we learn displacement fields on a voxel grid which is a more explicit representation of deformations intrinsic to NRPSR.



**Figure 3:** Overview of the forward and backward pass of trilinear interpolation on a voxel grid. Numbers (I), (II), (III) and (IV) indicate the sequence of steps performed in every training iteration, see Sec. 3.1 for more details.

### 3. The Proposed Approach

We propose a neural network-based method for NRPSR that takes a template and a reference point set and returns a displacement vector for each template point, see Fig. 2 for an overview. As described in Sec. 2, existing methods show lower accuracy in the presence of large deformations between point sets. We expect that neural-network-based methods are able to deal with such challenging cases since they learn class-specific priors implicitly during training. In NRPSR, the numbers of points in the template and reference are generally different. This inconsistency of the input dimensionality is problematic because we need to fix the number of neurons before training. To resolve this issue, we convert the point sets into a regular voxel-grid representation at the beginning of the pipeline, which makes our approach invariant with respect to the number and order of input points. Furthermore, due to the nature of convolutional layers, we expect a network with 3D convolutions to be robust to noises and outliers. Even though handling 3D voxel data is computationally demanding, modern hardware supports sufficiently fine-grained voxel grids which our approach relies on.

**Notations and Assumptions.** The inputs of our algorithm are two point sets: the reference  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T \in \mathbb{R}^{N \times D}$  and the template  $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_M)^T \in \mathbb{R}^{M \times D}$  which has to be non-rigidly matched to  $\mathbf{X}$ .  $N$  and  $M$  are the cardinalities of  $\mathbf{X}$  and  $\mathbf{Y}$  respectively, and  $D$  denotes the point set dimensionality. We assume the general case when  $M \neq N$  and  $D = 3$  in all experiments, although our method is directly applicable to  $D = 2$  and generalisable to  $D > 3$  if training data is available and a voxel grid is feasible in this dimension. Our objective is to find the displacement function (a vector field)  $v : \mathbb{R}^{M \times 3} \times \mathbb{R}^{N \times 3} \rightarrow \mathbb{R}^{M \times 3}$  so that  $\mathbf{Y} + v(\mathbf{Y}, \mathbf{X})$  matches  $\mathbf{X}$  as close as possible.

There is no universal criterion for optimal matching and it varies from scenario to scenario. We demand 1) that  $v$

ID	Layer	Output Size	Kernel	Padding/Stride	Concatenation
1	Input	64x64x64x2	-	-	-
2	3D Convolution	64x64x64x8	7x7x7	3/1	-
3	LeakyReLU	64x64x64x8	-	-	-
4	MaxPooling 3D	32x32x32x8	2x2x2	0/2	-
5	3D Convolution	32x32x32x16	5x5x5	2/1	-
6	LeakyReLU	32x32x32x16	-	-	-
7	MaxPooling 3D	16x16x16x16	2x2x2	0/2	-
8	3D Convolution	16x16x16x32	3x3x3	1/1	-
9	LeakyReLU	16x16x16x32	-	-	-
10	MaxPooling 3D	8x8x8x32	2x2x2	0/2	-
11	3D Convolution	8x8x8x64	3x3x3	1/1	-
12	LeakyReLU	8x8x8x64	-	-	-
13	3D Deconvolution	16x16x16x64	2x2x2	0/2	12 & 10
14	3D Deconvolution	16x16x16x64	3x3x3	1/1	-
15	LeakyReLU	16x16x16x64	-	-	-
16	3D Deconvolution	32x32x32x32	2x2x2	0/2	15 & 7
17	3D Deconvolution	32x32x32x32	5x5x5	2/1	-
18	LeakyReLU	32x32x32x32	-	-	-
19	3D Deconvolution	64x64x64x16	2x2x2	0/2	18 & 4
20	3D Deconvolution	64x64x64x16	7x7x7	3/1	-
21	LeakyReLU	64x64x64x16	-	-	-
22	3D Deconvolution	64x64x64x3	3x3x3	1/1	-

**Table 1:** U-Net-style architecture of DispVoxNet. The concatenation column contains the layer IDs whose outputs are concatenated and used as an input to the current layer. We use a negative slope for LeakyReLU of 0.01.

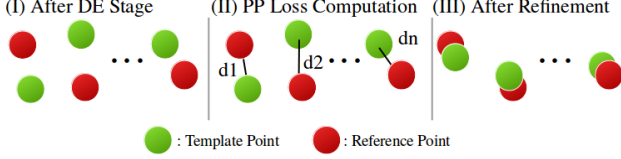
results in realistic class-specific object deformations so that the global alignment is recovered along with fine local deformations, and 2) that the template deformation preserves the point topology as far as possible. The first requirement remains very challenging for current general NRPSR methods. Either the shapes are globally matched while fine details are disregarded or the main deformation component is neglected which can lead to distorted registrations. The problem becomes even more ill-posed due to noise in the point sets. Some methods apply multiscale matching or parameter adjustment schemes [21, 18]. Even though a relaxation of the global topology-preserving constraint can lead to a finer local alignment, there is an increased chance of arriving at a local minimum and fitting to noise.

Let  $\mathbf{V}_X$  and  $\mathbf{V}_Y$  be the voxel grids, *i.e.*, the voxel-based proxies on which DispVoxNets regress deformations. Without loss of generality, we assume that  $\mathbf{V}_X$  and  $\mathbf{V}_Y$  are cubic and both have equal dimensions  $Q = 64$ . We propose to learn  $v$  as described next.

#### 3.1. Architecture

Our method is composed of *displacement estimation* and *refinement* stages. Each stage contains a DispVoxNet, *i.e.*, a 3D-convolutional neural network based on a U-Net architecture [48], which we also denote by  $\mathbf{D}_{vn}$ . See Figs. 2–3 and Table 1 for details on the network architecture.

**Displacement Estimation (DE) Stage.** We first discretise both  $\mathbf{Y}$  and  $\mathbf{X}$  on  $\mathbf{V}_Y$  and  $\mathbf{V}_X$ , respectively (*P2V* in Fig. 2). During the conversion, the point-to-voxel correspondences are stored in an affinity table. As voxel grids sample the space uniformly, each point of  $\mathbf{X}$  and  $\mathbf{Y}$  is mapped to one of the  $Q^3$  voxels in  $\mathbf{V}_X$  and  $\mathbf{V}_Y$ , respec-



**Figure 4:** (I) The displacement estimation stage regresses rough displacements between  $\mathbf{Y}$  and  $\mathbf{X}$ . (II) For all  $\mathbf{y}'$  in  $\mathbf{Y} + v(\mathbf{Y}, \mathbf{X})$ , we find the distance to the closest point  $\mathbf{x}_{\mathbf{y}'}$  among all points in  $\mathbf{X}$ . (III) At test time, the refined displacements yield a smoothed result  $\mathbf{Y} + v(\mathbf{Y}, \mathbf{X})$ .

tively. In  $\mathbf{V}_{\mathbf{X}}$  and  $\mathbf{V}_{\mathbf{Y}}$ , we represent  $\mathbf{X}$  and  $\mathbf{Y}$  as binary voxel occupancy indicators, *i.e.*, if at least one point falls into a voxel, the voxel’s occupancy is set to 1; otherwise, it equals to 0. Starting from  $\mathbf{V}_{\mathbf{X}}$  and  $\mathbf{V}_{\mathbf{Y}}$ , DispVoxNet regresses per-voxel displacements of the dimension  $Q^3 \times 3$ . During training, we penalise the discrepancy between the inferred voxel displacements and the ground truth displacements  $\mathbf{Z}$  using a mean squared error normalised by the number of voxels.  $\mathbf{Z}$  is obtained by converting the ground truth point correspondences to the voxel-based representation of dimensions  $Q^3 \times 3$  compatible with our architecture. The *point displacement* loss is given by:

$$\mathcal{L}_{\text{Disp.}}(\mathbf{Z}, \mathbf{V}_{\mathbf{Y}}, \mathbf{V}_{\mathbf{X}}) = \frac{1}{Q^3} \|\mathbf{Z} - \mathbf{D}_{\text{vn}}(\mathbf{V}_{\mathbf{Y}}, \mathbf{V}_{\mathbf{X}})\|_2^2. \quad (1)$$

Using the affinity table between  $\mathbf{Y}$  and  $\mathbf{V}_{\mathbf{Y}}$ , we determine each point’s displacement by applying trilinear interpolation on the eight nearest displacements in the voxel grid (*V2P* in Fig. 2 and see supplement). After adding the displacements to  $\mathbf{Y}$ , we observe that the resulting output after a single DispVoxNet bears some roughness. The refinement stage described in the following alleviates this problem.

