

# VIRDO: Visio-tactile Implicit Representations of Deformable Objects

Youngsun Wi<sup>1</sup>, Pete Florence<sup>2</sup>, Andy Zeng<sup>2</sup> and Nima Fazeli<sup>1</sup>

**Abstract**—Deformable object manipulation requires computationally efficient representations that are compatible with robotic sensing modalities. In this paper, we present VIRDO: an implicit, multi-modal, and continuous representation for deformable-elastic objects. VIRDO operates directly on visual (point cloud) and tactile (reaction forces) modalities and learns rich latent embeddings of contact locations and forces to predict object deformations subject to external contacts. Here, we demonstrate VIRDOs ability to: i) produce high-fidelity cross-modal reconstructions with dense unsupervised correspondences, ii) generalize to unseen contact formations, and iii) state-estimation with partial visio-tactile feedback.

## I. INTRODUCTION

Dexterous manipulation of deformable objects is an important open problem in robotics [1], [2]. These objects are ubiquitous in our day-to-day lives and play a key role in many applications including cooking, manufacturing, in-home assistive care, and surgery. Despite their prevalence and importance, deformable objects have received less attention than their rigid counterparts owing to their inherent complexities in modeling, perception, and controls [2]–[5]. To illustrate, the states of rigid bodies with known geometries can be succinctly represented with 6D pose and velocity. However, deformable objects have an infinite continuum of states and their representation and perception remains an open problem [2].

In this paper, we present VIRDO – an implicit, dense, cross modal, and continuous architecture that addresses these fundamental representation and perception challenges for the class of elastically deformable objects. The central feature of our method is learning deformation fields informed by cross modal visual and tactile cues of external contacts. We further contribute a dataset of elastically deformable objects with boundary conditions used to evaluate. This paper focuses on dense geometric representations because they can facilitate downstream tasks such as state-estimation from partial views and estimating dense correspondences, as we demonstrate, as well as bootstrapping keypoint/affordance learning.

**Problem Statement & Assumptions :** Our goal is to derive a computationally efficient and generative model that: 1) predicts object deformations subject to external forces; and 2) is compatible with common robotic sensors. We assume the object geometry is described by its point cloud: an unordered set  $\mathbf{P} := \{\mathbf{p} \in \mathbb{R}^3 : \text{SDF}(\mathbf{p}) = 0\}$  where SDF denotes the signed distance w.r.t. the surface of the object. Point clouds are obtained from commodity depth sensors or

<sup>1</sup> Youngsun Wi and Nima Fazeli are with the Robotics Institute at the University of Michigan, MI, USA <yswi, nfz>@umich.edu

<sup>2</sup> Pete Florence and Andy Zeng with Robotics at Google, Mountain View, CA, USA <peteflorence, andyzeng>@google.com

3D scanners commonly found in industry. Contact locations are also given as a set of points  $\mathbf{Q}$  which can be given by an upstream perception algorithm such as [6]–[8]. For the tactile input, net reaction force is given by  $\mathbf{u} \in \mathbb{R}^3$  at the wrist of the robot which can be measured by common industrial F/T sensors or recovered from joint torques. In this paper, we derive a continuous and implicit representation of the deformed object geometry ( $f(\mathbf{P}, \mathbf{Q}, \mathbf{u}) = s$ ) subject to external forces and their locations. Here the object geometry is given by the zero-level set of the implicit function; i.e.  $s = 0$ .

## II. REPRESENTING DEFORMABLE OBJECTS USING SDFs

At a high-level, VIRDO decomposes object representations into a nominal shape representation and a point-wise deformation field. Here, we choose signed distance fields as our underlying representation and discuss this choice in Sec. IV. The nominal shape representation decodes latent shape embeddings  $\boldsymbol{\alpha} \in \mathbb{R}^l$  into continuous signed-distance fields, the zero-level set of which is the undeformed object geometry – similar to the architectures proposed in [9]–[11]. The point-wise deformation field is produced using a summary of all boundary conditions (contact locations, reaction force, and fixed constraints) leveraging a permutation invariant set operator. The structure is fully differentiable and can be learned end-to-end. In the following, we discuss each component in more detail.

