



HAL
open science

Unsupervised Learning of Category-Specific Symmetric 3D Keypoints from Point Sets

Clara Fernandez-Labrador, Ajad Chhatkuli, Danda Pani Paudel, Jose J
Guerrero, Cédric Demonceaux, Luc Van Gool

► **To cite this version:**

Clara Fernandez-Labrador, Ajad Chhatkuli, Danda Pani Paudel, Jose J Guerrero, Cédric Demonceaux, et al.. Unsupervised Learning of Category-Specific Symmetric 3D Keypoints from Point Sets. 16TH EUROPEAN CONFERENCE ON COMPUTER VISION, ECCV 2020, Aug 2020, Glasgow, United Kingdom. hal-02916278

HAL Id: hal-02916278

<https://hal.science/hal-02916278>

Submitted on 17 Aug 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Unsupervised Learning of Category-Specific Symmetric 3D Keypoints from Point Sets

Clara Fernandez-Labrador^{1,2,3}, Ajad Chhatkuli³, Danda Pani Paudel³,
Jose J. Guerrero¹, Cédric Demonceaux², and Luc Van Gool^{3,4}

¹ I3A, University of Zaragoza, Spain

² VIBOT ERL CNRS 6000, ImViA, University of Bourgogne Franche-Comté, France

³ Computer Vision Lab, ETH Zürich, Switzerland

⁴ VISICS, ESAT/PSI, KU Leuven, Belgium

{cfernandez, josechu.guerrero}@unizar.es,

cedric.demonceaux@u-bourgogne.fr,

{ajad.chhatkuli,paudel,vangool}@vision.ee.ethz.ch

Abstract. Automatic discovery of category-specific 3D keypoints from a collection of objects of some category is a challenging problem. One reason is that not all objects in a category necessarily have the same semantic parts. The level of difficulty adds up further when objects are represented by 3D point clouds, with variations in shape and unknown coordinate frames. We define keypoints to be category-specific, if they meaningfully represent objects' shape and their correspondences can be simply established order-wise across all objects. This paper aims at learning category-specific 3D keypoints, in an unsupervised manner, using a collection of misaligned 3D point clouds of objects from an unknown category. In order to do so, we model shapes defined by the keypoints, within a category, using the symmetric linear basis shapes without assuming the plane of symmetry to be known. The usage of symmetry prior leads us to learn stable keypoints suitable for higher misalignments. To the best of our knowledge, this is the first work on learning such keypoints directly from 3D point clouds. Using categories from four benchmark datasets, we demonstrate the quality of our learned keypoints by quantitative and qualitative evaluations. Our experiments also show that the keypoints discovered by our method are geometrically and semantically consistent.



Fig. 1: **Category-specific 3D Keypoints.** The predicted keypoints follow the symmetric linear shape basis prior modeling all instances in a category under a common framework. They not only are consistent across different instances, but also are ordered and correspond to semantically meaningful locations.

1 Introduction

A set of keypoints representing any object is historically of large interest for geometric reasoning, due to their simplicity and ease to handle them. In fact, keypoints-based methods have been crucial to the success of many vision applications. A few examples include; 3D reconstruction [1–3], registration [4–7], human body pose [8–11], recognition [12, 13], and generation [14, 15]. That being said, many keypoints are defined manually, while considering their semantic locations such as facial landmarks and human body joints, to serve and simplify the problem at hand. To further benefit from their widespread utility, several attempts have been made on learning to detect keypoints [16–20], as well as on automatically discovering them [21–24]. In this regard, the task of learning to detect keypoints from several supervision examples, has achieved many successes. However, discovering them automatically from unlabeled data –such that they meaningfully represent shapes and semantics– so as to have a similar utility as those of manually defined, has received only limited attention due to its difficulty.

Therefore, it must not come as a surprise that keypoints defined in 3D space are preferred for geometric reasoning, where the objects of interest also reside. For given 3D keypoints, their counterparts in 2D images can be associated by merely using camera projection models [25–27]. By the above reasoning, it is natural to seek for 3D keypoints. In fact, one may infer 3D keypoints only using 2D images, 3D structures, or 2D-3D pairs. Learning such keypoints directly from 2D requires image correspondences and their poses to be known, along with the camera projection model. Generally, the difficulties of estimating camera pose and of establishing correspondences are avoided by using 2D-3D aligned pairs, where the primary interest is to infer 3D keypoints from 2D images. In this regard, a notable work [24] uses 3D structure and multiple associated 2D images with known poses. However in practice, multiple images together with aligned 3D structure may not always be available. In that context, one is left with the task of directly learning keypoints from 3D structures (or 2D images). In this work, we are interested on learning keypoints using only 3D structures. In fact, 3D structures with keypoints suffice for several applications including, registration [28], shape completion [29], and shape modeling [30]; without requiring their 2D counterparts.

When 3D objects go through shape variations, either because of being deformable or when two different objects of one category are compared, keypoints are desired to be consistent for meaningful geometric reasoning. Recall the examples of semantic keypoints such as facial landmarks and body joints. To serve a similar purpose, *can we automatically find keypoints that are consistent over inter-subject shape variations and intra-subject deformations in a category?* This is the primary question that we are interested to answer in this paper. Furthermore, we wish to discover such keypoints directly from 3D point clouds, in an unsupervised manner. We call these keypoints “category-specific”, which are expected to meaningfully represent objects’ shape and offer their correspondence order-wise across all objects. More formally, the desired properties of category-specific

keypoints are: i) generalizability over different shape instances and alignments in a category, ii) one-to-one ordered correspondences and semantic consistency, iii) representative of the shape as well as the category while preserving shape symmetry. These properties not only make the representation meaningful, but also tend to enhance the usefulness of keypoints. Learning category-specific keypoints on point clouds, however, is a challenging problem because not all the object parts are always present in a category. The challenges are exacerbated when the practical cases of misaligned data and unsupervised learning are considered. Related works do not address all these problems, but instead opt for; dropping category-specificity and using aligned data [23], employing manual supervision on 2D images [17], using aligned 3D and multiple 2D images with known pose [24]. The latter method achieves category-specificity without explicitly reasoning on the shapes.

In this paper, we show that the category-specific keypoints with the listed properties can be learned unsupervised by modeling them with non-rigidity based on linear basis shapes. We further model non-rigidity using reflective symmetry, with an instance-wise symmetry when available. For categories without instance-wise symmetry we propose the use of symmetric linear basis shapes in order to better model, what we define as symmetric deformation spaces, e.g., a human body deformations. This allows us to better constrain the pose and the shape coefficients prediction. Our proposed learning method does not assume aligned shapes [24], pre-computed basis shapes [17] or known planes of symmetry [31] and all quantities are learned in an end-to-end manner. Our symmetry modeling is powerful and more flexible compared to that of previous NRSfM methods [31, 32]. We achieve this by considering the shape basis for a category and the reflective plane of symmetry as the neural network weight variables, optimized during the training process. At inference time, the network predicts the basis coefficients and the pose in order to estimate the instance-specific keypoints. Fig. 2 shows the basic overview of our training strategy. Note that we do not require the Siamese-like architecture as in [4, 23]. Using multiple categories from four benchmark datasets, we evaluate the quality of our learned keypoints both quantitatively and with qualitative visualization. Our experiments also show that the keypoints discovered by our method are geometrically and semantically consistent, which are measured respectively by intra-category registration and semantic part-wise assignments. We further show that symmetric basis shapes can be used to model symmetric deformation space of categories such as human body.

