# Self-Supervised Learning of Part Mobility from Point Cloud Sequence

Yahao Shi, Xinyu Cao and Bin Zhou[†]

State Key Laboratory of Virtual Reality Technology and Systems, School of Computer Science and Engineering, Beihang University, Beijing, China

**Abstract**

*Part mobility analysis is a significant aspect required to achieve a functional understanding of 3D objects. It would be natural to obtain part mobility from the continuous part motion of 3D objects. In this study, we introduce a self-supervised method for segmenting motion parts and predicting their motion attributes from a point cloud sequence representing a dynamic object. To sufficiently utilize spatiotemporal information from the point cloud sequence, we generate trajectories by using correlations among successive frames of the sequence instead of directly processing the point clouds. We propose a novel neural network architecture called PointRNN to learn feature representations of trajectories along with their part rigid motions. We evaluate our method on various tasks including motion part segmentation, motion axis prediction, and motion range estimation. The results demonstrate that our method outperforms previous techniques on both synthetic and real datasets. Moreover, our method has the ability to generalize to new and unseen objects. It is important to emphasize that it is not required to know any prior shape structure, prior shape category information, or shape orientation. To the best of our knowledge, this is the first study on deep learning to extract part mobility from point cloud sequence of a dynamic object.*

**CCS Concepts**
• *Computing methodologies* → *Shape analysis;* *Animation; Neural networks;*
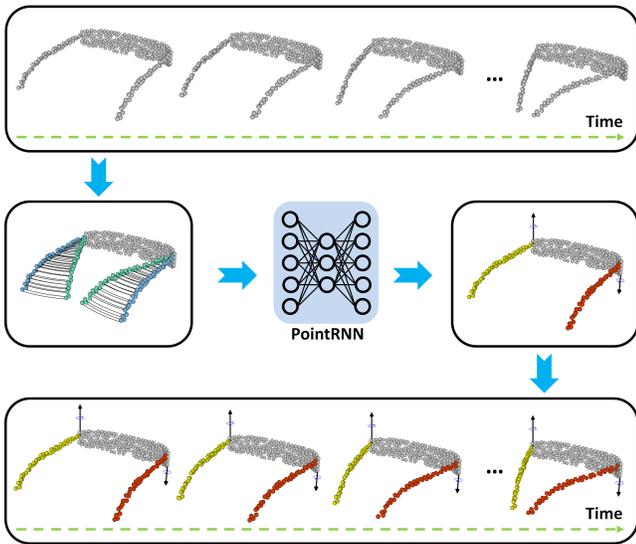
## 1. Introduction

In the real world, there often exist a number of dynamic and articulated objects, which can be directly operated using their moving parts. For example, we may need to open a refrigerator door to keep or remove an object. If we expect autonomous agents to interact with such objects correctly, they must have the ability to analyze which part of an object is a moving part. Hence, learning part mobility of 3D objects is beneficial for 3D computer vision [LWL[*]16] and robotics [HLRB13] and is closely related to the understanding of object affordances [MTFA15], functionality [HSvK18], and interaction using object recognition [LSZ[*]19] and human motion capture data [RLA[*]19].

Recently, with the emergence of large 3D labeled datasets and deep learning techniques, several studies have been conducted to make considerable progress in supervised semantic part segmentation, such as PointNet++ [QYSG17] and PointCNN [LBS[*]18]. However, these segments were often defined by personal subjective experience, which can easily lead to ambiguity. Thus, research on learning moveable part segmentation of 3D objects is of more practical significance, and it helps agents understand the essential features of dynamic objects. Moreover, current research on pars-

ing 3D shape representations into semantic parts faces challenges in processing novel object categories and discovering new functional parts. It is not conducive for agents to explore strange new worlds. Ideally, intelligent agents should be able to parse 3D shapes into previously unseen functional parts from observations of continuous part motion. Mobility-based shape parsing brings a novel perspective to the part determination problem and provides an unambiguous decomposition. Additionally, a mobility-based part segmentation algorithm enables intelligent agents to make better use of man-made objects designed to function or interact with other objects (including humans).