### A. Nominal Shape Representation

The nominal shape of an object is the geometry it takes in the absence of external contact forces. The nominal geometry is produced by the object module as  $\mathbf{O}(\mathbf{x}|\Psi_o(\boldsymbol{\alpha})) = s$  as a signed distance field, where  $\mathbf{x} = (x, y, z)$  is a query point, and  $s$  is the signed-distance. The purpose of the object code ( $\boldsymbol{\alpha}$ ) is to allow VIRDO to represent multiple objects. Here, we use a hyper-network  $\Psi_o(\boldsymbol{\alpha})$  to decode the object code into object module’s weights and biases  $\boldsymbol{\theta}_o$ , similar to [12].

We pre-train VIRDO on nominal shapes before training on deformations. During this stage, we initialize the object codes as  $\boldsymbol{\alpha} \sim N(0, 0.1^2)$  and update them with hyper-network  $\Psi_o$ . The loss function for pretraining nominal shapes is:

$$L_{nominal} = L_{sdf} + \lambda_2 L_{latent} + \lambda_3 L_{hyper}, \quad (1)$$

where  $L_{sdf}$  is defined as:

$$\begin{aligned} L_{sdf} = \sum_{i=1}^N \left( \sum_{\tilde{\mathbf{x}} \in \Omega} |clamp(\mathbf{O}^i(\mathbf{x}), \delta) - clamp(s^*, \delta)| \right. \\ \left. + \lambda \sum_{\tilde{\mathbf{x}} \in \Omega_0} (1 - \langle \nabla \mathbf{O}^i(\mathbf{x}), \mathbf{n}^* \rangle) \right). \end{aligned}$$

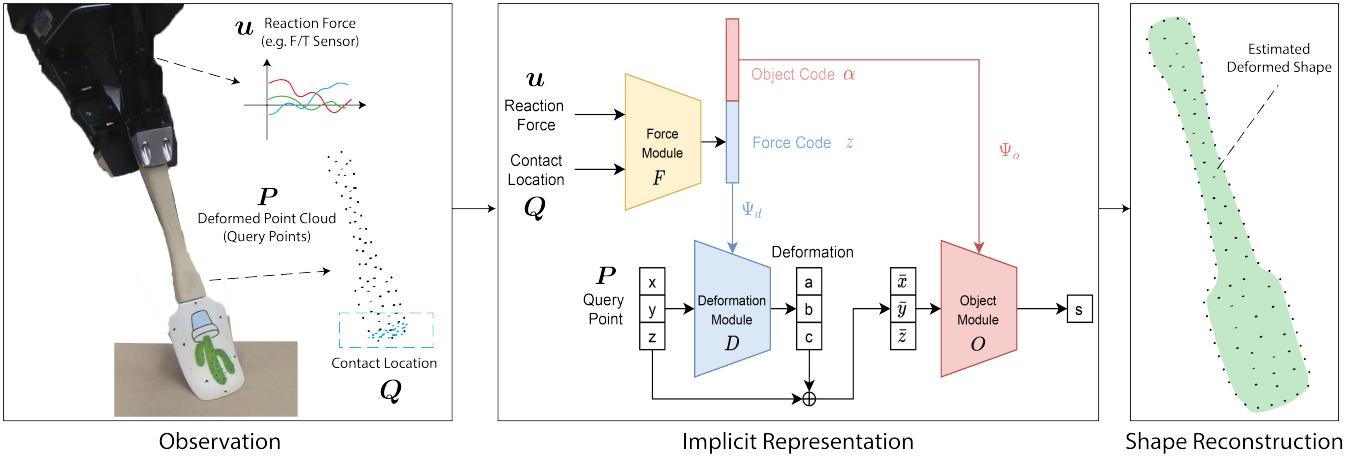


Fig. 1: **Representation Architecture:** The left panel depicts how visual data in the form of point clouds and tactile in the form of reaction forces may be collected in practical robotic settings. The middle panel depicts the network and how this information is processed to predict the implicit surface representation encoded as a signed-distance function. Finally, the right panel depicts the reconstruction of the estimated true surface given the external contacts and reaction force.