2 Related Work

Category-specific keypoints on objects have been extensively used in NRSfM methods, however, only few methods have tackled the problem of estimating them. In terms of the outcome, our work is closest to [24], which learns category-specific 3D keypoints by solving an auxiliary task of rigid registration between multiple renders of the same shape and by considering the category instances to

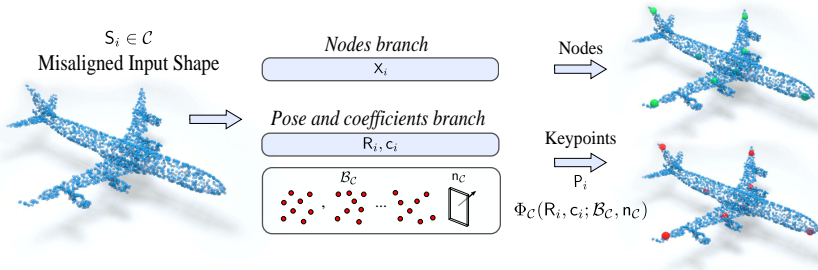


Fig. 2: **Overview:** the proposed learning strategy consists of two main branches that predict the instance-specific parameters while the common category parameters get optimized as network weights. Refer to Sec. 3 for the modeling, Sec. 4 for learning and supplementary material for more details.

be pre-aligned. Although the method shows promising results on 2D and 3D, it does so without explicitly modeling the shapes. Consequently, it requires renders of different instances to be pre-aligned to reason on keypoint correspondences between instances. A similar task is also solved in [17] for 6-degrees of freedom (DoF) estimation which uses low-rank shape prior to condition keypoints in 3D. Although, the low-rank shape modeling is a powerful tool, [17] requires supervision for heatmap prediction and relies on aligned shapes and pre-computed shape basis. [33] also predicts keypoints for categories with low-rank shape prior but the method is again trained on fully supervised manner. Moreover, all of the mentioned methods learn keypoints on images as heatmaps and thereafter lift them to 3D. Shape modeling of category shape instances has been widely explored in NRSfM works. Linear low-rank shape basis [2, 34, 35], low-rank trajectory basis [36], isometry or piece-wise rigidity [37, 38] are some of the different methods used for NRSfM. Recently, a few number of works have used low-rank shape basis in order to devise learned methods [1, 31, 33, 39]. Another useful tool in modeling shape category is the reflective symmetry, which is also directly related to the object pose. Although [32] showed that the low-rank shape basis can be formulated with unknown reflective symmetry, its adaptation to learned NRSfM methods is however, not trivial. Recent methods, in fact, assume that the plane of symmetry is one among a few known planes [31]. Moreover, none of the methods formulate symmetry applicable for non-rigidly deforming objects such as a human body.

While shape modeling is a key aspect of our work, another challenge is to infer ordered keypoints by learning on unordered point sets. While several advances have been made on deep neural networks for point sets [40–42], current achievement of learning on ordered structure such as images dwarfs those of learning on point sets. A related work learns to predict 3D keypoints unsupervised by again solving the auxiliary task of correctly estimating rotations in a Siamese architecture [43]. The keypoint prediction is done without order by pooling features of certain point neighborhoods. Another previous work [4] proposes learning point

features for matching, again using alignment as the auxiliary task. Matching such keypoints across shapes is not an easy task as the keypoints are not predicted in any order. In the following sections we show how one can model shape instances using the low-rank symmetric shape basis and use the shape modeling to predict ordered category-specific keypoints.

3 Background and Theory

Notations. We represent sets and matrices with special Latin characters (e.g., \mathcal{V}) or bold Latin characters (e.g., \mathbf{V}). Lower or uppercase normal fonts, e.g., K denote scalars. Lowercase bold Latin letters represent vectors as in \mathbf{v} . We use lowercase Latin letters to represent indices (e.g., i). Uppercase Greek letters represent mappings or functions (e.g., Π). Finally the operator $\text{vec}(\cdot)$ denotes the vectorize operation of a matrix.

3.1 Category-specific Shape and Keypoints

We represent shapes as sets of point coordinates, or point clouds, defined as an unordered set of points $\mathbf{S} = \{\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_M\}$, $\mathbf{s}_j \in \mathbb{R}^3$, $j \in \{1, 2, \dots, M\}$. The set of all such shapes in a category defines the category shape space \mathcal{C} . We write a particular i -th category-specific shape instance in \mathcal{C} as \mathbf{S}_i . For convenience, we will use the terms category-specific shape and shape interchangeably. The category shape space \mathcal{C} can be anything from a set of discrete shapes to a smooth manifold spanned by a deformation function $\Psi_{\mathcal{C}}$, whose co-domain consists of only category-specific shapes. The focus of the work is on learning meaningful 3D keypoints from the point set representation of \mathbf{S}_i . To that end, this section defines category-specific keypoints and develops their modeling.

Category-specific keypoints. We represent category-specific keypoints of a shape \mathbf{S}_i as a sparse tuple of points, $\mathbf{P}_i = (\mathbf{p}_{i1}, \mathbf{p}_{i2}, \dots, \mathbf{p}_{iN})$, $\mathbf{p}_{ij} \in \mathbb{R}^3$, $j \in \{1, 2, \dots, N\}$. Unlike the shape, its keypoints are represented as ordered points, forming a totally ordered set. Our objective is to learn a mapping $\Pi_{\mathcal{C}} : \mathbf{S}_i \rightarrow \mathbf{P}_i$ in order to obtain the category-specific keypoints from an input shape \mathbf{S}_i for a category shape space \mathcal{C} . Although not completely unambiguous, we can define the category-specific keypoints using the properties listed in Section 1. In mathematical notations they are:

- (i) Generalization: $\Pi_{\mathcal{C}}(\mathbf{S}_i) = \mathbf{P}_i$, $\forall \mathbf{S}_i \in \mathcal{C}$.
- (ii) Corresponding points and semantic consistency: $\mathbf{p}_{aj} \Leftrightarrow \mathbf{p}_{bj}$, $\mathbf{S}_a, \mathbf{S}_b \in \mathcal{C}$. For any $\mathbf{S}_a, \mathbf{S}_b \in \mathcal{C}$, \mathbf{p}_{aj} and \mathbf{p}_{bj} have the same semantics.
- (iii) Representative-ness: $\text{vol}(\mathbf{S}_i) = \text{vol}(\mathbf{P}_i)$ and $\mathbf{p}_{ij} \in \mathbf{S}_i$, where $\text{vol}(\cdot)$ is the Volume operator. If $\mathbf{S}_i \in \mathcal{C}$ has a reflective symmetry, \mathbf{P}_i should have the same symmetry.

3.2 Category-specific Shapes as Instances of Non-Rigidity

Several recent works have modeled shapes in a category as instances of non-rigid deformations [1, 31, 33, 39]. The motivation lies in the fact that such shapes often share similarities to a large extent. Consequently there, likely, exists a deformation function $\Psi_C : S_T \rightarrow S_i$, which can map a global shape property S_T to a category shape instance S_i . However, we argue that modeling Ψ_C is not trivial and in fact a convenient representation of Ψ_C may not exist in many cases. This observation, in fact, is what makes the dense Non-Rigid Structure-from-Motion (NRSfM) [31] so challenging. On the other hand, one can imagine a deformation function $\Phi_C : P_T \rightarrow P_i$, going from a global keypoints property P_T to the category-specific keypoints P_i . The deformation function Φ_C thus satisfies: $p_{ij} \in \Phi_C$ implies $p_{ij} \in \Psi_C$ and effectively, $\Phi_C \subset \Psi_C$, if the set order in P_i is ignored. Unlike Ψ_C , the deformation function Φ_C may be simple enough to model and use for estimating the category-specific keypoints P_i . We therefore, choose to seek the non-rigidity modeling in the space of keypoints $\mathcal{P} = \{P_1, P_2, \dots, P_L\}$, which functions as an abstraction of the space \mathcal{C} . Non-rigidity can be used to define the prediction function Π_C as below:

$$\Pi_C(S_i; \theta) = \Phi_C(r_i; \theta) = P_i, \quad (1)$$

where θ denotes the constant function parameters of Π_C and r_i is the predicted instance specific vector parameter. In our problem, we want to learn θ from the example shapes in \mathcal{C} without using the ground-truth labels, supervised by Φ_C . In the NRSfM literature, two most common approaches of modeling shapes as non-rigid deformations are the shape basis or low-rank shape prior [2, 34–36] and the isometric prior [37, 38]. In this paper, we investigate the modeling using a particular form of low-rank shape prior, i.e., the symmetric shape basis.

3.3 Low-Rank Non-rigid Representation of Keypoints

The NRSfM approach of low-rank shape basis comes as a natural extension of the rigid orthographic factorization prior [44] and was introduced by Bregler et al. [34]. The key idea is that a large number of object deformations can be explained by linearly combining a smaller number K of basis shapes at some pose. In the rigid case, this number is one, hence the rank is 3. In the non-rigid case, it can be higher, while the exact value depends on the complexity of the deformations. Consider F shape instances in \mathcal{C} and N points in each keypoints instance P_i . The following equation describes the projection with shape basis.

$$P_i = \Phi_C(r_i; \theta) = R_i \text{vec}(c_i \mathbf{1}^\top) \mathcal{B}_C. \quad (2)$$

where $\mathcal{B}_C = (\mathbf{B}_1, \dots, \mathbf{B}_K)$, $\mathcal{B}_C \in \mathbb{R}^{3K \times N}$ forms the low-rank shape basis. The rank is lower than the maximum possible rank of $3F$ or N for $3K < 3F$ or $3K < N$. The vector $c_i \in \mathbb{R}^K$ denotes the coefficients that linearly combines different basis for the keypoints instance i using the 3-vector of ones $\mathbf{1}$. Each keypoints instance is then completely parametrized by the basis \mathcal{B}_C and the coefficients c_i . Next, the

projection matrix $R_i \in \text{SO}_3$ is simply the rotation matrix for the shape instance i instead of a Stiefel matrix in the NRSfM problem.