**Refinement Stage.** Since the DE stage accounts for global deformations but misses some fine details, the unresolved residual displacements at the refinement stage are small. Recall that DispVoxNet is exposed to scenarios with small displacements during the training of the DE stage, since our datasets also contain similar (but not equal) states. Thus, assuming small displacements, we design a refinement stage as a combination of a pre-trained DispVoxNet and an additional unsupervised loss. Eventually, the refinement stage resolves the remaining small displacements and smooths the displacement field. To summarise, at the beginning of the refinement stage, the already deformed template point set is converted into the voxel representation  $\mathbf{V}_{\mathbf{Y}}^*$ . From  $\mathbf{V}_{\mathbf{Y}}^*$  and  $\mathbf{V}_{\mathbf{X}}$ , a pre-trained DispVoxNet learns to regress refined per-voxel displacements.

To apply the inferred voxel displacements to a template point at the end of the refinement stage (see Fig. 3), (I)<sup>2</sup> we

<sup>2</sup>(I), (II), (III) and (IV) refer to the steps in Fig. 3

compute the trilinear interpolation of the eight nearest displacements of  $\mathbf{y}_i, i \in \{1, \dots, M\}$ , and calculate a weighted consensus displacement for  $\mathbf{y}_i$ . (II) The weights and indices of the eight nearest voxels are saved in an affinity table. To further increase the accuracy, (III) we introduce the unsupervised, differentiable *point projection* (PP) loss between the final output  $\mathbf{Y} + v(\mathbf{Y}, \mathbf{X})$  and  $\mathbf{X}$ . The PP loss penalises the Euclidean distances between a point  $\mathbf{y}'$  in  $\mathbf{Y} + v(\mathbf{Y}, \mathbf{X})$  and its closest point  $\mathbf{x}_{\mathbf{y}'}$  in  $\mathbf{X}$ :

$$\mathcal{L}_{\text{PP}}(\mathbf{Y} + v(\mathbf{Y}, \mathbf{X}), \mathbf{X}) = \frac{1}{M} \sum_{i=1}^M \|\mathbf{y}'_i - \mathbf{x}_{\mathbf{y}'_i}\|_2. \quad (2)$$

We employ a  $k$ -d tree to determine  $\mathbf{x}_{\mathbf{y}'}$  for all  $\mathbf{y}'$  in  $\mathbf{Y} + v(\mathbf{Y}, \mathbf{X})$ , see Fig. 4 for a schematic visualisation.

Since the training is performed through backpropagation, we need to ensure the differentiability of all network stages. Our approach contains conversions from voxel to point-set representations and vice versa that are not fully differentiable. Thanks to the affinity table, we know the correspondences between points and voxels at the refinement stage. Therefore, (IV) gradients back-propagated from the PP loss can be distributed into the corresponding voxels in the voxel grid as shown in Fig. 3. As eight displacements contribute to the displacement of a point due to trilinear interpolation in the forward pass, the gradient of the point is back-propagated to the eight nearest voxels in the voxel grid according to the trilinear weights from the forward pass.

Two consecutive DispVoxNets in the DE and refinement stages implement a hierarchy with two granularity levels for the regression of point displacements. Combining more stages does not significantly improve the result, *i.e.*, we find two DispVoxNets necessary and sufficient.

### 3.2. Training Details

We use Adam [32] optimiser with a learning rate of  $3 \cdot 10^{-4}$ . As the number of points varies between the training pairs, we set the batch size to 1. We train the stages in two consecutive phases, starting with the DE stage using the displacement loss until convergence. This allows the network to regress rough displacements between  $\mathbf{Y}$  and  $\mathbf{X}$ . Then, another instance of DispVoxNet in the refinement stage is trained using only the PP loss. We initialise it with the weights from DispVoxNet of the DE stage. Since the PP loss considers only one nearest neighbour, we need to ensure that each output point from the DE stage is already close to its corresponding point in  $\mathbf{X}$ . Thus, we freeze the weights of DispVoxNet in the DE stage when training the refinement stage. See our supplement for training statistics.

To enhance robustness of the network to noise and missing points, 0 – 30% of the points are removed at random from  $\mathbf{Y}$  and  $\mathbf{X}$ , and uniform noise is added to both point sets. The number of added points ranges from 0% to 100%

		Ours	NR-ICP [9]	CPD [41]	GMR [31]
thin plate [17]	$e$	0.0103	0.0402	<b>0.0083</b> / 0.0192	0.2189
	$\sigma$	<b>0.0059</b>	0.0273	0.0102 / 0.0083	1.0121
FLAME [36]	$e$	0.0063	0.0588	<b>0.0043</b> / 0.0094	0.0056
	$\sigma$	0.0009	0.0454	0.0008 / <b>0.0005</b>	0.0007
DFAUST [5]	$e$	<b>0.0166</b>	0.0585	0.0683 / 0.0721	0.2357
	$\sigma$	<b>0.0020</b>	0.0215	0.0314 / 0.0258	0.8944
cloth [2]	$e$	<b>0.0080</b>	0.0225	0.0149 / 0.0138	0.2189
	$\sigma$	<b>0.0021</b>	0.0075	0.0066 / 0.0033	1.0121

**Table 2:** Comparison of registration errors for all tested methods. For CPD, we also report results with FGT (right-hand values).

	DE	DE + Ref. (nearest voxel)	Full: DE + Ref. (trilinear)
$e$	0.0100	0.0088	<b>0.0069</b>
$\sigma$	0.0021	0.0075	<b>0.0016</b>

**Table 3:** Ablation study highlighting the importance of the refinement stage and trilinear interpolation compared to a nearest voxel lookup.

of the number of points in the respective point set. The amount of noise per sample is determined randomly. When computing the PP loss, added noise is not considered.

## 4. Experiments

Our method is implemented in PyTorch [43]. The evaluation system contains two Intel(R) Xeon(R) E5-2687W v3 running at 3.10GHz and a NVIDIA GTX 1080Ti GPU. We compare DispVoxNets with four methods with publicly available code, *i.e.*, point-to-point non-rigid ICP (NR-ICP) [9], GMR [31], CPD [41] and CPD with Fast Gaussian Transform (FGT) [41]. FGT is a technique for the fast evaluation of Gaussians at multiple points [22].

### 4.1. Datasets

In total, we evaluate on four different datasets which represent various types of common 3D deformable objects, *i.e.*, *thin plate* [17], *FLAME* [36], *Dynamic FAUST (DFAUST)* [5] and *cloth* [2]. *Thin plate* contains 4648 states of a synthetic isometric surface. *FLAME* consists of a variety of captured human facial expressions (10k meshes in total). *DFAUST* is a scanned mesh dataset of human subjects in various poses (7948 meshes in total). Lastly, the *cloth* dataset contains 6238 captured states of a deformable sheet. Except for *FLAME*, the datasets are sequential and contain large non-linear deformations. Also, the deformation complexity in *FLAME* is lower than in the other datasets, *i.e.*, the deformations are mostly concentrated around the mouth area of the face scans.

We split the datasets into training and test subsets by considering blocks of a hundred point clouds. The first eighty samples from every block comprise the training set and the remaining twenty are included in the test set. For *FLAME*, we pick 20% of samples at random for testing and use the remaining ones for training. As all datasets have consistent topology, we directly obtain ground truth correspon-

dences which are necessary for training and error evaluation. We evaluate the registration accuracy of our method on clean samples (see Sec. 4.2) as well as in settings with uniform noise and clustered outliers (added sphere and removed chunk, see Sec. 4.3), since point cloud data captured by real sensors often contains noise and can be incomplete. In total, thirty template-reference pairs are randomly selected from each test dataset. We use the same pairs in all experiments. For the selected pairs, we report the average root-mean-square error (RMSE) between the references and aligned templates and standard deviation of RMSE, denoted by  $e$  and  $\sigma$  respectively:  $e = \frac{1}{M} \sum_{i=1}^M \frac{\|y_i - x_i\|_2}{\sqrt{D}}$ , with the template points  $y_i$  and corresponding reference points  $x_i$ .