For convenience, we introduce a shorthand  $\mathbf{O}^i(\mathbf{x}) = \mathbf{O}(\mathbf{x}|\Psi_o(\boldsymbol{\alpha}_i))$  to denote the signed-distance field corresponding to the  $i^{th}$  object code. Also,  $\Omega$  is the 3D querying space,  $\Omega_0$  is the zero-level surfaces in the querying space,  $s^*$  is the ground truth signed-distance,  $\mathbf{n}^*$  is the ground truth normal, and  $\nabla$  denotes the spatial gradient.  $\delta$  is for clamping off-surface signed-distances to concentrate network capacity on details near the surface, as demonstrated in [9]. For latent space regularization, we impose a Gaussian prior on the object code  $L_{latent}(\boldsymbol{\alpha}) = \sum_{i=1}^N \|\frac{1}{l_i} \boldsymbol{\alpha}_i\|_2$  and weights of the network  $L_{hyper}(\boldsymbol{\theta}_o) = \sum_{i=1}^N \frac{1}{l_o} \|\boldsymbol{\theta}_o^i\|_2$  where  $l_o$  is the length of  $\boldsymbol{\theta}_o$ . After the pre-training, we fix the object codes as constant vectors, but continue updating the object module and the hyper-network with low learning rate( $\sim 1e-8$ ) for additional refinement.

### B. Deformed Object Representation

The Force Module  $\mathbf{F}$  is an encoder that summarizes the contact locations and reaction force  $(\mathbf{Q}, \mathbf{u})$  into a force code  $\mathbf{z} = \mathbf{F}(\mathbf{Q}, \mathbf{u})$ . We assume that the contact location set  $\mathbf{Q}$  is given as a subset of the nominal point cloud ( $\mathbf{Q} \subset \mathbf{P}$ ) and the reaction force  $\mathbf{u} \in \mathbb{R}^3$  is directly measured at the robot's wrist. Point clouds, including the contact location set  $\mathbf{Q}$ , are unordered and variable in length. Our contact location encoder utilizes the PointNet architecture [13].

We define the deformation field as a 3D vector field that pushes a deformed object back to its original (nominal) shape. VIRDO recovers the signed-distance field of the nominal shape by adding the deformation field to the SDF of the deformed shape. We highlight that VIRDO has learned to focus the deformation field around the boundary of the object with magnitude reflecting the amount of deformation.

VIRDO represents the deformation field as  $\mathbf{D}(\mathbf{x}|\Psi_d(\mathbf{z}, \boldsymbol{\alpha}))$ , where the deformation module  $\mathbf{D}$  shares the same structure as the object module  $\mathbf{O}$  with parameters  $\boldsymbol{\theta}_d$  predicted by the hyper-network  $\Psi_d$ . We highlight that  $\boldsymbol{\theta}_d$  is conditioned on the latent code pair  $(\mathbf{z}, \boldsymbol{\alpha}) \in \mathbb{R}^{l+m}$  to capture the underlying object-specific deformation behavior.

This results in  $\mathbf{D}$  predicting different deformation fields for different objects despite similar contact locations and reaction force measurements. This is desirable because objects may be geometrically similar but deform differently due to varying material properties. We will demonstrate examples of this in Sec. III-B.

Using the shorthand  $\mathbf{D}\Psi_d(\mathbf{x}) = \mathbf{D}(\mathbf{x}|\Psi_d(\mathbf{z}, \boldsymbol{\alpha}))$ , we can relate the signed-distance fields of the deformed and nominal object via  $s = \text{SDF}(\mathbf{x}) = \mathbf{O}_{\Phi_o}(\mathbf{x} + \mathbf{D}\Psi_d(\mathbf{x}))$ . We note that the deformed point cloud  $\mathbf{p} \in \mathbf{P}$  satisfies  $\text{SDF}(\mathbf{p}) = \mathbf{O}_{\Phi_o}(\mathbf{p} + \mathbf{D}\Psi_d(\mathbf{p})) = 0$ . To learn this mapping, we solve the optimization problem  $\underset{\boldsymbol{\theta}_d}{\text{argmin}} f_c(\mathbf{P})$ , where

$$f_c(\mathbf{P}) = \left( \text{CD}(\mathbf{P} + \mathbf{D}\Psi_d(\mathbf{P}), \bar{\mathbf{P}}^*) + \lambda_c \frac{1}{\bar{\mathbf{P}}} \sum_{\mathbf{p} \in \mathbf{P}} \|\mathbf{D}\Psi_d(\mathbf{p})\|_2 \right) \quad (2)$$

where  $\bar{\mathbf{P}}^*$  is the true nominal point cloud of the same length as  $\mathbf{P}$ ,  $\lambda_c$  is a weighting for the minimal correction prior similar to [14] and CD is the Chamfer Distance measure between two point clouds.