Unlike in NRSfM, the problem of computing the category-specific keypoints, has P_i as unknown. Similar to NRSfM, the rest of the quantities in Eq. (2) – c_i , \mathcal{B}_C and R_i are also unknown. This fact makes our problem doubly hard. First the problem becomes more than just lifting the 2D keypoints to 3D and second, the order of keypoints present in the NRSfM measurements matrix is no longer available. Eq. (2) can be directly related to the deformation representation of Φ_C in Eq. (1), where θ includes the global parameters or basis \mathcal{B}_C and r_i includes the instance-wise pose R_i and coefficients c_i . A remark here is necessary regarding the number of shape basis K . Although we choose $K \approx N$, Eq. (2) still reduces the solution space for P_i . However, one issue remains related to the ambiguities of R_i when the number of basis shapes K is large. We propose to address the problem by also computing the reflective plane of symmetry of the category.

3.4 Modeling Symmetry with Non-Rigidity

Many object categories have shapes which exhibit a fixed reflective symmetry over the whole category. To discover and use symmetry, we consider two different priors: instance-wise symmetry and symmetric deformation space.

Instance-wise symmetry. Instance-wise reflective symmetry about a fixed plane is common in a large majority of rigid object categories in ShapeNet [45] and ModelNet [46] datasets. Instance-wise symmetry has been previously combined with the shape basis prior in NRSfM [32], however, a convenient representation for learning both the symmetry and the shapes have not been explored yet. A recent learning-based method [31] uses the symmetry prior by performing an exhaustive search over a few planes in order to predict symmetric dense non-rigid shapes. However, such a strategy may not work when the shapes are not perfectly aligned. Instance-wise symmetry can be included by re-writing Eq. (2) as follows:

$$P_{i\frac{1}{2}} = R_i \text{vec}(c_i \mathbf{1}^\top) \mathcal{B}_{C\frac{1}{2}}, \quad P_i = [P_{i\frac{1}{2}} \quad A_C P_{i\frac{1}{2}}], \quad (3)$$

where $P_{i\frac{1}{2}} \in \mathbb{R}^{3 \times N/2}$ represents one half of the category-specific keypoints. $P_{i\frac{1}{2}}$ is reflected using $A_C \in \mathbb{R}^{3 \times 3}$ and concatenated to obtain the final keypoints. Due to the exact instance-wise symmetry, we similarly can parametrize the basis as $\mathcal{B}_{C\frac{1}{2}} \in \mathbb{R}^{3K \times N/2}$ to denote the shape basis for the first half of the keypoints. The reflection operator A_C is parametrized by a normal vector $n_C \in \mathbb{R}^3$ of the plane of symmetry passing through the origin. The most important advantage going from Eq. (2) to Eq. (3) is the reduced dimensionality of the unknowns in \mathcal{B}_C as well as the additional second equality constraint of Eq. (3) which reduces the ambiguities in NRSfM [32].

Symmetric deformation space. In many non-rigid objects, shape instances are not symmetric. However, symmetry may still exist in the deformation space,

e.g., in a human body. Suppose that the shape instance $S_k \in \mathcal{C}$ has the reflective symmetry about $n_{\mathcal{C}}$, which allows us to define its two halves: $S_{k\frac{1}{2}}$ and $S'_{k\frac{1}{2}}$ and thus correspondingly for all shape instances.

Definition 1 (Symmetric deformation space). \mathcal{C} is a symmetric deformation space if for every half shape deformation instance $S_{i\frac{1}{2}}$, there exists any shape instance $S_j \in \mathcal{C}$ such that the $S'_{j\frac{1}{2}}$ is symmetric to $S_{i\frac{1}{2}}$.

The above definition also applies for the keypoints shape space \mathcal{P} . The instance-wise symmetric space is a specific case of the above. However, Eq. (3) cannot model the keypoints instances in the symmetric deformation space. We model such keypoints by introducing symmetric basis that can be weighted asymmetrically, thereby, obtaining the following:

$$P_i = R_i \left[\text{vec}(c_i \mathbf{1}^\top) \mathcal{B}_{\mathcal{C}\frac{1}{2}} \text{vec}(c'_i \mathbf{1}^\top) \mathcal{B}'_{\mathcal{C}\frac{1}{2}} \right] \quad (4)$$

where $\mathcal{B}'_{\mathcal{C}\frac{1}{2}}$ is obtained by reflecting $\mathcal{B}_{\mathcal{C}\frac{1}{2}}$ with $A_{\mathcal{C}}$ and $c'_i \in \mathbb{R}^K$ forms the coefficients for the second half of the basis. Although Eq. (4) increases the dimension of the unknowns in the coefficients over Eq. (2), the added modeling of the symmetry of the deformation space and the reduced dimensionality of the basis can improve the final keypoints estimate. This brings us to the following proposition.

Proposition 1. *Provided that $\mathcal{B}_{\mathcal{C}\frac{1}{2}}$ and $\mathcal{B}'_{\mathcal{C}\frac{1}{2}}$ are symmetric about a plane, Eq. (4) approximates a symmetric deformation space if the estimates of c_i and c'_i come from the same probabilistic distribution.*

Proof. The proof is straightforward and provided in the supplementary material.

As a consequence of Proposition 1, we can model keypoints in non-rigid symmetric objects with Eq. (4), while also tightly modeling the symmetry as long as we maintain the distribution of c and c' to be the same.

4 Learning Category-specific Keypoints

In this section, we use the modeling of $\Phi_{\mathcal{C}}$ to describe the unsupervised learning process of the category-specific keypoints. More precisely, we want to learn the function $\Pi_{\mathcal{C}} : S_i \rightarrow P_i$ as a neural network of parameters θ , using the supervisory signal from $\Phi_{\mathcal{C}}$. In regard to learning keypoints on pointsets, recent work [23] trains a Siamese network to predict order-agnostic keypoints stable to rotations for rigid objects [23]. We use a similar network architecture as in [23] that is based on PointNet [40] but we do not use the Siamese training. The overview of the network is shown in Fig. 2, whose input consists of a single shape S_i misaligned in SO_2 . This is reasonable since point clouds are usually aligned to the gravity direction. Our learning strategy consists of two main branches that will predict the instance-specific parameters while optimizing the network weights as category-specific parameters. The complete network architecture is provided in the supplementary material. We describe the branches below.

Nodes branch. This branch estimates nodes that are potentially category-specific keypoints but are not ordered. We denote them as $\mathbf{X}_i = \{\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iN}\}$, $\mathbf{x}_{ij} \in \mathbb{R}^3$ and $j \in \{1, 2, \dots, N\}$. Initially, a predefined number of nodes N are sampled from the input shape using the Farthest Point Sampling (FPS) and a local neighborhood of points is built for each node with point-to-node grouping [23,47], creating N clusters which are mean normalized inside the network. The number of nodes corresponds to the desired number of category-specific keypoints, and every point in S_i is associated with one of these nodes. The branch consists of two PointNet-like [40] networks followed by a kNN grouping layer that uses the initial sampled nodes to achieve hierarchical information aggregation. Finally, the local feature vectors are fed into a Multi-Layer Perceptron (MLP) that directly outputs the final nodes.

Pose and coefficients branch. The quantities \mathbf{R}_i and \mathbf{c}_i are learned and predicted by this branch. We use a single rotation angle to parametrize \mathbf{R}_i . The branch consists of an MLP that estimates the mentioned parameters. The output size will vary depending on whether we are interested in symmetric shape instances as in Eq.(3) or symmetric basis as in Eq. (4), the size being double in the latter.

Additional learnable parameters. Several quantities in Eq. (3) or (4) are constant for a category shape space \mathcal{C} . Such quantities need not be predicted instance-wise. We rather choose to optimize them as part of the network parameters θ . These parameters are $\mathcal{B}_{\mathcal{C}}$ of dimension $3K \times N$ and the plane of symmetry of dimension 3 with unitary constraint. We choose $5 \leq K \leq 10$. Depending upon the problem, alternate parametrization can be considered for $\mathbf{n}_{\mathcal{C}}$, e.g., Euler angles.

4.1 Training Losses

In order to adhere to the definitions of the category-specific keypoints introduced in Section 1 as well as our shape modeling, several key loss functions are used for the training process. We list these loss functions and define them below.