### 4.2. Noiseless Data

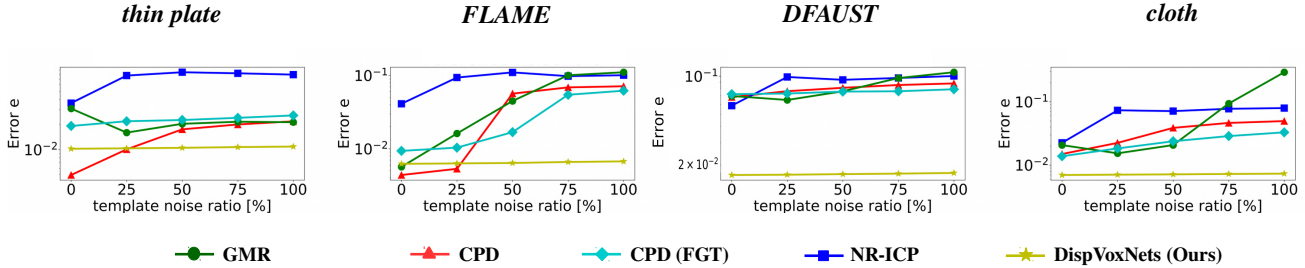
We first evaluate the registration accuracy of our method and several baselines on noiseless data. Table 2 and Fig. 1 summarise the results. Our approach significantly outperforms other methods on the *DFAUST* and *cloth* datasets which contain articulated motion and large non-linear deformations between the template and reference. On *thin plate*, DispVoxNets perform on par with CPD (CPD with FGT) and show a lower  $e$  in three cases out of four. On *FLAME*, which contains localised and small deformations, our approach achieves  $e$  of the same order of magnitude as CPD and GMR. CPD, GMR and DispVoxNets outperform NR-ICP in all cases. The experiment confirms the advantage of our approach in aligning point sets with large global non-linear deformations (additionally, see the supplement).

**Ablation Study.** We conduct an ablation study on the *cloth* dataset to test the influence of each component of our architecture. The tested cases include 1) only DE stage, 2) DE and refinement stages with nearest voxel lookup (the naïve alternative to trilinear interpolation) and 3) our entire architecture with DE and refinement stages, plus trilinear interpolation. Quantitative results are shown in Table 3. The full model with trilinear interpolation reduces the error by more than 30% over the DE only setting.

### 4.3. Deteriorated Data

The experiments with deteriorated data follow the evaluation protocol of Sec. 4.2. We introduce clustered outliers and uniform noise to the data.

**Structured Outliers.** We evaluate the robustness of our method to added structured outliers and missing data. We either add a sphere-like object to the inputs or arbitrarily remove a chunk from one of the point sets. As summarised in Tables 4–5, our approach shows the highest accuracy among all methods on *thin plate*, *DFAUST* and *cloth*, even though DispVoxNets were not trained with clustered outliers and have not been exposed to sphere-like structures.



**Figure 5:** Comparison of registration errors in the scenario with uniform noise.  $p\%$  is the ratio between the number of noise points added to the template and the number of points in the unperturbed template.  $e$  statistics of our approach is barely affected by the noise level. See the supplement for more results.

		Ours	NR-ICP [9]	CPD [41]	GMR [31]	
<i>thin plate</i> [17]	ref.	$e$	<b>0.0151</b>	0.1267 / 0.1136	0.6332	
		$\sigma$	<b>0.0117</b>	0.0302	0.0224 / 0.0211	1.5749
	temp.	$e$	<b>0.0150</b>	0.0509	0.0304 / 0.0636	0.0528
		$\sigma$	<b>0.0106</b>	0.0406	0.0200 / 0.0149	0.0300
<i>FLAME</i> [36]	ref.	$e$	0.0098	<b>0.0039</b>	0.0492 / 0.0617	0.0577
		$\sigma$	0.0034	<b>0.0007</b>	0.0301 / 0.0218	0.0205
	temp.	$e$	0.0073	0.0566	<b>0.0072</b> / 0.0246	0.0309
		$\sigma$	<b>0.0015</b>	0.0334	0.0070 / 0.0142	0.0117
<i>DFAUST</i> [5]	ref.	$e$	<b>0.0308</b>	0.0605	0.1127 / 0.1151	0.9730
		$\sigma$	<b>0.0111</b>	0.0226	0.0308 / 0.0295	2.2267
	temp.	$e$	<b>0.0190</b>	0.0669	0.0791 / 0.0775	0.0845
		$\sigma$	<b>0.0036</b>	0.0187	0.0304 / 0.0220	0.0295
<i>cloth</i> [2]	ref.	$e$	<b>0.0213</b>	0.0248	0.1081 / 0.1096	0.1098
		$\sigma$	<b>0.0091</b>	0.0095	0.0235 / 0.0223	0.0234
	temp.	$e$	0.0649	<b>0.0296</b>	0.0408 / 0.0522	0.0476
		$\sigma$	0.0395	<b>0.0081</b>	0.0115 / 0.0114	0.0223

**Table 4:** Registration errors for the case with clustered outliers. For CPD, we also report results in the mode with FGT (right-hand values). “ref.” and “temp.” denote whether outliers are added to  $\mathbf{X}$  or  $\mathbf{Y}$ , respectively.

		Ours	NR-ICP [9]	CPD [41]	GMR [31]	
<i>thin plate</i> [17]	ref.	$e$	<b>0.0107</b>	0.0668	0.0218 / 0.0386	0.4415
		$\sigma$	<b>0.0061</b>	0.0352	0.0148 / 0.0067	1.4632
	temp.	$e$	<b>0.0108</b>	0.0334	0.0479 / 0.0471	0.4287
		$\sigma$	0.0062	0.0281	0.0101 / <b>0.0038</b>	1.3832
<i>FLAME</i> [36]	ref.	$e$	0.0084	0.0519	<b>0.0046</b> / 0.0140	0.0193
		$\sigma$	0.0010	0.0451	0.0009 / <b>0.0006</b>	0.0008
	temp.	$e$	0.0088	0.0215	<b>0.0076</b> / 0.0201	0.0274
		$\sigma$	<b>0.0010</b>	0.0219	<b>0.0010</b> / 0.0016	0.0019
<i>DFAUST</i> [5]	ref.	$e$	<b>0.0167</b>	0.0463	0.0562 / 0.0636	0.0714
		$\sigma$	<b>0.0029</b>	0.0195	0.0308 / 0.0216	0.0282
	temp.	$e$	<b>0.0169</b>	0.0426	0.0672 / 0.0710	0.0737
		$\sigma$	<b>0.0033</b>	0.0194	0.0291 / 0.0229	0.0243
<i>cloth</i> [2]	ref.	$e$	<b>0.0090</b>	0.0455	0.0248 / 0.0315	0.0288
		$\sigma$	<b>0.0018</b>	0.0061	0.0056 / 0.0027	0.0087
	temp.	$e$	<b>0.0132</b>	0.0208	0.0486 / 0.0347	0.0397
		$\sigma$	0.0019	0.0087	0.0077 / <b>0.0014</b>	0.0092

**Table 5:** Registration errors for the case with missing parts. For CPD, we also report results in the mode with FGT (right-hand values). “ref.” and “temp.” denote whether parts are removed from  $\mathbf{X}$  or  $\mathbf{Y}$ , respectively.

Tables 4–5 report results for both the cases with modified references and templates, and Fig. 6 compares results qualitatively. DispVoxNets are less influenced by the outliers and produce visually accurate alignments of regions with correspondences in both point sets. CPD, GMR and NR-ICP suffer from various effects, *i.e.*, their alignments are severely influenced by outliers and the target regions are corrupted in many cases. We hypothesise that convolutional layers in DispVoxNet learn to extract informative features from the input points set and ignore noise. Furthermore, the network learns a class-specific deformation model which further enhances the robustness to outliers.