In this study, we are interested in discovering the part mobility of 3D objects by observing their continuous part motion. We define part mobility as motion part segmentation, motion axis prediction, and motion range estimation. In previous studies, motion parts have been typically extracted from one or two single static snapshots of dynamic objects. By contrast, we aim to deduce motion part structure from observations of continuous articulation states of an object. The reason is that speculating the motion parts and confirming motion vectors by observing only a few motion states can easily cause confusion. Moreover, point cloud sequences reveal sufficient spatiotemporal information, and they are easier to obtain than large well-annotated datasets, owing to recent advances in the techniques of real-time 3D acquisition such as commercial RGB-D

---

[†] Corresponding author: zhoubin@buaa.edu.cn

**Figure 1:** *In our method, we first convert a raw point cloud sequence into trajectories. Then, we adopt PointRNN to extract varieties of part rigid motion hypotheses by taking these trajectories as input. Finally, we merge these hypotheses and segment motion parts using an iterative algorithm. The results can propagate from the first frame to other frames.*

cameras. Therefore, we adopt point cloud sequences rather than a single point cloud to represent dynamic objects.

Automatic motion part induction from point cloud sequences is challenging for several reasons. First, objects differ significantly because of their geometry and pose. Motion directions and motion ranges of motion parts have noticeable differences. Second, a point cloud is unordered. There is no tight point correspondence between adjacent frames of point cloud sequences. Third, we must consider the influence factors of the acquisition quality of scan data, including noisy and missing data.

To address the above mentioned problems, we propose a novel self-supervised deep neural network-based method (See Fig. 1), inspired by the traditional methods of Yan and Pollefeys [YP06]. We observed that points on the same part produce similar motion trajectories. Therefore, we treat motion part segmentation as a trajectory clustering problem, and we can calculate rigid motion hypotheses based on these trajectory sets. However, directly handling point cloud sequences can be difficult in a network without point correspondence. Thus, we first transform a point cloud sequence into a bunch of trajectories. Then, we design a neural network called PointRNN to process these trajectories, and it can learn latent trajectory feature representations and generate candidate rigid motion hypotheses represented by motion axes and motion ranges. Finally, we merge the trajectories belonging to a similar rigid motion to achieve motion part segmentation.

Our method is verified through qualitative and quantitative evaluations on both synthetic and real datasets. Additionally, we con-

duct ablation experiments to confirm the influence of different loss designs, hyperparameter settings, etc. The experimental results demonstrate that our method achieves better performance and requires less time than the traditional method [YLX*16] and deep learning methods [WZS*19, YHL*18, YHY*19]. Moreover, results also demonstrate that our approach can tolerate some noises and has the ability to apply to unseen categories.

In summary, a new self-supervised approach is proposed to parse 3D shapes into moving parts, motion axes, and motion ranges from point cloud sequences without labeling or any prior knowledge. Specifically, our method makes three key contributions:

- We introduce a self-supervised approach to learning motion part segmentation, motion attribute estimation, and motion flow prediction.
- We propose a novel neural network called PointRNN, which has the capability to process trajectories and extract a feature representation.
- Experimental results demonstrate that the performance of our network is superior to state-of-the-art methods, and it can apply to unseen object categories.

## 2. Related Work

Deformable shape registration [SA07] is a fundamental problem in computational geometry and widely applied to various fields such as computer vision. There are various deformable registration algorithms for point clouds. Many methods are extensions of the classic ICP algorithm [BM92]. Papazov and Burschka [PB11] proposed an algorithm that computes shape transitions based on local similarity transforms, thereby allowing it to model not only as-rigid-as-possible deformations but also shapes on local and global scales.

Many 3D shape segmentation approaches have been proposed in previous studies to extract moving parts from an RGB-D sequence, a single point cloud, a point cloud sequence, and a mesh model. Given a specified RGB-D sequence that contains a dynamic object, attempts have been made in previous studies to recover the 3D scene flow, thus discovering moving parts of 3D objects. Jaimez et al. [JSGC15] proposed a primal–dual algorithm to compute RGB-D flow for estimating heterogeneous and non-rigid motion at a high frame rate. Vogel et al. [VRS14] introduced a method to recover dense 3D scene flow from multiple consecutive frames in a sliding temporal window. Motion part segmentation can be extracted by point correspondence, as shown by Fayad et al. [FRA11], and then they can utilize these motion parts to reconstruct the articulated structure. However, common defects in these approaches are their reliance on RGB color to compute scene flow and inability to handle complex structures or large motions.