Chamfer loss with symmetry and non-Rigidity. Eq. (1) suggests that an ℓ_2 loss between the neural network predictions \mathbf{X}_i and the deformation function $\mathbf{P}_i = \Phi_{\mathcal{C}}(\mathbf{R}_i, \mathbf{c}_i; \mathcal{B}_{\mathcal{C}}, \mathbf{n}_{\mathcal{C}})$ should be enough to supervise the neural network $\Psi_{\mathcal{C}}$ in order to satisfy $\mathbf{P}_i = \mathbf{X}_i$. However, as confirmed by our evaluations in our model as well as in [23], the ℓ_2 loss does not converge as predicting the order of points adds much more complexity to the network. Alternatively, the Chamfer loss [48] does converge, minimizing the distance between each point \mathbf{x}_{ik} in the first set \mathbf{X}_i and its nearest neighbor \mathbf{p}_{ij} in the second set \mathbf{P}_i . We define it as follows:

$$\mathcal{L}_{chf} = \sum_{k=1}^N \min_{\mathbf{p}_{ij} \in \mathbf{P}_i} \|\mathbf{x}_{ik} - \mathbf{p}_{ij}\|_2^2 + \sum_{j=1}^N \min_{\mathbf{x}_{ik} \in \mathbf{X}_i} \|\mathbf{x}_{ik} - \mathbf{p}_{ij}\|_2^2, \quad (5)$$

Chamfer loss ensures that the learned keypoints follow a generalizable category-specific property – that they are a linear combination of common basis

learned specifically for the category. To additionally model symmetry, Eq. (3) or (4) is directly used in Eq. (5). Therefore, two different Chamfer losses are possible modeling two different types of symmetries. We further add unitary vector constraint on the global variable \mathbf{n}_c .

Coverage and inclusivity loss. The Chamfer loss, however, does not ensure that the keypoints follow the object shape. Loss terms that guarantee this property can be designed by having the following conditions: a) the keypoints cover the whole category shape (coverage loss), b) the keypoints are not far from the point cloud (inclusivity loss). The coverage loss can be defined as a Huber loss comparing the singular values Λ of the nodes \mathbf{X}_i with respect to those of the input shape \mathbf{S}_i . However, for the sake of efficiency, we reformulate it to compare the 3D bounding boxes defined by these set of points instead. This improves the training speed and based on our initial evaluations showed better accuracy. The final loss is as follows:

$$\mathcal{L}_{cov} = \|\text{vol}(\mathbf{X}_i) - \text{vol}(\mathbf{S}_i)\| \quad (6)$$

The inclusivity loss is formulated as a single side Chamfer loss [49] which penalizes nodes in \mathbf{X}_i that are far from the original shape \mathbf{S}_i , similarly to Eq. (5):

$$\mathcal{L}_{inc} = \sum_{k=1}^N \min_{\mathbf{s}_{ij} \in \mathbf{S}_i} \|\mathbf{x}_{ik} - \mathbf{s}_{ij}\|_2^2. \quad (7)$$

5 Experimental Results

We conduct experiments to evaluate the desired properties of the proposed category-specific keypoints and show their generalization over indoor/outdoor objects and rigid/non-rigid objects with four different datasets in total (Sec. 5.1 5.2). All these properties are also compared with a proposed baseline. We then evaluate the practical use of our keypoints for intra-category shapes registration (Sec. 5.3), analyzing the influence of symmetry. Additional qualitative results are shown in Fig. 1 and the supplementary material.

Datasets. We use four main datasets. These include ModelNet10 [46], ShapeNet parts [45], Dynamic FAUST [50] and Basel Face Model 2017 [51]. Since our method is category-specific, we require separate training data for each class in the datasets. For indoor rigid objects, we choose three categories from ModelNet10 [46], including chair, table and bed. Three outdoor rigid object categories, including airplane, car and motorbike, are evaluated from ShapeNet part dataset [45]. For non-rigid objects, we randomly choose a sequence of the Dynamic Faust [50] dataset that provides high-resolution 4D scans of human subjects in motion. Finally, we generate shape models of faces using the Basel Face Model 2017 [51] combining 50 different shapes and 20 different expressions. All the models are normalized so that the longest dimension lies in $[-1,1]$ and are randomly misaligned within a 45 degrees range on each axis.

Baseline. Since this is the first work computing category-specific keypoints from point sets, we construct our own baseline based on the very recent work USIP [23]. The method detects stable interest points in 3D point clouds under arbitrary transformations and is also unsupervised, which makes it the closest method for comparison. The USIP detector is not category-based, so we train the network per category to create the baseline. Additionally, we adapt the number of predicted keypoints so that the results are directly comparable to ours. While training some of the categories with this detector, specifically car and bed, we observe that predicting lower number of keypoints can lead to some degeneracies in the results [23], which is also mentioned in the paper.

Implementation details. Input point clouds of dimension 3×2000 are used. We implement the network in Pytorch [52] and train it end-to-end from scratch using the Adam optimizer [53]. The initial learning rate is 10^{-3} , which is exponentially decayed by a rate of 0.5 every 40 epochs. We use a batch size of 32 and train each model until convergence, around 200 epochs. The final loss function combines the three training losses mentioned above and are weighted as follows: $w_{chf} = w_{cov} = 1$ and $w_{inc} = 2$. For the ModelNet10 and ShapeNet parts datasets, we use the training and testing split provided by the authors. For the Basel Face Model 2017 we follow the common practice and split the 1000 generated faces in 85% training and 15% test. We use the same split strategy for the sequence used from the Dynamic Fuaust dataset, which is ‘50009_jiggle_on_toes’ and contains 244 examples.

5.1 Desired Properties Analysis

As described in Sections 1 and 3, the category-specific keypoints satisfy certain desired properties. We propose six different metrics to evaluate the properties which are also used for comparison against the baseline. All the results are presented in Table 1, and are averaged across the test samples.

Coverage: According to property iii), we seek keypoints that are representative of each instance shape as well as of the category itself. To measure it, we calculate the percentage of the input shape 3D bounding box covered by the keypoints. On average, we achieve a 29.4% more of coverage with respect to our baseline.

Model Error: This metric refers to the Chamfer distance between the estimated nodes and the learned category-specific keypoints, normalized by the model’s scale. We get a very low error, meaning that the network satisfactorily manages to generalize, describing the nodes with the symmetric non-rigidity modeling (Property i) and iii)).

Correspondence: We measure the ability of the model to find the same set of keypoints on different instances of a given category (Property ii)). We first use K-means clustering to show this property in comparison to our baseline in Fig. S3, rest of categories are provided in the supplementary material. One can see at a glance how our keypoints are neatly clustered, whereas the ones of our baseline USIP get mixed. Numerically, we show the % occurrence of each specific keypoint belonging to the same cluster across instances. This property can just be evaluated when the interest points are ordered, therefore the USIP keypoints

are ordered per instance according to the performed clustering. Obeying the low-rank non-rigidity prior, this property is fully satisfied by our learned keypoints, meaning that they are consistent across shapes of the same class, in contrast to the baseline keypoints.

Inclusivity: We measure the percentage of keypoints that lie inside the point cloud (of scale 2) within a certain threshold chosen as 0.015, which also proves property iii). This is the only metric in which our method doesn’t outperform the baseline in all the cases.

Symmetry: The metric shows the angle error of the predicted reflective plane of symmetry. We obtain highly accurate prediction for rigid categories. In the non-rigid human body shape however, the ambiguities are severe. Despite that, the learned keypoints satisfy the other properties, particularly that of semantic correspondence. Both of these facts can be observed in Fig. 1.

Definition: As soon as the number of predefined nodes is increased, the predicted keypoints get grouped into clusters to follow the low-rank constrain. By detecting those clusters, we obtain the number of points in the keypoints.

| <i>Category</i> | <i>Coverage</i> | <i>ModelErr</i> | <i>Correspondence</i> | <i>Inclusivity</i> | <i>SymErr</i> | <i>Definition</i> |
|-----------------|-----------------|-----------------|-----------------------|--------------------|---------------|-------------------|
| | % | % | % | % | ° | |
| chair | 88.83 | 0.72 | 100 | 90.46 | 0.40 | 10 |
| table | 93.33 | 0.99 | 100 | 93.38 | 2.86 | 6 |
| bed | 80.31 | 0.94 | 100 | 95.33 | 0.13 | 6 |
| airplane | 89.15 | 0.64 | 100 | 96.35 | 0.20 | 8 |
| car | 92.39 | 0.72 | 100 | 97.77 | 2.21 | 8 |
| motorbike | 96.13 | 0.79 | 100 | 90.53 | 1.42 | 8 |
| human body | 85.59 | 0.72 | 100 | 97.73 | 33.30 | 11 |
| faces | 97.93 | 0.41 | 100 | 100 | 0.15 | 9 |
| chair | 79.73 | – | 55.6 | 98.50 | – | 10 |
| table | 79.72 | – | 34.5 | 99.83 | – | 6 |
| bed | 42.18 | – | 49.33 | 70.00 | – | 6 |
| airplane | 69.24 | – | 47.5 | 87.13 | – | 8 |
| car | 26.87 | – | 32.18 | 74.0 | – | 8 |
| motorbike | 75.29 | – | 48.14 | 84.57 | – | 8 |
| human body | 72.66 | – | 50.45 | 100 | – | 11 |
| faces | 42.98 | – | 30.11 | 100 | – | 9 |

Table 1: **Properties Analysis:** First (top) and second (bottom) block of the table present our and baseline results [23] respectively. For coverage, correspondence and inclusivity *bigger is better*, whereas for model error and symmetry error *smaller is better*. We not only demonstrate the desired properties of our keypoints, but also show the generalization of our method over indoor/outdoor, rigid/non-rigid objects and over four different datasets. Best results are in bold.