**Uniform Noises.** Next, we augment templates with uniform noise and repeat the experiment. Fig. 5 reports metrics for different amounts of added noise. Note that CPD, GMR and NR-ICP fail multiple times, and we define the success criterion in this experiment as  $e < (4 \times \text{median})$  followed by  $e < 4.0$ . DispVoxNets show stable accuracy across different noise ratios and datasets, while the error of other approaches increases significantly (up to 100 times) with an increasing amount of noise. Only our approach is agnostic to large amount of noise, despite CPD explicitly modeling a uniform noise component. For a qualitative comparison, see

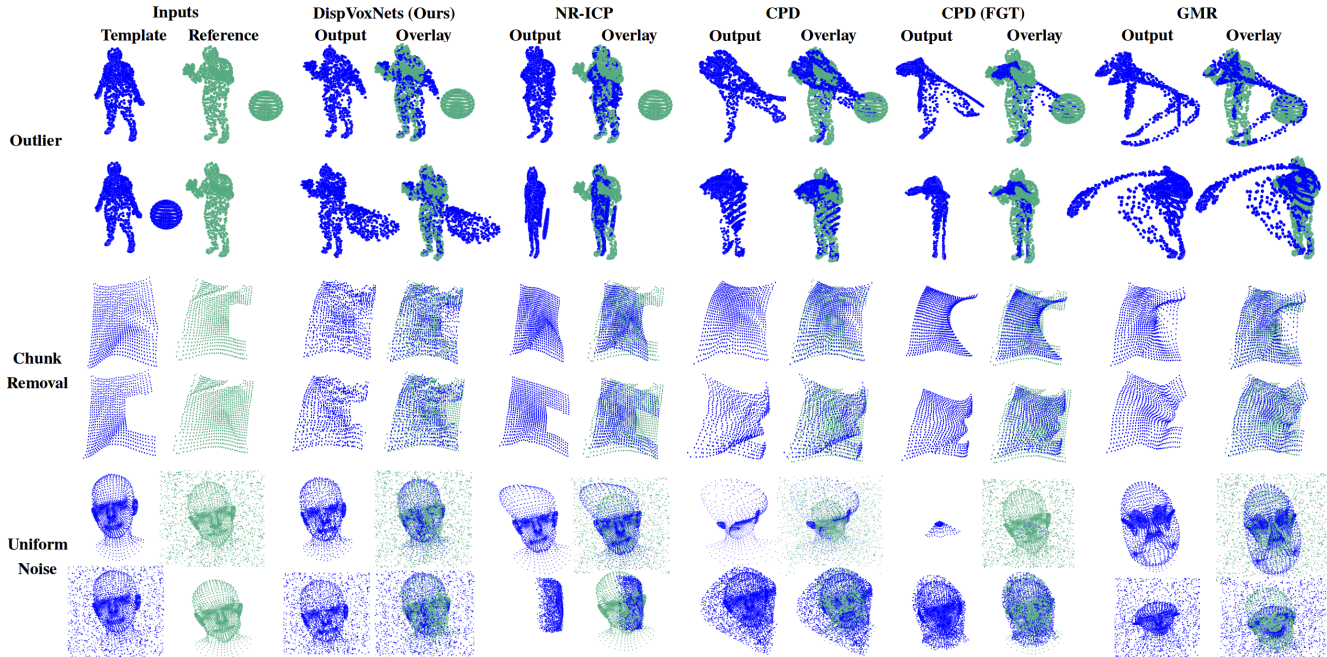
the fifth and sixth rows in Fig. 6 as well as the supplement.

#### 4.4. Runtime Analysis

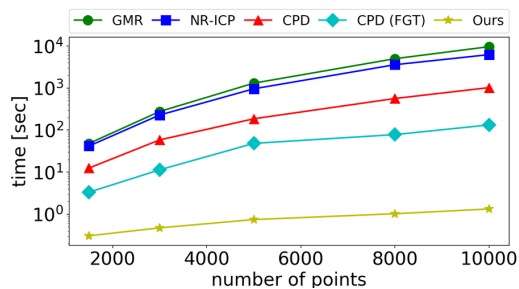
We prepare five point set pairs out of *DFAUST* dataset where the number of points varies from 1.5k to 10k. The runtime plots for different number of points are shown in Fig. 7. The numbers of points in both point sets are kept equal in this experiment. For 10k points, GMR, NR-ICP, CPD and CPD (FGT) take about 2 hours, 1.5 hours, 15 minutes, and 2 minutes per registration, respectively. Our approach requires only 1.5 seconds, which suggests its potential for applications with real-time constraints.

#### 4.5. Real Scans

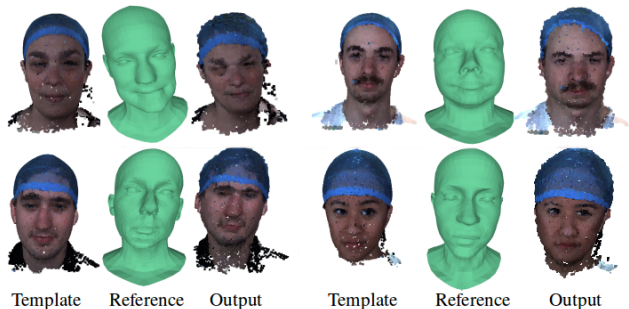
In this section, we demonstrate the generalisability of DispVoxNets to real-world scans. We test on 3D head point sets from the data collection of Dai *et al.* [10]. Since facial expression do not vary much in it, we use a reference from *FLAME* [36] and a template from [10]. Registration results can be seen in Fig. 8. Although some distortion in the output shape is recognisable, DispVoxNets transfer the reference facial expression to the template. Even though the network has only seen *FLAME* at training time, it is able to align two point sets of different cardinalities and origins.



**Figure 6:** Qualitative comparison of our approach to other methods in challenging scenarios with added clustered outliers (sphere added to either template or reference; the first and second rows), removed parts (the third and fourth rows) and 50% of added uniform noise (the fifth and sixth rows).



**Figure 7:** Runtime comparison for different number of points. The horizontal axis shows the number of points in the template and reference point sets. The vertical axis shows the  $\log_{10}$ -scaled inference time in seconds.



**Figure 8:** Aligning a real scan from [10] to a *FLAME* [36] reference. The references are shown as meshes for visualisation purposes.

## 5. Conclusions

We propose a DNN architecture with DispVoxNets — the first NRPSR approach with a supervised learning proxy — which regresses 3D displacements on regularly sampled voxel grids. Thanks to two consecutive DispVoxNets with trilinear interpolation and point projection, our approach outperforms other NRPSR methods by a large margin in scenarios with large deformations, articulations, clustered outliers and large amount of noise, while not relying on engineered class-specific priors. The runtime of DispVoxNets is around one second whereas other methods can take a few hours per registration, which suggests that our approach can be used in interactive applications. We show a high degree of generalisability of our architecture to real scans.

We believe that the direction revealed in this paper has a high potential for further investigation. Even though DispVoxNets represent a promising step towards accurate non-rigid point set alignment, its accuracy is limited by the resolution of the voxel grid and composition of training datasets. In future work, our method could be extended for operation on non-uniform voxel grids, and other types of losses could be investigated. Embedding of further alignment cues such as point normals, curvature, colours as well as sparse prior matches is a logical next step. We also expect to see extensions of DispVoxNets for point sets with small overlaps and adaptations of the proposed architecture for learning-based depth map fusion systems.