In the case of a raw 3D point cloud sequence, many studies aimed to establish a point-wise correspondence between consecutive point clouds of an articulated shape [CZ08], [PB11]. Yan and Pollefeys [YP06] cast the problem of motion segmentation of feature trajectories as linear manifold finding problems and proposed a general framework for motion segmentation under affine projections. Besides, Kim et al. [KLAK16] considered that trajectories could be grouped by clustering to separate different motion parts.

Yuan et al. [YLX*16] proposed a local-to-global approach to co-segment point cloud sequences of articulated objects into near-rigid moving parts. Most of the methods mentioned above require a large amount of computation to achieve better performance. These approaches also require considerable effort in threshold setting from case to case.

Several approaches parse mesh models into moving parts. Mitra et al. [MYY*13] inferred the motion of individual parts and the interactions among parts based on their geometry and a few user-specified constraints. They utilized the results to illustrate the motion of mechanical assemblies. Hu et al. [HLVK*17] introduced a data-driven approach for learning part mobility based on a defined motion pattern from a single static state of 3D objects. However, these methods are based on a well-defined motion pattern and well-segmented 3D objects.

Recent supervised learning approaches in 3D shape segmentation employ deep network architecture to train a classifier on labeled data represented by a point cloud [QYSG17], volumetric grids [MS15] or spatial data structures [KL17]. Subsequent attempts have been made to extract part rigid motions using new deep learning technology. Wang et al. [WZS*19] proposed the shape2motion network architecture, which takes a single point cloud as input. Their network aims to simultaneously segment motion parts and predict motion axes based on a large, well-annotated dataset. However, their approach has difficulty in discovering new motion parts and generalizing them to novel categories. Their network requires a considerable amount of time for supervised learning using a large dataset. Another type of method, such as that proposed by Yi et al. [YHL*18], infers part motion flow and estimates part correspondence by comparing two different motion states. Unfortunately, they did not perform a specific analysis of motion attributes such as motion axes and motion ranges. Behl et al. [BPDG19] and Liu et al. [LQG19] promoted this class of algorithm to identify motion flow in scene data. However, these methods still rely on supervised learning. Yan et al. [YHY*19] introduced RPM-Net to infer movable parts of a single point cloud. They adopted a recurrent neural network (RNN) to movable segment parts by forecasting a temporal sequence motion of 3D objects. In contrast, our network takes point cloud sequences as input instead of a single point cloud. Furthermore, we employ an RNN to process each trajectory locally rather than the entire point cloud directly and globally.

## 3. Method

Our method comprises three steps for learning part mobility, as shown in Fig. 1. First, we adopt a deformable registration algorithm using point correspondence to generate trajectories (Sec. 3.1). Second, we propose a network architecture called PointRNN to produce candidate motion hypotheses extracted from trajectories (Sec. 3.2). Third, we design an iterative algorithm for removing redundant motion hypotheses and segmenting a point cloud into several motion parts while considering the matching degree between trajectories and motion hypotheses (Sec. 3.3). The second step of our method occurs within the learning pipeline, and the other two steps occur outside it. In this section, we discuss the modules of our approach in detail.

### 3.1. Trajectory Generation

Two adjacent frames of the raw point cloud sequence do not have point-wise correlations because of unordered points. Therefore, we can not directly concatenate all the raw point cloud sequence frames and take them as input to the neural network. To employ the deep learning approach, we first convert the raw point cloud sequence into trajectories. Moreover, we adopt the method of Papazov and Burschka [PB11] based on local similarity transforms to generate motion trajectories to ensure a fair comparison with Yuan et al. [YLX*16], and this method is an alternative module to generate trajectories for our pipeline. All the trajectories are cropped to an equal length to facilitate deep neural network processing.