5.2 Semantic Consistency

We make use of the ShapeNet part dataset to show the semantic consistency of the proposed keypoints. Following the low-rank non-rigidity prior, the keypoints correspond to geometrically meaningful locations. The idea of the experiment is to measure keypoint-semantic relationship for every keypoint across

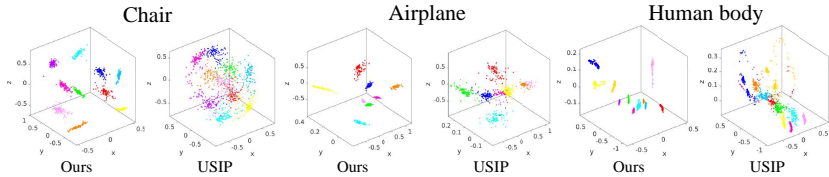


Fig. 3: **Keypoints correspondence across instances.** We cluster the keypoints predicted for all the instances of a category to show their geometric consistency. Note how our keypoints get neatly clustered creating a general 3D shape template.

instances of the category. The results are presented in Fig. 4 as covariance matrices, along with qualitative result per category for our method. On average, the proposed keypoints have consistency of 93% across instances, which means that our category-specific keypoints preserve the semantic relation across instances, despite the large appearance differences and intra-category variability. The same experiment is performed for our baseline and presented in the bottom part of Fig. 4. Again, we follow the approach mentioned in Sec. 5.1 for matching keypoints. Here, the degeneracy in the case of the car causes all the keypoints to approach the object centroid. Nonetheless, one can observe there is less clear semantic consistency even in Airplane category without degeneracies. As can be seen in the results, our model, aiming for a common representation for all the instances of the category, avoids placing keypoints in less representative parts or unique parts, e.g., arm rests in chairs (in Fig. 1), engines in airplanes or gas tank in motorbikes. This highlights a significant difficulty in modeling, requiring well-constrained and effective learning mechanism in order to achieve robustness.

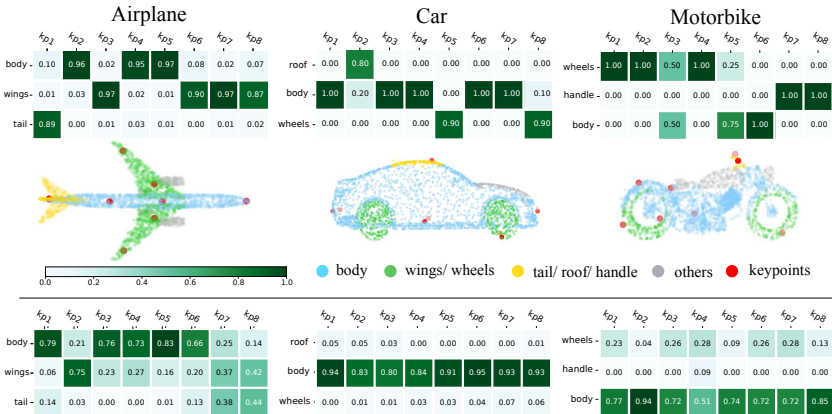


Fig. 4: **Semantic part correspondence.** First block, including qualitative results, presents the semantic correspondence for our category-specific keypoints, whereas the second block presents USIP results. Our predicted keypoints correspond to semantically meaningful locations across the category.

5.3 Objects Pose and Intra-category Registration

Previous methods do not handle misaligned data due to the obvious difficulty it poses to unsupervised learning. This deserves special attention since real data is never aligned. In this section we evaluate the intra-category registration performance of our model and show the impact of the different symmetry models proposed. These results implicitly measure the object poses estimated as well.

Rotation Ambiguities: Recent unsupervised approaches for keypoint detection actually self-supervise rotation during training, e.g., [23, 24], and highlight that it is crucial for achieving a good performance. In our case, we do not supervise the learned rotations and discover that different combination of basis shapes can result in different alignments. This means that computing P_i with the deformation function $\Phi_{\mathcal{C}}$ will give a correct set of keypoints, but the predicted rotations alone are not meaningful for a quantitative evaluation.

Experimental setup: Despite the above ambiguity, an important characteristic of the proposed keypoints is that they are ordered, which empowers direct inter-instances registration since no extra descriptors are needed for matching. We perform experiments for the chair category, using 10 keypoints (Table 1) and a misalignment of 45 degrees range on each axis. Three different models are compared. The first one is trained without symmetry awareness following Eq. (2). A second one uses shape symmetry during training as shown in Eq. (3). The last model is trained with basis symmetry as in Eq. (4). We attempt to register keypoints in each instance to those of randomly chosen three aligned templates by computing a similarity transformation and observe the mean. Fig. 5 shows that symmetry helps to have more control over the rotations and tackle higher misalignment. More results and analysis are provided in the supplementary.

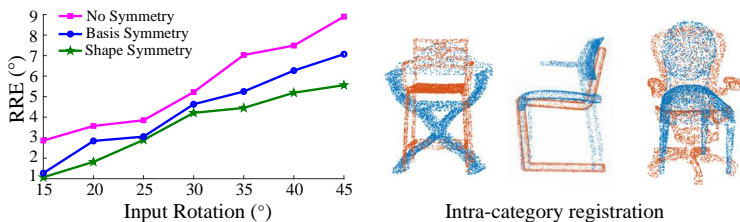


Fig. 5: Left: Relative rotation error for different symmetry modelings. Right: 3 examples of registration between different instances of the same category.

6 Conclusions

This paper investigates automatic discovery of keypoints in 3D misaligned point clouds that are consistent over inter-subject shape variations and intra-subject deformations in a category. We find that this can be solved, with unsupervised learning, by modeling keypoints with non-rigidity, based on symmetric linear basis shapes. Additionally, the proposed category-specific keypoints have one-to-one ordered correspondences and semantic consistency. Applications for the learned keypoints include registration, recognition, generation, shape completion

and many more. Our experiments showed that high quality keypoints can be obtained using the proposed methods and that the method can be extended to complex non-rigid deformations. Future work could focus on better modeling complex deformations with non-linear approaches.

References

1. Novotny, D., Ravi, N., Graham, B., Neverova, N., Vedaldi, A.: C3dpo: Canonical 3d pose networks for non-rigid structure from motion. In: Proceedings of the IEEE International Conference on Computer Vision. (2019) 7688–7697
2. Dai, Y., Li, H., He, M.: A simple prior-free method for non-rigid structure-from-motion factorization. In: CVPR. (2012)
3. Snavely, N., Seitz, S.M., Szeliski, R.: Modeling the world from internet photo collections. *Int. J. Comput. Vision* **80**(2) (2007) 189–210
4. Yew, Z.J., Lee, G.H.: 3dfeat-net: Weakly supervised local 3d features for point cloud registration. In: European Conference on Computer Vision, Springer (2018) 630–646
5. Kneip, L., Li, H., Seo, Y.: Upnp: An optimal $\mathcal{O}(n)$ solution to the absolute pose problem with universal applicability. In Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T., eds.: *Computer Vision – ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part I.* (2014)
6. Luong, Q.T., Faugeras, O.: The fundamental matrix: theory, algorithms, and stability analysis. *International Journal of Computer Vision* **17** (1995) 43–75
7. Loper, M., Mahmood, N., Romero, J., Pons-Moll, G., Black, M.J.: SMPL: A skinned multi-person linear model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)* **34**(6) (2015) 248:1–248:16
8. Shotton, J., Fitzgibbon, A., Cook, M., Sharp, T., Finocchio, M., Moore, R., Kipman, A., Blake, A.: Real-time human pose recognition in parts from single depth images. In: CVPR 2011, Ieee (2011) 1297–1304
9. Moreno-Noguer, F.: 3d human pose estimation from a single image via distance matrix regression. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2017) 2823–2832
10. Cao, Z., Simon, T., Wei, S.E., Sheikh, Y.: Realtime multi-person 2d pose estimation using part affinity fields. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2017) 7291–7299
11. Bogo, F., Kanazawa, A., Lassner, C., Gehler, P., Romero, J., Black, M.J.: Keep it smpl: Automatic estimation of 3d human pose and shape from a single image. In: European Conference on Computer Vision, Springer (2016) 561–578
12. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: Proceedings of the IEEE international conference on computer vision. (2017) 2961–2969
13. Sattler, T., Leibe, B., Kobbelt, L.: Fast image-based localization using direct 2d-to-3d matching. In: 2011 International Conference on Computer Vision, IEEE (2011) 667–674
14. Tang, H., Xu, D., Liu, G., Wang, W., Sebe, N., Yan, Y.: Cycle in cycle generative adversarial networks for keypoint-guided image generation. In: Proceedings of the 27th ACM International Conference on Multimedia. (2019) 2052–2060
15. Zafeiriou, S., Chrysos, G.G., Roussos, A., Ververas, E., Deng, J., Trigeorgis, G.: The 3d menpo facial landmark tracking challenge. In: Proceedings of the IEEE International Conference on Computer Vision Workshops. (2017) 2503–2511