## References

- [1] S. A. Ali, V. Golyanik, and D. Stricker. Nrga: Gravitational approach for non-rigid point set registration. In *International Conference on 3D Vision (3DV)*, pages 756–765, 2018. [3](#), [13](#)
- [2] J. Bednařík, P. Fua, and M. Salzmann. Learning to reconstruct texture-less deformable surfaces from a single view. In *International Conference on 3D Vision (3DV)*, 2018. [2](#), [6](#), [7](#), [13](#), [14](#)
- [3] P. J. Besl and N. D. McKay. A method for registration of 3-d shapes. *Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 14(2):239–256, 1992. [1](#), [2](#), [12](#)
- [4] R. Bhotika, D. J. Fleet, and K. N. Kutulakos. A probabilistic theory of occupancy and emptiness. In *European Conference on Computer Vision (ECCV)*, pages 112–130, 2002. [3](#)
- [5] F. Bogo, J. Romero, G. Pons-Moll, and M. J. Black. Dynamic FAUST: Registering human bodies in motion. In *Computer Vision and Pattern Recognition (CVPR)*, 2017. [1](#), [2](#), [6](#), [7](#), [13](#), [14](#)
- [6] A. Broadhurst, T. W. Drummond, and R. Cipolla. A probabilistic framework for space carving. In *International Conference on Computer Vision (ICCV)*, 2001. [3](#)
- [7] Y. Chen and G. Medioni. Object modeling by registration of multiple range images. In *International Conference on Robotics and Automation (ICRA)*, 1991. [2](#)
- [8] C. B. Choy, D. Xu, J. Gwak, K. Chen, and S. Savarese. 3D-R2N2: A unified approach for single and multi-view 3d object reconstruction. In *European Conference on Computer Vision (ECCV)*, 2016. [3](#)
- [9] H. Chui and A. Rangarajan. A new point matching algorithm for non-rigid registration. *Computer Vision and Image Understanding (CVIU)*, 89(2-3):114–141, 2003. [1](#), [2](#), [6](#), [7](#), [13](#)
- [10] H. Dai, N. Pears, W. A. P. Smith, and C. Duncan. A 3d morphable model of craniofacial shape and texture variation. In *International Conference on Computer Vision (ICCV)*, 2017. [7](#), [8](#)
- [11] Y. Deng, A. Rangarajan, S. Eisenschenk, and B. C. Vemuri. A riemannian framework for matching point clouds represented by the schrödinger distance transform. In *Computer Vision and Pattern Recognition (CVPR)*, 2014. [3](#)
- [12] A. Fitzgibbon. Robust registration of 2d and 3d point sets. In *British Machine Vision Conference (BMVC)*, pages 1145–1153, 2003. [2](#)
- [13] S. Ge, G. Fan, and M. Ding. Non-rigid point set registration with global-local topology preservation. In *Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2014. [1](#), [3](#)
- [14] S. Gold, A. Rangarajan, C.-P. Lu, and E. Mjolsness. New algorithms for 2d and 3d point matching: Pose estimation and correspondence. *Pattern Recognition*, 31:957–964, 1997. [2](#)
- [15] V. Golyanik, S. A. Ali, and D. Stricker. Gravitational approach for point set registration. In *Computer Vision and Pattern Recognition (CVPR)*, 2016. [3](#)
- [16] V. Golyanik, G. Reis, B. Taetz, and D. Stricker. A framework for an accurate point cloud based registration of full 3d human body scans. In *Machine Vision Applications (MVA)*, pages 67–72, 2017. [3](#)
- [17] V. Golyanik, S. Shimada, K. Varanasi, and D. Stricker. Hdmnet: Monocular non-rigid 3d reconstruction with learned deformation model. In *Virtual Reality and Augmented Reality (EuroVR)*, pages 51–72, 2018. [1](#), [2](#), [3](#), [6](#), [7](#), [13](#), [14](#)
- [18] V. Golyanik, B. Taetz, G. Reis, and D. Stricker. Extended coherent point drift algorithm with correspondence priors and optimal subsampling. In *Winter Conference on Applications of Computer Vision (WACV)*, 2016. [1](#), [4](#)
- [19] V. Golyanik, C. Theobalt, and D. Stricker. Accelerated gravitational point set alignment with altered physical laws. In *International Conference on Computer Vision (ICCV)*, 2019. [3](#)
- [20] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Neural Information Processing Systems (NIPS)*, pages 2672–2680, 2014. [3](#)
- [21] S. Granger and X. Pennec. Multi-scale em-icp: A fast and robust approach for surface registration. In *European Conference on Computer Vision (ECCV)*, 2002. [2](#), [4](#)
- [22] L. Greengard and J. Strain. The fast gauss transform. *SIAM Journal on Scientific and Statistical Computing*, 12(1):79–94, 1991. [6](#)
- [23] M. Greenspan and G. Godin. A nearest neighbor method for efficient icp. In *International Conference on 3D Digital Imaging and Modeling (3DIM)*, pages 161–168, 2001. [2](#)
- [24] D. Hähnel, S. Thrun, and W. Burgard. An extension of the icp algorithm for modeling nonrigid objects with mobile robots. In *International Joint Conference on Artificial Intelligence (IJCAI)*, pages 915–920, 2003. [2](#)
- [25] R. Hanocka, N. Fish, Z. Wang, R. Giryes, S. Fleishman, and D. Cohen-Or. Alignet: Partial-shape agnostic alignment via unsupervised learning. *SIGGRAPH*, 38(1), 2018. [3](#)
- [26] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. [3](#)
- [27] B. K. P. Horn, H. M. Hilden, and S. Negahdaripour. Closed-form solution of absolute orientation using orthonormal matrices. *Journal of the Optical Society of America A (JOSA A)*, 5(7):1127–1135, 1988. [2](#)
- [28] D. Jack, J. K. Pontes, S. Sridharan, C. Fookes, S. Shirazi, F. Maire, and A. Eriksson. Learning free-form deformations for 3d object reconstruction. In *Asian Conference on Computer Vision (ACCV)*, pages 317–333, 2018. [3](#)
- [29] M. Jaderberg, K. Simonyan, A. Zisserman, and K. Kavukcuoglu. Spatial transformer networks. In *Neural Information Processing Systems (NIPS)*. 2015. [3](#)
- [30] P. Jauer, I. Kuhlemann, R. Bruder, A. Schweikard, and F. Ernst. Efficient registration of high-resolution feature enhanced point clouds. *Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 41(5):1102–1115, 2019. [3](#)
- [31] B. Jian and B. C. Vemuri. A robust algorithm for point set registration using mixture of gaussians. In *International Conference for Computer Vision (ICCV)*, pages 1246–1251, 2005. [1](#), [3](#), [6](#), [7](#), [12](#), [13](#)
- [32] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint 1412.6980*, 2014. [5](#)