Given two adjacent frames, $PC_s$ and $PC_t$, of a point cloud sequence, the registration algorithm includes iterative steps of correspondence and deformation. In terms of the correspondence step, we find the closest point in $p' \in PC_t$ for each point $p \in PC_s$. We optimize $p'$ by searching over its $k$-nearest neighbors to make the match error smooth. Formally, we compute a mapping, $\pi_{s \to t}$, by minimizing the Laplacian smoothness energy of the residual field $\{\pi_{s \to t}(p) - p'\}$:

$$\mathcal{E}_{smooth}(\pi_{s \to t}) = \sum_{p \in PC_s} \|\delta(p)\|^2, \quad (1)$$

where the residual Laplacian, $\delta(p)$, is defined as

$$\delta(p) = \frac{1}{|\mathcal{N}(p)|} \sum_{q \in \mathcal{N}(p)} [(\pi_{s \to t}(p) - p') - (\pi_{s \to t}(q) - q')]. \quad (2)$$

Here, $\mathcal{N}(p)$ is the $k$-nearest neighbors of $p$. In terms of the deformation step, we minimize a fitting error between the local neighborhood of $p \in PC_s$ and its matched counterpart to estimate a similarity transformation, denoted as $(\mathbf{R}, \mathbf{t})$, by using the singular value decomposition [SMW06].
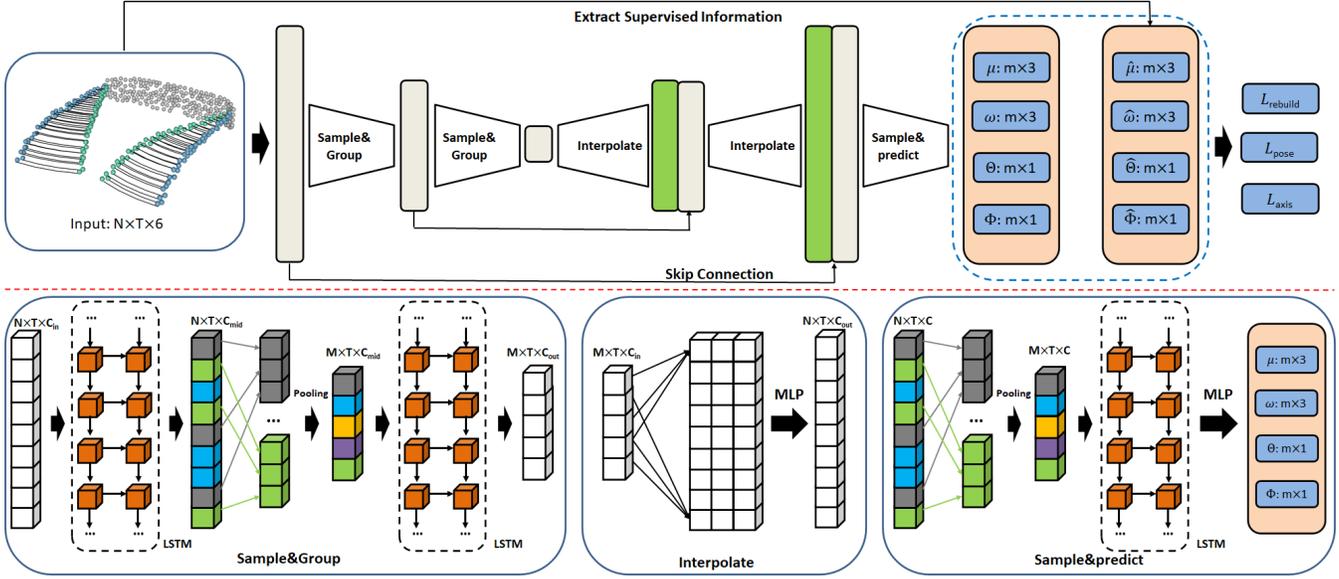
$$\mathcal{E}_{fitting}(\mathbf{R}, \mathbf{t}) = \sum_{q \in \mathcal{N}(p)} \|\pi_{s \to t}(q) - (\mathbf{R}q + \mathbf{t})\|^2. \quad (3)$$

### 3.2. Part Rigid Motion Hypothsis Extraction

We propose a novel neural network called PointRNN, which has the capacity to process trajectories and extract the feature representation, as illustrated in Fig. 2 (Sec. 3.2.1). Moreover, we design a self-supervised loss function to train the PointRNN (Sec. 3.2.2).