To learn the full model, we train the deformation, force, and pre-trained object modules end-to-end using the loss function:

$$L_{deformed} = f_c(\{\mathbf{x} | \mathbf{x} \in \Omega_0\}) + \lambda_1 L_{sdf} + \lambda_3 L_{hyper} + \lambda_4 L_{latent}. \quad (3)$$

The first term solves the optimization in Eq. 2 with on-surface points. The second loss term  $L_{sdf}$  couples the deformation field to the signed-distance field of the object module:

$$L_{sdf} = \sum_{i=1}^M \left( \sum_{\mathbf{x} \in \Omega} |\text{clamp}(\mathbf{O}_{\Phi_o}(\mathbf{x} + \mathbf{D}\Psi_d(\mathbf{x}), \delta) - \text{clamp}(s^*, \delta)| \right. \\ \left. + \sum_{\mathbf{x} \in \Omega_0} \lambda_n (1 - \langle \nabla \mathbf{O}^i(\mathbf{x}), \mathbf{n}^* \rangle) \right) \quad (4)$$

where  $M$  is the total number of deformed objects. The  $L_{hyper}(\mathbf{x})$  and  $L_{latent}(\mathbf{x})$  are the Gaussian prior on the latent code and the network parameters. Since we are updating two

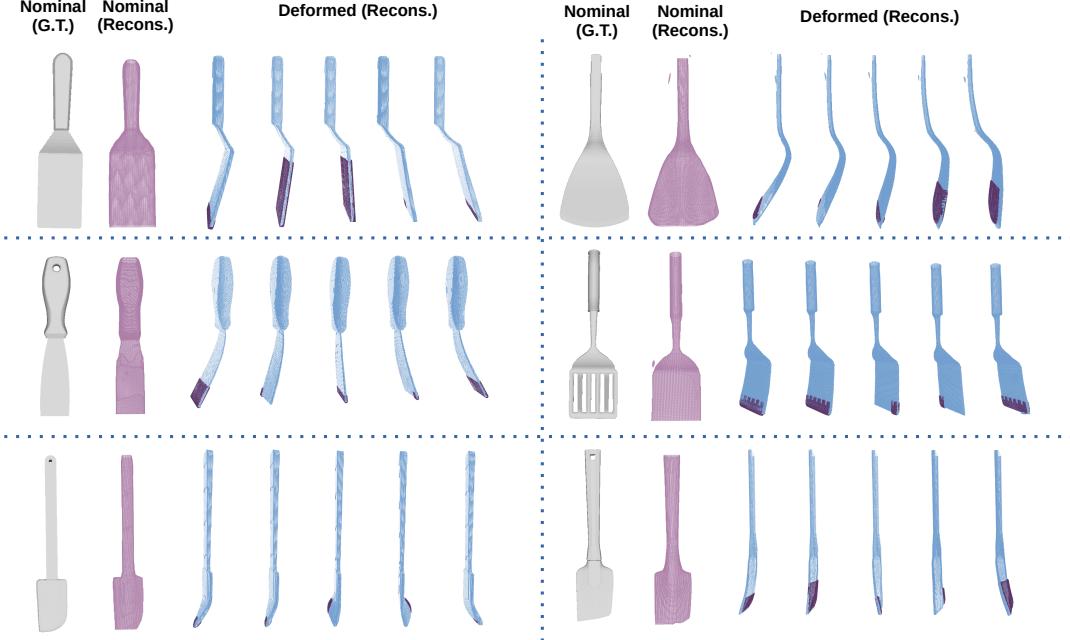


Fig. 2: **Reconstruction Results:** Reconstructions of multiple nominal shapes and their deformations, learned simultaneously by VIRDO. Marching Cube algorithm is used for the reconstruction, where we highlighted the contact location as purple region.

hyper-networks,  $L_{\text{hyper}}$  is the weighted sum of  $\Psi_o$  and  $\Psi_d$  and the  $L_{\text{latent}}$  is  $\sum_{i=1}^N \|\frac{1}{T}\mathbf{z}\|_2$ .

### III. EXPERIMENTS AND RESULTS

#### A. Data Preparation

In total, we generated 6 objects each with 24 unique boundary conditions using MATLAB PDE toolbox. The 3D meshes were collected from open-sourced 3D model repository. For the training, we normalized the point cloud with the geometric center at [0,0,0].