- [33] K. Kolev, M. Klodt, T. Brox, and D. Cremers. Continuous global optimization in multiview 3d reconstruction. *International Journal of Computer Vision (IJCV)*, 84(1):80–96, 2009. 3
- [34] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Neural Information Processing Systems (NIPS)*, 2012. 3
- [35] A. Kurenkov et al. Deformnet: Free-form deformation network for 3d shape reconstruction from a single image. In *Winter Conference on Applications of Computer Vision (WACV)*, pages 858–866, 2018. 3
- [36] T. Li, T. Bolkart, M. J. Black, H. Li, and J. Romero. Learning a model of facial shape and expression from 4D scans. *SIGGRAPH Asia*, 36(6), 2017. 2, 6, 7, 8, 13, 14
- [37] S. Liu and D. B. Cooper. Statistical inverse ray tracing for image-based 3d modeling. *Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 36(10):2074–2088, 2014. 3
- [38] J. Ma, J. Zhao, J. Jiang, and H. Zhou. Non-rigid point set registration with robust transformation estimation under manifold regularization. In *Conference on Artificial Intelligence (AAAI)*, pages 4218–4224, 2017. 13
- [39] J. Ma, J. Zhao, J. Tian, Z. Tu, and A. L. Yuille. Robust estimation of nonrigid transformation for point set registration. In *Computer Vision and Pattern Recognition (CVPR)*, pages 2147–2154, 2013. 3
- [40] L. Mundermann, S. Corazza, and T. P. Andriacchi. Accurately measuring human movement using articulated icp with soft-joint constraints and a repository of articulated models. In *Computer Vision and Pattern Recognition (CVPR)*, 2007. 1, 3
- [41] A. Myronenko and X. Song. Point-set registration: Coherent point drift. *Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2010. 1, 2, 3, 6, 7, 12, 13
- [42] R. A. Newcombe, D. Fox, and S. M. Seitz. Dynamicfusion: Reconstruction and tracking of non-rigid scenes in real-time. In *Computer Vision and Pattern Recognition (CVPR)*, 2015. 3
- [43] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer. Automatic differentiation in pytorch. In *Neural Information Processing Systems (NIPS) Workshops*, 2017. 6
- [44] S. Pellegrini, K. Schindler, and D. Nardi. A generalisation of the icp algorithm for articulated bodies. In *British Machine Vision Conference (BMVC)*, pages 87.1–87.10, 2008. 1, 3
- [45] T. Pollard and J. L. Mundy. Change detection in a 3-d world. In *Computer Vision and Pattern Recognition (CVPR)*, 2007. 3
- [46] A. Pumarola, A. Agudo, L. Porzi, A. Sanfeliu, V. Lepetit, and F. Moreno-Noguer. Geometry-aware network for non-rigid shape prediction from a single view. In *Computer Vision and Pattern Recognition (CVPR)*, 2018. 3
- [47] G. Riegler, O. Ulusoy, and A. Geiger. Octnet: Learning deep 3d representations at high resolutions. In *Computer Vision and Pattern Recognition (CVPR) 2017*, 2017. 3
- [48] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention (MICCAI)*, 2015. 4
- [49] S. Rusinkiewicz and M. Levoy. Efficient variants of the ICP algorithm. In *International Conference on 3-D Digital Imaging and Modeling (3DIM)*, 2001. 2
- [50] N. Savinov, L. Ladický, C. Häne, and M. Pollefeys. Discrete optimization of ray potentials for semantic 3d reconstruction. In *Computer Vision and Pattern Recognition (CVPR)*, pages 5511–5518, 2015. 3
- [51] S. M. Seitz and C. R. Dyer. Photorealistic scene reconstruction by voxel coloring. *International Journal of Computer Vision (IJCV)*, 35(2):151–173, 1999. 3
- [52] S. Shimada, V. Golyanik, C. Theobalt, and D. Stricker. Ismog: Adversarial learning for monocular non-rigid 3d reconstruction. In *Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2019. 3
- [53] R. W. Sumner, J. Schmid, and M. Pauly. Embedded deformation for shape manipulation. In *SIGGRAPH*, 2007. 3
- [54] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *Computer Vision and Pattern Recognition (CVPR)*, 2015. 3
- [55] A. Tagliasacchi, M. Schröder, A. Tkach, S. Bouaziz, M. Botsch, and M. Pauly. Robust articulated-icp for real-time hand tracking. *Computer Graphics Forum (Symposium on Geometry Processing)*, 34(5), 2015. 1, 3
- [56] G. K. Tam et al. Registration of 3d point clouds and meshes: A survey from rigid to nonrigid. *Transactions on Visualization and Computer Graphics (TVCG)*, 19(7):1199–1217, 2013. 3
- [57] J. Taylor et al. Efficient and precise interactive hand tracking through joint, continuous optimization of pose and correspondences. *SIGGRAPH*, 35, 2016. 1
- [58] E. Tretschk, A. Tewari, M. Zollhöfer, V. Golyanik, and C. Theobalt. DEMEA: Deep Mesh Autoencoders for Non-Rigidly Deforming Objects. *arXiv e-prints*, 2019. 3
- [59] A. O. Ulusoy, A. Geiger, and M. J. Black. Towards probabilistic volumetric reconstruction using ray potentials. In *International Conference on 3D Vision (3DV)*, 2015. 3
- [60] Z. Wang, S. Rosa, and A. Markham. Learning the intuitive physics of non-rigid object deformations. In *Neural Information Processing Systems (NIPS) Workshops*, 2018. 3
- [61] W. Xu, A. Chatterjee, M. Zollhöfer, H. Rhodin, D. Mehta, H.-P. Seidel, and C. Theobalt. Monoperfcap: Human performance capture from monocular video. *SIGGRAPH*, 37(2):27:1–27:15, 2018. 3
- [62] Y. Zheng and D. Doermann. Robust point matching for non-rigid shapes by preserving local neighborhood structures. *Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 28(4):643–649, 2006. 2
- [63] Z. Zhou, J. Zheng, Y. Dai, Z. Zhou, and S. Chen. Robust non-rigid point set registration using student’s-t mixture model. *PLOS ONE*, 9, 03 2014. 3
- [64] H. Zhu, B. Guo, K. Zou, Y. Li, K.-V. Yuen, L. Mihaylova, and H. Leung. A review of point set registration: From pairwise registration to groupwise registration. *Sensors*, 19(5), 2019. 3

## A. Appendix

In this supplement, we provide details on the interpolation of the coarse displacement field (Sec. A.1) and report training statistics (Sec. A.2). We show more qualitative comparisons (Sec. A.3) as well as graphs for further cases with uniform noise (Sec. A.4).

### A.1. Interpolation of the 3D Displacement Field

Due to the limited resolution of the voxel grid, we apply trilinear interpolation to obtain displacements for every template point at sub-voxel precision. Note that in DE stage, interpolation is applied only in the forward pass. In the refinement stage, it is applied in the forward pass, and the computed trilinear weights are used during backpropagation to weight the gradients.

Suppose  $\bar{\mathbf{D}}: \mathbb{Z}^3 \rightarrow \mathbb{R}^3$  is the initial regressed 3D displacement field on a regular lattice induced by the voxel grid. Suppose the template point of interest after the DE stage  $\mathbf{y}_j^* = (x_j, y_j, z_j)$ ,  $j \in \{1, \dots, M\}$ , falls into a neighbourhood cube between eight displacement values of  $\bar{\mathbf{D}}$ . We denote these boundary displacements compactly by  $\mathbf{d} = \{\mathbf{d}_{abc}\}$ ,  $a, b, c \in \{0, 1\}$  on a unit cube<sup>3</sup> in a local coordinate system, see Fig. I for a schematic visualisation. In the refinement stage, we store for every  $\mathbf{y}_j^*$  the index of the voxel it belongs to, the indexes of the eight nearest displacements as well as the corresponding trilinear interpolation weights  $\mathbf{w} \in \mathbb{R}^8$  in the point affinity table. The latter is then used in the backward pass of the refinement stage.

Let  $x_{\max}, y_{\max}, z_{\max}$  and  $x_{\min}, y_{\min}, z_{\min}$  be the maximum and minimum  $x$ -,  $y$ - and  $z$ -values among the eight nearest lattice point coordinates, respectively. To convert  $\mathbf{y}_j^*$  from the coordinate system of the lattice to the local coordinate system, we calculate normalised distances  $l_x, l_y$  and  $l_z$ :

$$l_x = \frac{x_j - x_{\min}}{x_{\max} - x_{\min}}, l_y = \frac{y_j - y_{\min}}{y_{\max} - y_{\min}} \text{ and } l_z = \frac{z_j - z_{\min}}{z_{\max} - z_{\min}}. \quad (3)$$

The individual displacement  $\bar{\mathbf{v}}_j$  of  $\mathbf{y}_j^*$  is obtained by trilinear interpolation of the eight nearest displacements, *i.e.*, as an inner product of  $\mathbf{w}$  and  $x$ -,  $y$ - and  $z$ -components of  $\mathbf{d}$ :

$$\bar{\mathbf{v}}_{j,x} = \mathbf{w}^\top \mathbf{d}_x = \begin{bmatrix} (1-l_x)(1-l_y)(1-l_z) \\ (1-l_x)(1-l_y)l_z \\ (1-l_x)l_y(1-l_z) \\ l_x(1-l_y)(1-l_z) \\ (1-l_x)l_y l_z \\ l_x l_y (1-l_z) \\ l_x(1-l_y)l_z \\ l_x l_y l_z \end{bmatrix}^\top \begin{bmatrix} \mathbf{d}_{000,x} \\ \mathbf{d}_{001,x} \\ \mathbf{d}_{010,x} \\ \mathbf{d}_{100,x} \\ \mathbf{d}_{011,x} \\ \mathbf{d}_{110,x} \\ \mathbf{d}_{101,x} \\ \mathbf{d}_{111,x} \end{bmatrix}, \quad (4)$$

<sup>3</sup> $\mathbf{d}_{abc}$  is a shorthand notation for the displacement at point  $(x, y, z)$  in the local coordinate system, *i.e.*, at  $(0, 0, 0)$ ,  $(0, 0, 1)$ ,  $(0, 1, 0)$ , *etc.*

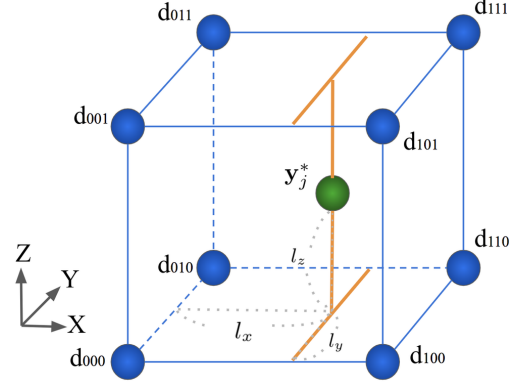


Figure I: Schematic visualisation of trilinear interpolation for a given  $\mathbf{y}_j^*$ .