### 3.2.1. PointRNN Network Architecture

Given a set of trajectories, $\{\tau_1, \tau_2, ..., \tau_n\}$ ($n = 2048$ by default), $\tau_i$ has a point set, $\{x_{i,1}, x_{i,2}, ..., x_{i,t}\}$ ($t = 11$ by default), with $x_{i,j} \in \mathbb{R}^3$ and a motion direction set, $\{d_{i,1}, d_{i,2}, ..., d_{i,t}\}$ ($t = 11$ by default), with $d_{i,j} \in \mathbb{R}^3$. There is a one-to-one relationship between these two sets, such as $(x_{i,1}, d_{i,1})$. For each trajectory, we employ a Long Short-Term Memory (LSTM) network [HS97] to encode time information, because LSTM has demonstrated high accuracy in tasks

**Figure 2:** *PointRNN network architecture consists of three parts. The Sample&Group Module is employed to learn temporal and spatial features from trajectories. We adopt the Interpolate Module to recover original resolution features from global features. We feed these features to the Sample&Predict Module to estimate motion parameters. Our supervised information is automatically obtained from the input. $c_{in}$, $c_{mid}$ and $c_{out}$ are the number of channels, which are different in different layers. Specifically, $c_{in} = 6, c_{mid} = 128, c_{out} = 64$ in the first Sample&Group Module. $c_{in} = 64, c_{mid} = 128, c_{out} = 256$ in the second Sample&Group Module. These are hyperparameters in our method.*

related to the processing of sequential and time-series data. In previous work, Yan et al. [YHY*19] used the PointNet++ encoder to create a global feature for the input point cloud. Then, they fed this feature vector to the LSTM network. Compared to their approach, ours pays more attention to local and detailed features for each trajectory.

We assume that similar trajectories should belong to one motion part, as inspired by the traditional clustering method [YP06]. Thus, we design a Sample&Group Module to learn the spatial and temporal features with their underlying motion from these similar trajectories. To this end, our network need to partition the set of trajectories into overlapping local similar trajectory partitions by the distance metric and then learns their features. Therefore, our network first performs a downsampling process on the set of trajectories. We adopt iterative farthest point sampling to obtain a subset of trajectories, $\{\tau_1, \tau_2, ..., \tau_i\}$, such that $\tau_j$ is the most distant trajectory (in metric distance) from set $\{\tau_1, \tau_2, ..., \tau_{j-1}\}$ ($j \leq i$). The advantage of the farthest point sampled algorithm is that it can cover the trajectory set in the space. Moreover, the sampling trajectories are isometry-invariant. The farthest point sampling algorithm is widely used in deep learning to process point cloud, such as PointNet++ [QYSG17]. The metric distance function between two trajectories, $(\tau_i, \tau_j)$, is as follows:

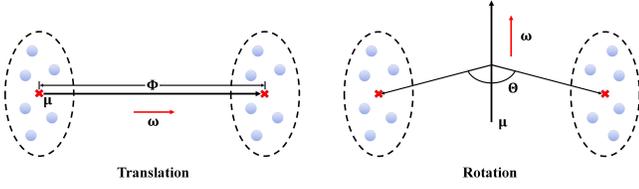$$dis_{(\tau_i, \tau_j)} = \frac{1}{t} \sum_{1 \leq m \leq t} \|x_{i,m} - x_{j,m}\|^2. \tag{4}$$

We adopt Euclidean distance rather than Hausdorff distance or

Fréchet distance because the latter two metrics require significantly more computations without yielding improved performance. Next, we search a neighborhood trajectory set, $\{\tau_{i_1}, \tau_{i_2}, ..., \tau_{i_k}\}$, for each sampled trajectory, $\tau_i$, using the $k$-nearest neighbor method, where the metric distance between $\tau_i$ and $\tau_{i_m}$ ($1 \leq m \leq k$) does not exceed a certain threshold $\varepsilon$. Then, we obtain several overlapping partitions of trajectories.

Next, our network learns motion features from these partitions. Specifically, we employ an LSTM network to encode time information for each point trajectory. This feature contains the positions of each point at different times and the movement trend of each point. We utilize two LSTM networks before and after sampling. To encode spatial information, we aggregate the features of the trajectory partition by a max-pooling operation. This feature contains all trajectory features from each trajectory partition. We treat this feature as regional spatial features and express them as features of $\tau_i$. In our method, $k$ and $\varepsilon$ are hyperparameters in different layers.

We can obtain a global feature after Sample&Group Module, and then we need to propagate global features to single trajectory features. To this end, we design an Interpolate Module with distance-based interpolation and skip-connection. We achieve feature propagation by interpolating features from the $k$-nearest trajectory set, $\{\tau_{j_1}, \tau_{j_2}, ..., \tau_{j_k}\}$ ($k = 3$ by default), to each trajectory, $\tau_j$. We adopt the inverse distance-weighted average based on the $k$-nearest neighbors in interpolation. Finally, we feed these interpolated trajectory features into a Multilayer Perceptron (MLP) Network to encode single trajectory features.