#### B. Representing Known Deformable Shapes

The average reconstruction accuracy is  $0.3474 \times 10^3$  in Chamfer Distance(CD). CD is measured between reconstructions and query points unseen during training in average, where we utilized Marching Cube algorithm [15] for the reconstruction. We emphasize that only one neural network model was used for the entire data-set.

#### C. Deformation Field Inference

We test the model’s ability to infer a deformation field given reaction force, partial pointcloud, and object code, seen from the training. Here, we infer the contact feature to estimate deformation. First, we randomly initialize the contact feature from  $N(0, 0.01^2)$ . Then, we update the feature with an L1 loss which only consumes a partial zero-level set:  $L_{\text{infer}} = \sum_{\mathbf{x} \in \Omega_o} |\text{clamp}(\mathbf{O}(\mathbf{x}), \delta)|$ . The loss encourages VIRDO to update the contact feature by minimizing the mismatch between the initial guess and the partial observation. Fig. 3 is a partial pointcloud where the handle and the tip are occluded, rendered in simulation with a single pinhole camera. At epoch 0, the model already makes deformation field fairly close to the ground truth. This shows VIRDO’s ability to perform state estimation when the vision is missing. As the gradient descent progresses on the contact feature, the

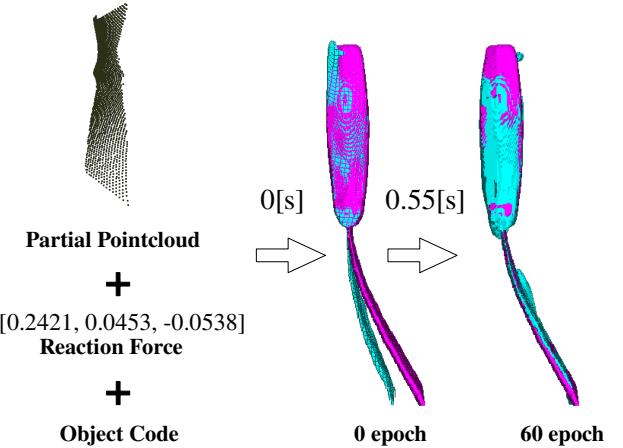


Fig. 3: **Inference:** Reconstructions with inferred deformation field (cyan), ground truth deformed object (magenta)

estimated deformation converges towards the ground truth. We note that only the on-surface points are used for this experiment, since an RGBD camera would feasibly only give on-surface points in real-world experiments; however, it is also possible boost the inference performance by collecting off-surface samples along the camera ray and augment the partial observation.

#### D. Generalization and Code Interpolation/Extrapolation

In this section, we evaluate VIRDO’s ability to generalize to unseen contact formations. This functionality is important for robotic applications given the wide variety of contact interactions. To this end, we use the model from Sec. III-B (trained on 144 deformations) and evaluate the reconstruction accuracy of 6 unseen contact formations for the object

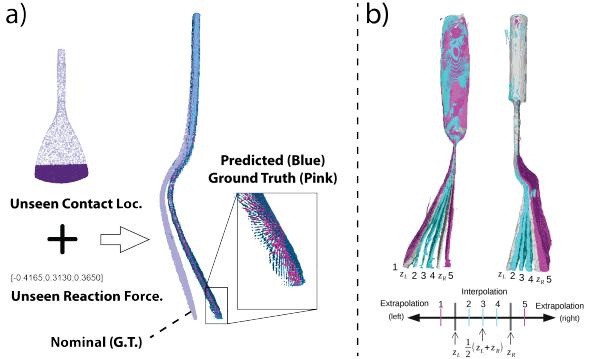


Fig. 4: a) **Generalization:** Example of shape estimation given unseen contact formation during training. b) **Latent Code Interpolation and Extrapolation:** Two trained force codes  $z_l$  and  $z_r$  are interpolated and extrapolated evenly as shown where 1-5 indicate tested force code and corresponding reconstruction result.

depicted in Fig.4a). To better understand the generalization properties of VIRDO, we interpolate and extrapolate in latent force code space. This task evaluates the continuity and semantic meaning of the latent space; e.g. whether the reconstructions are similar to direct inter/extrapolation of two deformations. Given an object, we pick  $z_l$  and  $z_r$  among the successfully trained force codes;  $z_l$  bears the maximum deflection to the left (or  $-x$  direction) and  $z_r$  generates moderate deflections in a random direction. We then linearly interpolate and extrapolate the 32 dimensional force code. Fig.4b) shows the resulting reconstructions. We note the interpolations and extrapolations are smooth, even, and continuous. This suggests a well-formed latent force space and explains the effectiveness of the model in generalizing to unseen contact formations.