	<i>thin plate</i>	<i>FLAME</i>	<i>DFAUST</i>	<i>cloth</i>
DE stage	530k	400k	715k	500k
refinement	14k	20k	24k	12k

Table I: Number of training iterations for DispVoxNets in the DE and refinement stages for the tested datasets.

$$\bar{\mathbf{v}}_{j,y} = \mathbf{w}^\top \mathbf{d}_y = \begin{bmatrix} (1-l_x)(1-l_y)(1-l_z) \\ (1-l_x)(1-l_y)l_z \\ (1-l_x)l_y(1-l_z) \\ l_x(1-l_y)(1-l_z) \\ (1-l_x)l_y l_z \\ l_x l_y (1-l_z) \\ l_x(1-l_y)l_z \\ l_x l_y l_z \end{bmatrix}^\top \begin{bmatrix} \mathbf{d}_{000,y} \\ \mathbf{d}_{001,y} \\ \mathbf{d}_{010,y} \\ \mathbf{d}_{100,y} \\ \mathbf{d}_{011,y} \\ \mathbf{d}_{110,y} \\ \mathbf{d}_{101,y} \\ \mathbf{d}_{111,y} \end{bmatrix}, \quad (5)$$

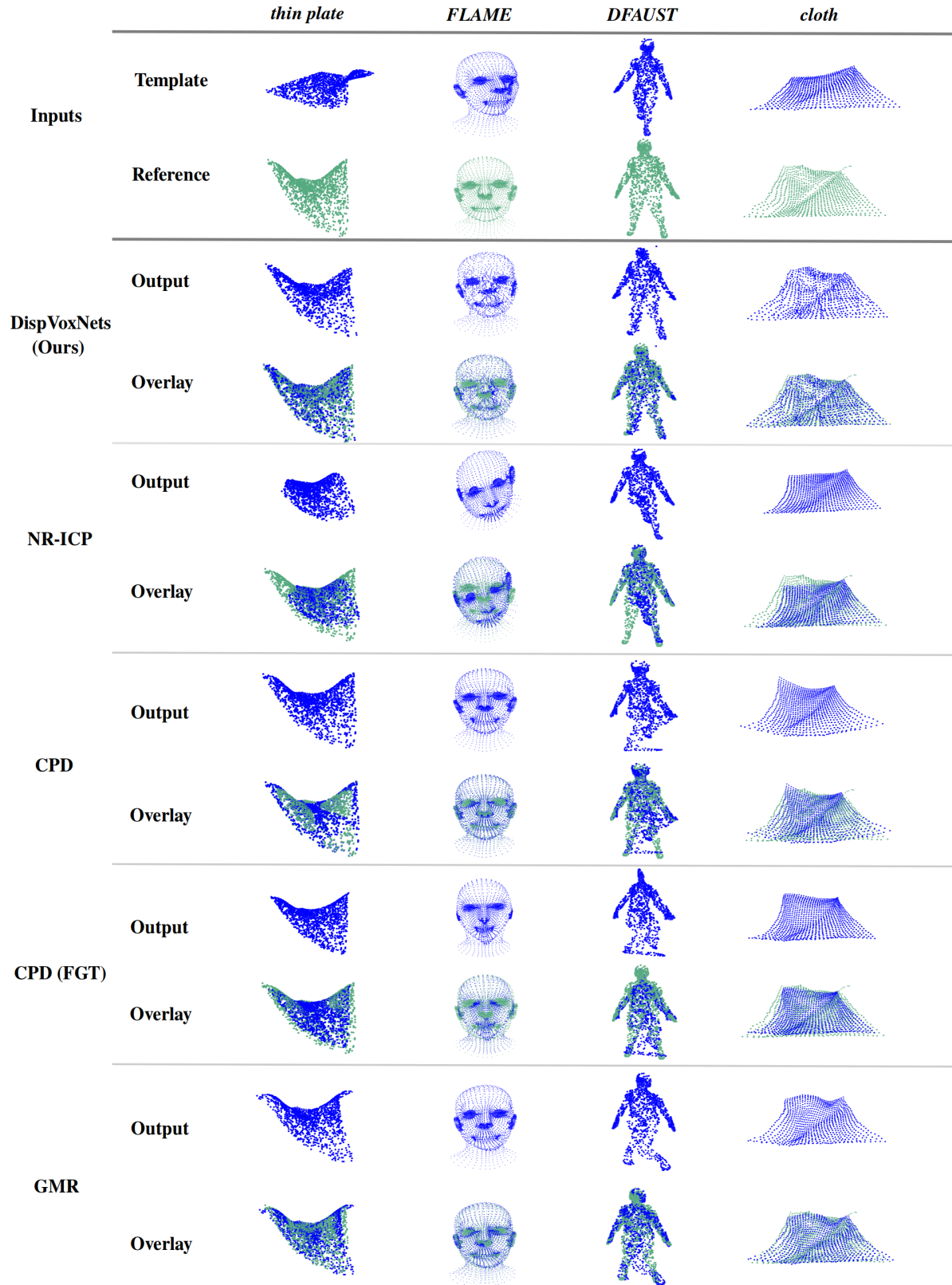
and

$$\bar{\mathbf{v}}_{j,z} = \mathbf{w}^\top \mathbf{d}_z = \begin{bmatrix} (1-l_x)(1-l_y)(1-l_z) \\ (1-l_x)(1-l_y)l_z \\ (1-l_x)l_y(1-l_z) \\ l_x(1-l_y)(1-l_z) \\ (1-l_x)l_y l_z \\ l_x l_y (1-l_z) \\ l_x(1-l_y)l_z \\ l_x l_y l_z \end{bmatrix}^\top \begin{bmatrix} \mathbf{d}_{000,z} \\ \mathbf{d}_{001,z} \\ \mathbf{d}_{010,z} \\ \mathbf{d}_{100,z} \\ \mathbf{d}_{011,z} \\ \mathbf{d}_{110,z} \\ \mathbf{d}_{101,z} \\ \mathbf{d}_{111,z} \end{bmatrix}. \quad (6)$$

Note that  $\mathbf{w}$ ,  $l_x$ ,  $l_y$  and  $l_z$  are shared across all dimensions.

### A.2. Training Statistics

Table I shows the number of training iterations until convergence for each dataset. Since *DFAUST* contains relatively large displacements between point sets, it requires the highest number of iterations followed by *thin plate* and *cloth*. On the contrary, *FLAME* contains only small displacements, and the network requires fewer parameter updates to converge compared to other datasets.



**Figure II:** Qualitative comparison of our DispVoxNets approach and other methods (NR-ICP [3], CPD/CPD with FGT [41] and GMR [31]). The input samples from each dataset are shown in the top rows, followed by the results (aligned templates and overlaid samples) for every method.

### A.3. Qualitative Analysis and Observations

In this section, we provide additional qualitative results. In Fig. II, we show selected registrations by our approach and other tested methods (NR-ICP [9], CPD/CPD with FGT [41], and GMR [31]) on the tested datasets (*thin plate* [17], *FLAME* [36], *DFAUST* [5] and *cloth* [2]).