We can obtain the final trajectory features after Interpolate Mod-

**Figure 3:** *Motion parameters (μ, ω, Θ, and Φ). μ is the start of the motion axis. The motion direction (ω) is a unit vector. Φ is the translation distance, and Θ is the rotation angle.*

ule. Then, we need to utilize these features to estimate the motion attribute for each motion part. Our network will produce abundant motion hypotheses because the number of motion part is uncertain for each object. Therefore, we design a Sample&Predict Module to generate part rigid motion hypotheses. By observing most objects in the real world, we found that the number of motion parts is finite and small. We assume that the number of motion parts is no more than ten based on our dataset. Similar to the Sample&Group Module, we obtain $m$ ($m = 64$ by default) trajectory partitions using farthest point sampling and $k$-nearest neighbor method ($k = 32$ by default). We consider that these sets can cover all the motion parts of 3D shapes. We utilize a max-pooling operation to obtain aggregated spatial features for each trajectory set. Then, we obtain temporal features using an LSTM network. Finally, we feed their temporal and spatial features to an MLP Network. The network outputs a motion axis including a start point ($\mu$) and direction ($\omega$), motion ranges including a shifted distance ($\Phi$), and a rotation angle ($\Theta$) (Fig. 3).

### 3.2.2. Self-Supervised Loss Function Design.

We assume that trajectories contain rich motion information. Therefore, we can utilize estimated motion parameters to confirm whether it can reconstruct the motion sequence. Moreover, we can calculate the approximate parameters from motion sequence and take them as supervised information.

To reconstruct the motion sequences, we compute the motion matrix, $M \in \mathbb{R}^{4 \times 4}$, by taking advantage of network output. The motion matrix is composed of a rotation matrix and a translation matrix ($M = R + t$). $t$ is easy to generate by using $\omega$ and $\Phi$. The problem of solving $R$ can be reduced to transforming the point set rotating around the given axis (including $\mu$ and $\omega$) into the point set rotating around the z-axis of the standard coordinate system, $O_{xyz}$. The basic idea of deriving this matrix is to divide the problem into a few known simple steps:

- We align the start of the given axis with the coordinate origin and move the point set.
- We rotate the given axis, and the point set such that the axis lies in the $yOz$ coordinate planes.
- We rotate the given axis, and the point set such that the axis is aligned with the z-axis.
- We utilize one of the fundamental rotation matrices to rotate the point sets using $\Theta$. Finally, we undo steps three to one.

We can rebuild the last frame ($PC_t$) from the first frame ($PC_1$) by pre-multiplication ($M$ on the left) in our neural network. Moreover, we implement the partial derivative for each network output in our neural network to apply the backpropagation algorithm. Therefore, we define the first loss function, called the rebuild loss function, as follows:

$$\mathcal{L}_{rebuild} = \|PC_t - PC_1 \cdot M\|^2. \tag{5}$$

For each trajectory set in the Sample&Predict Module, we solve an absolute orientation problem to obtain an approximate solution of rigid motion ($\hat{M}$). Our solution for the absolute orientation problem is based on the method of Myronenko and Song [MS09]. Besides, we guarantee that the determinant of $\hat{M}$ is 1 to ensure that $\hat{M}$ has a rotation matrix other than the reflection matrix. Then, we extract rotation matrix $R \in \mathbb{R}^{3 \times 3}$ from $M$. Similarly, we can obtain $\hat{R}$ from $\hat{M}$. We define the relative pose estimation loss function, which was used to measure the angular distance between $\hat{R}$ and $R$ in the study by Suwajanakorn et al. [SSTN18].

$$\mathcal{L}_{pose} = 2 \arcsin\left(\frac{1}{2\sqrt{2}}\|\hat{R} - R\|^2\right). \tag{6}$$