16. Huang, S., Gong, M., Tao, D.: A coarse-fine network for keypoint localization. In: Proceedings of the IEEE International Conference on Computer Vision. (2017) 3028–3037
17. Pavlakos, G., Zhou, X., Chan, A., Derpanis, K.G., Daniilidis, K.: 6-dof object pose from semantic keypoints. In: ICRA. (2017)
18. Zhang, Z., Luo, P., Loy, C.C., Tang, X.: Facial landmark detection by deep multi-task learning. In: European conference on computer vision, Springer (2014) 94–108
19. Dong, X., Yan, Y., Ouyang, W., Yang, Y.: Style aggregated network for facial landmark detection. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2018) 379–388
20. Yu, X., Zhou, F., Chandraker, M.: Deep deformation network for object landmark localization. In: European Conference on Computer Vision, Springer (2016) 52–70
21. Alahi, A., Ortiz, R., Vandergheynst, P.: Freak: Fast retina keypoint. In: 2012 IEEE Conference on Computer Vision and Pattern Recognition, Ieee (2012) 510–517
22. Li, Y.: A novel fast retina keypoint extraction algorithm for multispectral images using geometric algebra. *IEEE Access* **7** (2019) 167895–167903
23. Li, J., Lee, G.H.: Usip: Unsupervised stable interest point detection from 3d point clouds. In: Proceedings of the IEEE International Conference on Computer Vision. (2019) 361–370
24. Suwajanakorn, S., Snavely, N., Tompson, J.J., Norouzi, M.: Discovery of latent 3d keypoints via end-to-end geometric reasoning. In: Advances in Neural Information Processing Systems. (2018) 2059–2070
25. Yang, H., Carlone, L.: In perfect shape: Certifiably optimal 3d shape reconstruction from 2d landmarks. *arXiv preprint arXiv:1911.11924* (2019)
26. Hejrati, M., Ramanan, D.: Analyzing 3d objects in cluttered images. In: Advances in Neural Information Processing Systems. (2012) 593–601
27. Wang, C., Wang, Y., Lin, Z., Yuille, A.L., Gao, W.: Robust estimation of 3d human poses from a single image. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2014) 2361–2368
28. Persad, R.A., Armenakis, C.: Automatic 3d surface co-registration using keypoint matching. *Photogrammetric engineering & remote sensing* **83**(2) (2017) 137–151
29. Mitra, N.J., Wand, M., Zhang, H., Cohen-Or, D., Kim, V., Huang, Q.X.: Structure-aware shape processing. In: ACM SIGGRAPH 2014 Courses. (2014) 1–21
30. Reed, M.P.: Modeling body shape from surface landmark configurations. In: International Conference on Digital Human Modeling and Applications in Health, Safety, Ergonomics and Risk Management, Springer (2013) 376–383
31. Sridhar, S., Rempe, D., Valentin, J., Sofien, B., Guibas, L.J.: Multiview aggregation for learning category-specific shape reconstruction. In: Advances in Neural Information Processing Systems. (2019) 2348–2359
32. Gao, Y., Yuille, A.L.: Symmetric non-rigid structure from motion for category-specific object structure estimation. In: European Conference on Computer Vision, Springer (2016) 408–424
33. Wu, J., Xue, T., Lim, J.J., Tian, Y., Tenenbaum, J.B., Torralba, A., Freeman, W.T.: Single image 3d interpreter network. In: European Conference on Computer Vision, Springer (2016) 365–382
34. Bregler, C., Hertzmann, A., Biermann, H.: Recovering non-rigid 3D shape from image streams. In: CVPR. (2000)
35. Torresani, L., Hertzmann, A., Bregler, C.: Nonrigid structure-from-motion: Estimating shape and motion with hierarchical priors. *IEEE Trans. Pattern Anal. Mach. Intell.* **30**(5) (2008) 878–892

36. Akhter, I., Sheikh, Y., Khan, S., Kanade, T.: Nonrigid structure from motion in trajectory space. In: NIPS. (2008)
37. Taylor, J., Jepson, A.D., Kutulakos, K.N.: Non-rigid structure from locally-rigid motion. In: CVPR. (2010)
38. Parashar, S., Pizarro, D., Bartoli, A.: Isometric non-rigid shape-from-motion in linear time. In: CVPR. (2016)
39. Kong, C., Lucey, S.: Deep non-rigid structure from motion. In: Proceedings of the IEEE International Conference on Computer Vision. (2019) 1558–1567
40. Qi, C.R., Su, H., Mo, K., Guibas, L.J.: Pointnet: Deep learning on point sets for 3d classification and segmentation. In: CVPR. (2017)
41. Qi, C.R., Yi, L., Su, H., Guibas, L.J.: Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In: Advances in neural information processing systems. (2017)
42. Verma, N., Boyer, E., Verbeek, J.: Feastnet: Feature-steered graph convolutions for 3d shape analysis. In: CVPR. (2018)
43. Bromley, J., Guyon, I., LeCun, Y., Säckinger, E., Shah, R.: Signature verification using a” siamese” time delay neural network. In: Advances in neural information processing systems. (1994) 737–744
44. Tomasi, C., Kanade, T.: Shape and motion from image streams under orthography: a factorization method. *International journal of computer vision* **9**(2) (1992) 137–154
45. Yi, L., Kim, V.G., Ceylan, D., Shen, I.C., Yan, M., Su, H., Lu, C., Huang, Q., Sheffer, A., Guibas, L.: A scalable active framework for region annotation in 3d shape collections. *ACM Transactions on Graphics (TOG)* **35**(6) (2016) 1–12
46. Wu, Z., Song, S., Khosla, A., Yu, F., Zhang, L., Tang, X., Xiao, J.: 3d shapenets: A deep representation for volumetric shapes. In: CVPR. (2015) 1912–1920
47. Li, J., Chen, B.M., Hee Lee, G.: So-net: Self-organizing network for point cloud analysis. In: Proceedings of the IEEE conference on computer vision and pattern recognition. (2018) 9397–9406
48. Fan, H., Su, H., Guibas, L.J.: A point set generation network for 3d object reconstruction from a single image. In: Proceedings of the IEEE conference on computer vision and pattern recognition. (2017) 605–613
49. Besl, P.J., McKay, N.D.: Method for registration of 3-d shapes. In: Sensor fusion IV: control paradigms and data structures. Volume 1611., International Society for Optics and Photonics (1992) 586–606
50. Bogo, F., Romero, J., Pons-Moll, G., Black, M.J.: Dynamic faust: Registering human bodies in motion. In: CVPR. (2017) 6233–6242
51. Gerig, T., Morel-Forster, A., Blumer, C., Egger, B., Luthi, M., Schönborn, S., Vetter, T.: Morphable face models-an open framework. In: 2018 13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018), IEEE (2018) 75–82
52. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: Pytorch: An imperative style, high-performance deep learning library. In: NIPS. (2019)
53. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014)
54. Song, S., Lichtenberg, S.P., Xiao, J.: Sun rgb-d: A rgb-d scene understanding benchmark suite. In: Proceedings of the IEEE conference on computer vision and pattern recognition. (2015) 567–576

Supplementary Material

Abstract. In this **supplementary document**, we provide more details on our network architecture. We also give more insights regarding symmetry, including the proof of Proposition 1. Furthermore, experiments showing the generalization of our method on real data is included, as well as some results for the segmentation label transfer task. Finally additional qualitative results are presented on the four datasets evaluated in the main paper at the end of the document.