On the *thin plate* — due to the rather simple object structure — all approaches except NR-ICP align the point sets reasonably accurate. CPD and DispVoxNets produce qualitatively similar results in the shown example. All methods show similar qualitative accuracy on the *cloth* dataset, while differences are noticeable in the corners and areas with large wrinkles. At the same time, only our approach simultaneously captures both small and large wrinkles. Thus, many fine foldings present in the reference surface are not well recognisable in the aligned templates in the case of NR-ICP, CPD/CPD with FGT and GMR. All in all, results of these methods appear to be oversmoothed.

In the absence of large displacements between the point sets — which is the case with *FLAME* dataset — model-based approaches CPD and GMR regress the displacements most accurately. The result of DispVoxNets is of comparable quality, though the deformed template is perceptually rougher and the points are arranged less regularly. This is due to the intermediate conversion steps from the point cloud representation to the voxel grid and back. We see that for small displacements, the limited resolution of the voxel grid is a more influential factor on the accuracy than the deformation prior learned from the data. With an increase of the voxel grid resolution, we expect our approach to come closer to CPD and GMR, up to the complete elimination of the accuracy gap (this is the matter of future work; currently, our focus is handling of large deformations which is a more challenging problem).

Next, we see that model-based approaches with global regularisers often fail on the *FAUST* dataset, while the proposed approach demonstrates superior quantitative and visual accuracy. Even though the surface produced by DispVoxNets after the refinement stage can still seem coarse at some parts, the overall pose and shape are correctly and realistically inferred as we expect, despite substantial differences between the template and reference in the feet area (a subject standing on one foot and a subject standing on both feet respectively). Thus, model-based methods have difficulty in aligning the feet.

Overall, the qualitative results in Fig. II demonstrate the advantages of DispVoxNets for non-rigid point set alignment over classic, non-supervised learning-based approaches. Since our technique learns class-specific priors implicitly during training, it is successful in registering samples with large displacements and articulations.

### A.4. Additional Experiments with Noisy Data

We present further experimental results with uniform noise in this section. Fig. III shows RMSE graphs for various combinations of uniform noise ratios in the reference and template for all four datasets (*thin plate* [17], *FLAME* [36], *DFAUST* [5] and *cloth* [2]).

For previous methods, we observe the tendency that adding uniform noise to both the template and the reference can result in a lower error than only adding it to one of them. It is reasonable to assume that two point sets certainly differ more if noise is added to only one of them. Thus, when inputs contain a similar amount of noise (we can say that the noise levels correlate), we observe the tendency that the alignment error becomes lower (e.g., see the seventh row, third column), i.e., some graphs roughly show a U-curve bottoming out at around 50% of the added noise in the template. We hypothesise that this is due to what we call the *mutual noise compensation effect*. Further study is required to clarify a more precise reason (it is possible that our observations are dataset-specific). Note that adding noise to both point sets is not a common evaluation setting. Usually, either template or reference is augmented with noise (cf. experimental sections in [31, 41, 38, 1]). With our experiment, we go beyond the prevalent evaluation methodology with noisy point sets.

On the one hand, CPD has the most stable error curve among the four model-based approaches, followed by NR-ICP and GMR. GMR shows higher errors when the noise is only in the template rather than in the reference, and CPD with FGT is the least stable as the noise ratio increases. Moreover, we observe that the relative performance of NR-ICP increases with the added noise. Thus, NR-ICP outperforms CPD only on *DFAUST* according to Table 2 (the experiment with no added noise). In Fig. III, we recognise multiple cases when NR-ICP outperforms CPD also on *FLAME* (the blue curve is below the red curve).

On the other hand, our approach with DispVoxNets shows almost constant error through all noise ratio combinations and all datasets. Compared to the case without noise, it even achieves the lowest RMSE on *FLAME* for multiple noise level combinations (~40% of the cases). As our network becomes aware of class-specific features after the training and learns to ignore noise, it can distinguish the meaningful shapes from noise, which contributes to its overall robustness. To the best of our belief, it is for the first time that a NRPSR method is so stable, even in the presence of large amount of noise in the data. Recall that we follow a simple noise augmentation policy for the training data (Sec. 3.2). Thus, our framework seemingly learns to filter uniform noise. Another factor could be that individual unstructured points cause neuron deactivations. In future work, it could be interesting to study augmentation policies for further types of noise (e.g., noise along the surfaces).

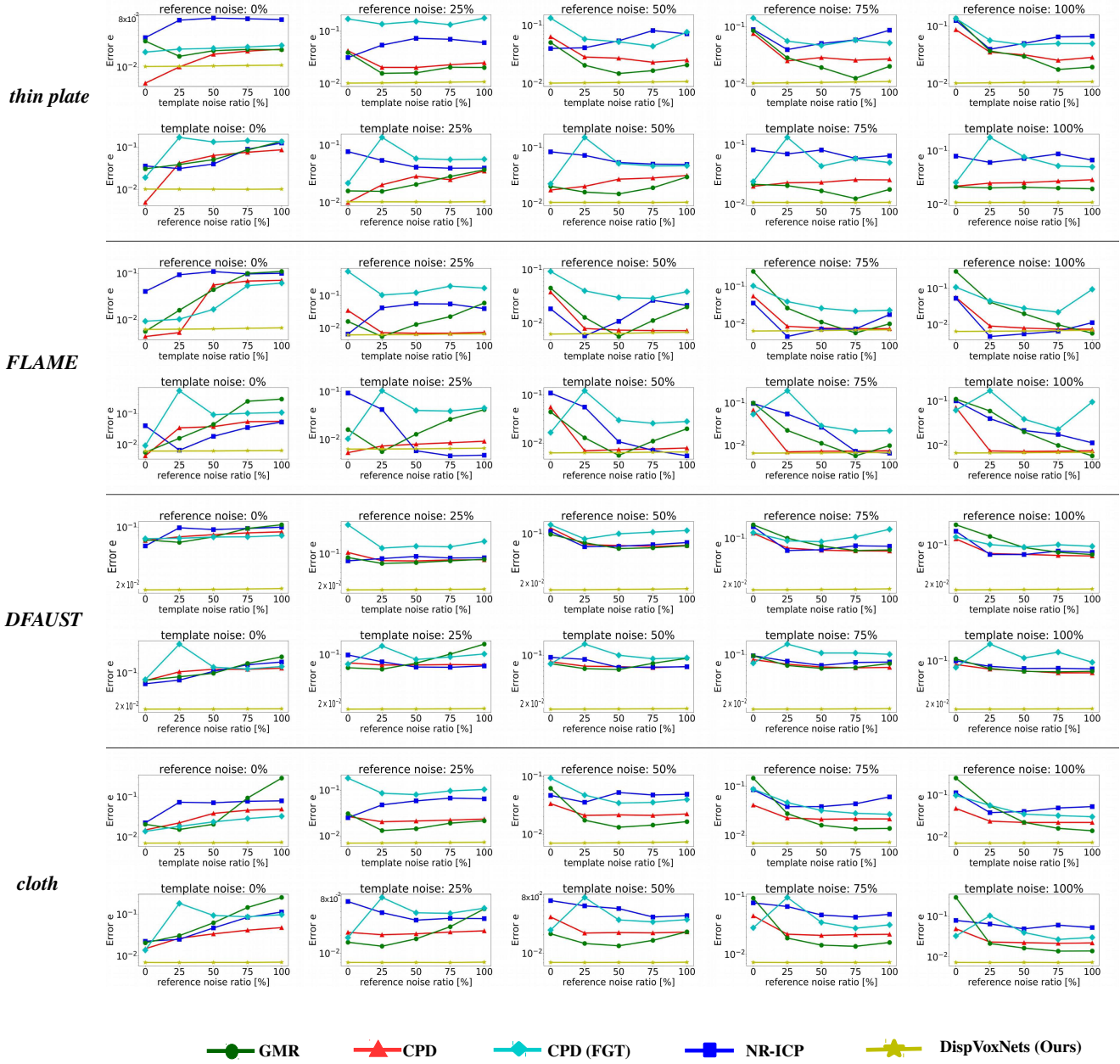


Figure III: RMSE ( $e$ ) graphs for additional experiments with uniform noise on *thin plate* [17], *FLAME* [36], *DFAUST* [5] and *cloth* [2] datasets.  $p\%$  is the ratio between the number of added points and the number of points in the sample. In this experiment, both reference and template are augmented with noise.