We note that the reconstruction for extrapolation (1) is omitted from the object on the right in Fig. 4. This is because of a failure case resulting in a noisy and poor quality reconstruction. We found that the extrapolation results can be inaccurate if it exceeds the bound of maximum deflections seen in the training data.

#### IV. DISCUSSION & LIMITATIONS

The fundamental principle driving VIRDO is the ability to learn deformation fields informed by visuo-tactile sensing. VIRDO is the first learned implicit method to integrate tactile and visual feedback while modeling object deformations subject to external contacts. The representation has arbitrary resolution and is cheap to evaluate for point-wise sampling. Additionally, the latent code is well-behaved and can be used for inference.

#### V. ACKNOWLEDGEMENT

This research is partly supported by Robotics at Google.

#### REFERENCES

- [1] A. Billard and D. Kragic, “Trends and challenges in robot manipulation,” *Science*, vol. 364, no. 6446, 2019.
- [2] H. Yin, A. Varava, and D. Kragic, “Modeling, learning, perception, and control methods for deformable object manipulation,” *Science Robotics*, vol. 6, no. 54, 2021.
- [3] J. Sanchez, J.-A. Corrales, B.-C. Bouzgarrou, and Y. Mezouar, “Robotic manipulation and sensing of deformable objects in domestic and industrial applications: A survey,” *The International Journal of Robotics Research*, vol. 37, no. 7, pp. 688–716, 2018.
- [4] V. E. Arriola-Rios, P. Guler, F. Ficuciello, D. Kragic, B. Siciliano, and J. L. Wyatt, “Modeling of deformable objects for robotic manipulation: A tutorial and review,” *Frontiers in Robotics and AI*, vol. 7, p. 82, 2020.
- [5] Y. C. Hou, K. S. M. Sahari, and D. N. T. How, “A review on modeling of flexible deformable object for dexterous robotic manipulation,” *International Journal of Advanced Robotic Systems*, vol. 16, no. 3, p. 1729881419848894, 2019.
- [6] T. Hermans, F. Li, J. M. Rehg, and A. F. Bobick, “Learning contact locations for pushing and orienting unknown objects,” in *2013 13th IEEE-RAS international conference on humanoid robots (humanoids)*, IEEE, 2013, pp. 435–442.
- [7] M. Sharma and O. Kroemer, “Relational learning for skill preconditions,” *arXiv preprint arXiv:2012.01693*, 2020.
- [8] Q. V. Le, D. Kamm, A. F. Kara, and A. Y. Ng, “Learning to grasp objects with multiple contact points,” in *2010 IEEE International Conference on Robotics and Automation*, IEEE, 2010, pp. 5062–5069.
- [9] J. J. Park, P. Florence, J. Straub, R. Newcombe, and S. Lovegrove, “Deepsdf: Learning continuous signed distance functions for shape representation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 165–174.
- [10] V. Sitzmann, J. Martel, A. Bergman, D. Lindell, and G. Wetzstein, “Implicit neural representations with periodic activation functions,” *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [11] M. Tancik, P. P. Srinivasan, B. Mildenhall, *et al.*, “Fourier features let networks learn high frequency functions in low dimensional domains,” *arXiv preprint arXiv:2006.10739*, 2020.
- [12] Y. Deng, J. Yang, and X. Tong, “Deformed implicit field: Modeling 3d shapes with learned dense correspondence,” *arXiv preprint arXiv:2011.13650*, 2020.
- [13] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, “Pointnet: Deep learning on point sets for 3d classification and segmentation,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 652–660.
- [14] Y. Deng, J. Yang, and X. Tong, “Deformed implicit field: Modeling 3d shapes with learned dense correspondence,” *CoRR*, vol. abs/2011.13650, 2020. arXiv: 2011 . 13650. [Online]. Available: <https://arxiv.org/abs/2011.13650>.

- [15] W. E. Lorensen and H. E. Cline, “Marching cubes: A high resolution 3d surface construction algorithm,” *ACM siggraph computer graphics*, vol. 21, no. 4, pp. 163–169, 1987.