S1 Network Architecture Details

In this section, we detail the network architecture used throughout our experimental evaluations. Fig. S1 illustrates our network architecture. Our network comprises of two main branches that predict the instance-specific parameters; one branch that estimates a sparse tuple of unordered nodes that are potentially category-specific keypoints, X_i , and a second branch that predicts the parameters of the non-rigid shape (rotation, R_i , and basis coefficients, c_i). During training, we additionally learn the basis shapes, \mathcal{B}_C , and the normal direction of the plane of symmetry, n_C , that are common for the category shape space \mathcal{C} . Such quantities are optimized as network parameters.

The learning strategy consists of learning the parameters of the deformation function, Φ_C , which is able to explain the estimated nodes obtained from the neural network Π_C , as a linear combination of basis shapes. The final sparse tuple of ordered points that we name category-specific keypoints, P_i , are obtained from the deformation function Φ_C .

In our experiments, we consider a single point cloud as an input and the sequence of layers comprising the two main branches are depicted in Fig. S1.

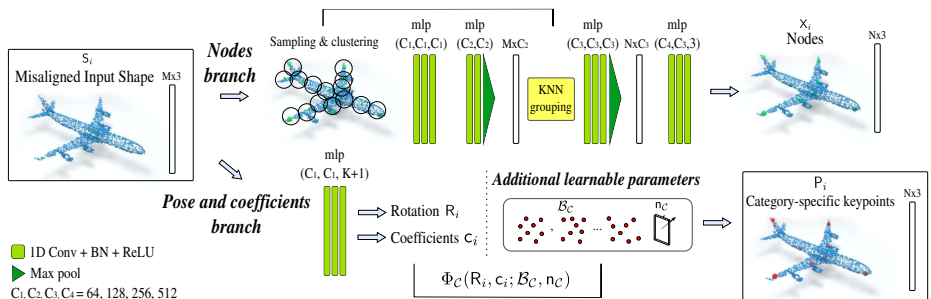


Fig. S1: **Illustration of our network architecture.** The *pose and coefficients branch* and the *additional learnable parameters* generate the output category-specific keypoints. The *nodes branch* estimates the nodes that guide the learning process. “mlp” stands for multi-layer perceptron.

S2 Symmetry

S2.1 Symmetric deformation space.

Proof. The two linear spaces due to the two basis $\mathcal{B}_{\mathcal{C}_{\frac{1}{2}}}$ and $\mathcal{B}'_{\mathcal{C}'_{\frac{1}{2}}}$ are symmetric by Definition 1 as $\mathcal{B}_{\mathcal{C}_{\frac{1}{2}}}$ is symmetric to $\mathcal{B}'_{\mathcal{C}'_{\frac{1}{2}}}$ for any $K \in \mathbb{Z}$. Let $\mathbf{c}_i \in \mathcal{L}$ and $\mathbf{c}'_i \in \mathcal{L}'$, such that \mathcal{L} and \mathcal{L}' define the spaces of the predicted coefficient vectors. Consequently, the actual deformation spaces are symmetric to one another if \mathcal{L} and \mathcal{L}' are equal. We define $p : p(\mathbf{c}_i)$ as the probability distribution of \mathbf{c}_i and $q : q(\mathbf{c}'_i)$ as the probability distribution of \mathbf{c}'_i . If p and q come from the same distribution, we approach $p = q$. Then we have:

$$\begin{aligned}
 & \text{if } \mathbf{c}_i = \mathbf{c}'_i, \\
 & \text{either, } p(\mathbf{c}_i) = q(\mathbf{c}'_i) = 0, \\
 & \text{or, } p(\mathbf{c}_i) > 0 \text{ and } q(\mathbf{c}'_i) > 0 \\
 & \text{for all, } \mathbf{c}_i \in \mathcal{L}, \mathbf{c}'_i \in \mathcal{L}'.
 \end{aligned} \tag{S1}$$

Condition (S1) guarantees that $\mathcal{L} = \mathcal{L}'$ and thus we obtain a symmetric deformation space. \square

Note that for condition (S1) to be true, we do not require the two distributions to be equal, however, it is sufficient and desirable to have so. Therefore, Proposition 1 in the main text highlights such sufficient and desirable case. It is particularly meaningful when we are learning to predict the coefficients through stochastic methods such as a neural network training. In the network architecture of Fig. S1, indeed one can expect the distributions of these two vectors to be similar given the data exhibits such a symmetric deformation space, since the prediction branches of \mathbf{c}_i and \mathbf{c}'_i are very similar. Alternatively, one may also try to enforce the condition using a KL divergence loss.

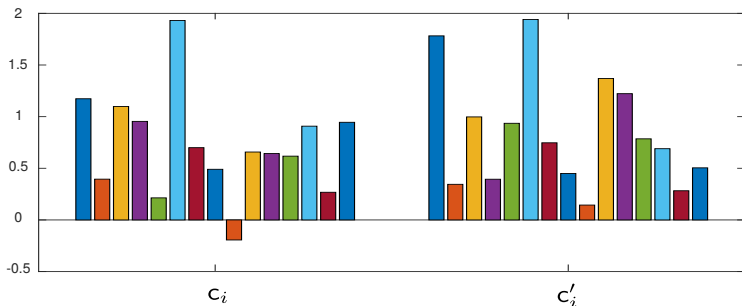


Fig. S2: **Coefficients distribution.** Mean values of \mathbf{c}_i components (left) and \mathbf{c}'_i components (right) for the Dynamic FAUST [50]. The mean of the variances for the different components are: $\mathbf{c}_i : 0.54$, $\mathbf{c}'_i : 0.50$. The figure shows that the network learns similar distribution for the coefficients \mathbf{c}_i and \mathbf{c}'_i .

S2.2 Symmetry Plane Parametrization.

As mentioned in Sec. 5.3 in the main paper, we observe that handling misaligned data with unsupervised methods can lead to some rotation ambiguities. More specifically, we observe that different combination of basis shapes can result in different alignments.

As we show in Fig. 5 in the text, predicting the symmetry plane of the object category allows to have more control over the predicted instance poses. We came up with the idea of learning an additional common parameter, R_C , which is directly related to the symmetry plane. By adding this category-specific parameter, the network learns a common rotation for all the objects in the category. As a consequence, the instance-wise rotation, R_i , can be thought like an offset from the reference basis alignment. Several evaluations confirmed that this strategy helps the learning process, reducing the rotation ambiguities.

S3 Additional Experiments

S3.1 Keypoints correspondence

We provide a complete overview for all the object categories evaluated regarding the keypoints correspondences across instances in Fig. S3. This demonstrates the ability of our model to capture and model the inter-subject shape variations and intra-subject deformations in a category.

S3.2 Segmentation Label Transfer

As demonstrated in Sec. 5.2 in the main paper, our predicted keypoints correspond to semantically meaningful locations. Therefore, here we explore the utility of the proposed category-specific keypoints for the segmentation label transfer task. In this experiment, for every point in the original shape $\mathbf{s}_{ij} \in \mathcal{S}_i$, we find its closest category-specific keypoint $\mathbf{p}_{ik} \in \mathcal{P}_i$, and transfer the corresponding semantic label to it. We assume the keypoints labels are known and correspond to those in Fig. 4 in the paper.

Some qualitative results are shown in Fig. S4. Our method achieves full correspondence between instances, therefore avoiding placing keypoints in less representative parts. An example is the engine, in grey, in the case of airplanes. This is reflected in the label transfer since there is no distinction of these parts. Besides that, only with eight keypoints in the example, we achieve reasonable results, close to the ground truth data.

S3.3 Real Data

In this section, we show the performance of our method for real data in Fig. S5. For this experiment, the network is trained on the chair category from the ModelNet10 dataset [46] and tested on real chairs from the SUNRGBD dataset [54]. To generate the real data dataset from [54], we crop the points inside the ground

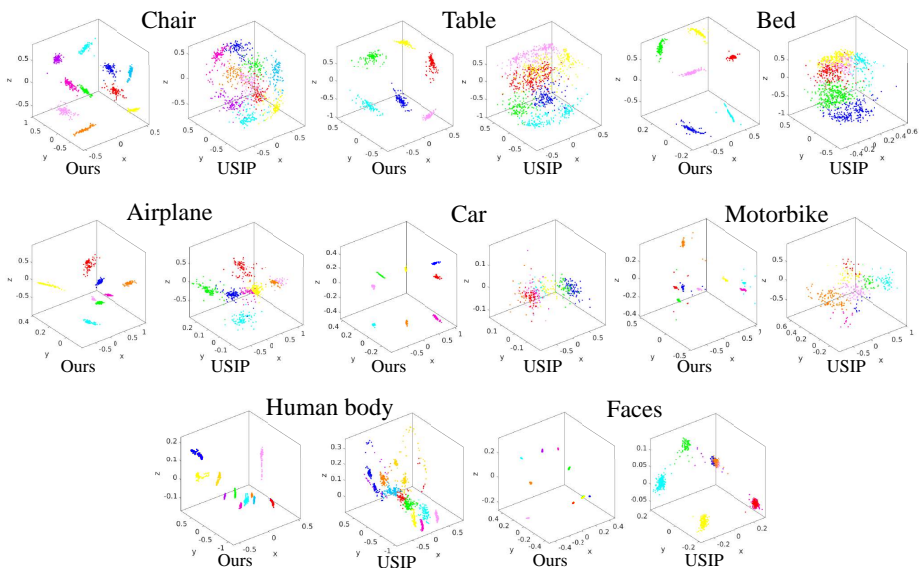


Fig. S3: **Keypoints correspondence across instances.** We cluster the keypoints predicted for all the instances of a category to show their geometric consistency. Note how our keypoints get neatly clustered creating a general 3D shape template.

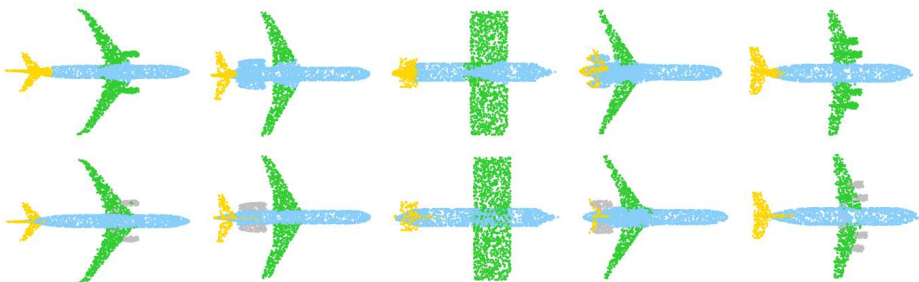


Fig. S4: **First row:** results of performing semantic label transfer with our keypoints. **Second row:** ground truth. This is evaluated in ShapeNet part dataset [45] using eight keypoints for the label transfer.

truth 3D bounding boxes provided by the authors. Real data entail additional challenges. This is not only because shapes appear incomplete and noisy, but also because other objects may cause occlusions, e.g. part of a table occluding a chair. As illustrated in Fig. S5, even though real data is fairly challenging, our network can still produce corresponding meaningful keypoints.

Being able to generalize to previously unseen real objects as demonstrated in Fig. S5 is crucial and really useful for many tasks such as guide for shape completion or shape generation.

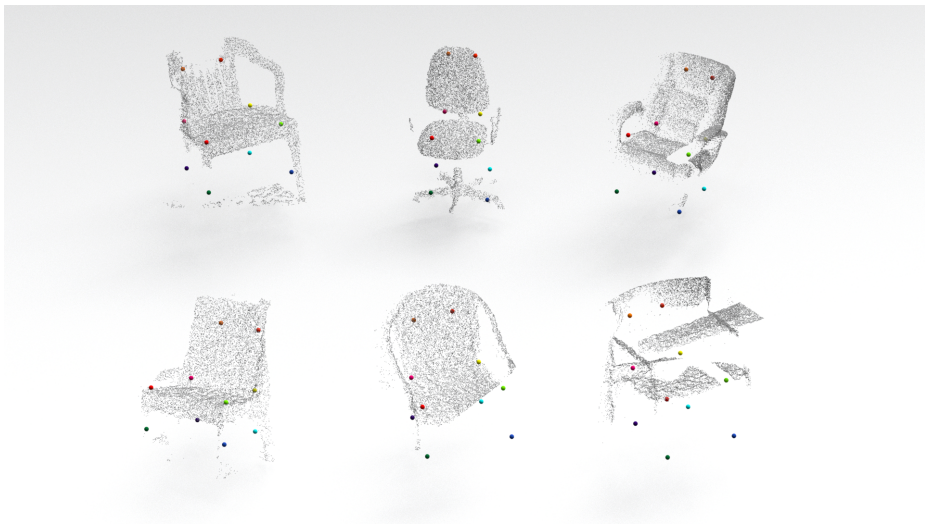


Fig. S5: Results in real chairs from SUNRGBD dataset [54] training with CAD chairs from ModelNet10 dataset [46].

S4 Qualitative results

In this section, we provide additional qualitative results on various object categories from the datasets evaluated in the paper; ModelNet10 [46] in Fig. S6, ShapeNet parts [45] in Fig. S7, Dynamic FAUST [50] in Fig. S8 and Basel Face Model 2017 [51] in Fig. S9.

Again, we note that our network predicts corresponding keypoints between instances of the same category and consistently associates the same keypoint with the same semantic part. For instance, for the chair object category, the keypoint colored in pink is always associated with the chair back, the keypoint colored in cyan is associated with the front left leg, etc.



Fig. S6: Qualitative results in table, chair and bed categories from ModelNet10 dataset [46].

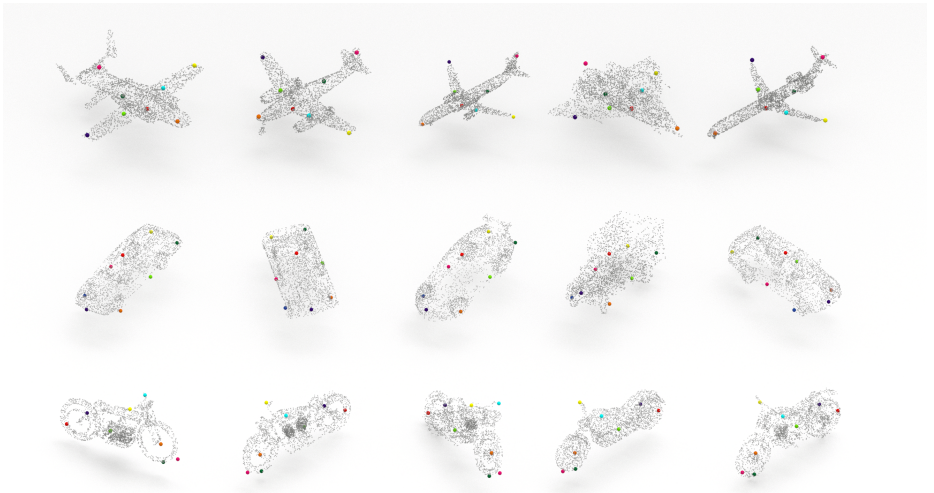


Fig. S7: Qualitative results in airplane, car and motorbike categories from ShapeNet parts dataset [45].

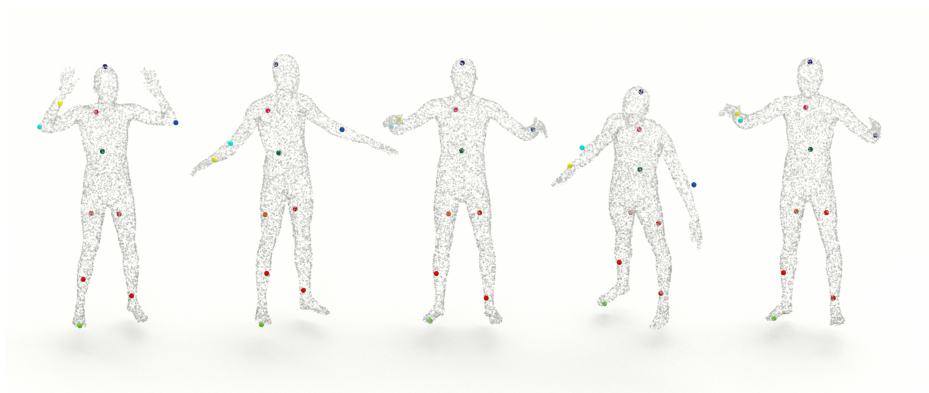


Fig. S8: Qualitative results in human bodies from Dynamic FAUST dataset [50].

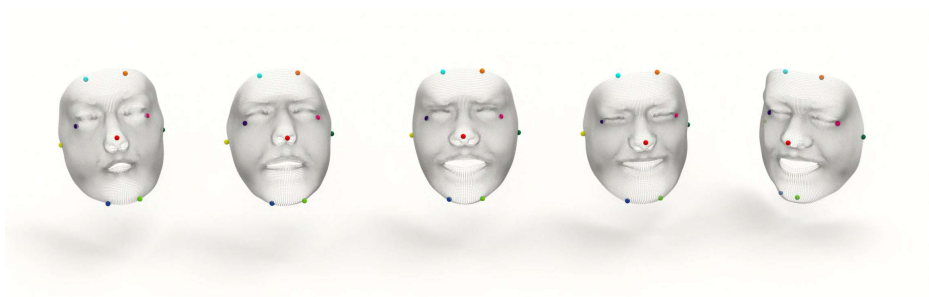


Fig. S9: Qualitative results in faces from Basel Face Model 2017 dataset [51].