We found that the network has a slow convergence and easily falls into a local optimum solution using only $\mathcal{L}_{rebuild}$ and $\mathcal{L}_{pose}$. The reason for this is that it is possible to generate the same $M$ by using two different parameter sets of $\mu$, $\omega$, $\Phi$, and $\Theta$. To solve this problem, we compute an approximate solution of $\mu$, $\omega$, $\Phi$, and $\Theta$ denoted as $\hat{\mu}$, $\hat{\omega}$, $\hat{\Phi}$, and $\hat{\Theta}$, respectively. In general, $\hat{R}$ can be written more concisely as Rodrigues' rotation formula [Bel] as follows:

$$\hat{R} = (\cos\hat{\Theta})I + (\sin\hat{\Theta})[\hat{\omega}]_{\times} + (1 - \cos\hat{\Theta})(\hat{\omega} \otimes \hat{\omega}). \tag{7}$$

where $[\hat{\omega}]_{\times}$ is the cross product matrix of $\hat{\omega}$, $\hat{\omega} \otimes \hat{\omega}$ is the outer product, and $I$ is the identity matrix. We can easily obtain $\hat{\Theta}$ and $\hat{\omega}$ from equation (7).

$$\hat{\Theta} = \arccos\left(\frac{tr(\hat{R}) - 1}{2}\right), \tag{8}$$

where $tr(\hat{R})$ is the trace of $\hat{R}$.

$$\hat{\omega} = \frac{1}{2\sin\Theta}\begin{pmatrix} \hat{R}[2,1] - \hat{R}[1,2] \\ \hat{R}[0,2] - \hat{R}[2,0] \\ \hat{R}[1,0] - \hat{R}[0,1] \end{pmatrix}. \tag{9}$$
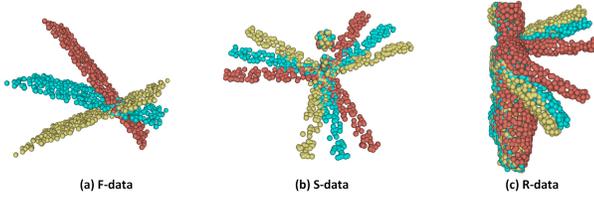
We extract $\hat{t} \in \mathbb{R}^3$ from $\hat{M}$, where $\hat{t}$ is a vector from the centroid of $PC_1$ to that of $PC_t$. $\hat{\Phi}$ is the length of $\hat{t}$ projected to $\hat{\omega}$.

$$\hat{\Phi} = \left|\frac{\hat{\omega} \cdot \hat{t}}{\|\hat{\omega}\|^2}\right|. \tag{10}$$

Finally, we solve linear algebraic equations to obtain $\hat{\mu}$, where $I$ is the identity Matrix.

$$\hat{\mu} = (\hat{R} - I)^{-1} \cdot (\hat{\Phi} \cdot \hat{\omega} - \hat{t}). \tag{11}$$

In the case of the general situation in which the start of an axis

**(a) F-data**  **(b) S-data**  **(c) R-data**

**Figure 4:** *Three frames of point cloud sequences in each of the three datasets used in our experiments.*

is not at the coordinate origin, we can obtain the same result. Then, we define the axis loss function as follows:

$$\mathcal{L}_{axis} = \|\hat{\Theta} - \Theta\|^2 + \|\hat{\Phi} - \Phi\|^2 + cosine(\hat{\omega}, \omega) \\ + \frac{1}{3}\|(\hat{\boldsymbol{R}} - I) \cdot \mu - (\hat{\Phi} \cdot \hat{\omega} - \hat{t})\|^2, \quad (12)$$

where $cosine(\hat{\omega}, \omega)$ is the cosine distance between $\hat{\omega}$ and $\omega$. Given that $(\hat{\boldsymbol{R}} - I)^{-1}$ is difficult to compute, network output $\mu$ should satisfy $(\hat{\boldsymbol{R}} - I) \cdot \mu - (\hat{\Phi} \cdot \hat{\omega} - \hat{t}) = 0$ rather than directly regressing $\mu$ by $\hat{\mu}$.

The final loss function is the weighted sum of $\mathcal{L}_{rebuild}$ (5), $\mathcal{L}_{pose}$ (6), and $\mathcal{L}_{axis}$ (12).

$$\mathcal{L}_{total} = \alpha_1 \cdot \mathcal{L}_{rebuild} + \alpha_2 \cdot \mathcal{L}_{pose} + \alpha_3 \cdot \mathcal{L}_{axis}. \quad (13)$$

In our implementation, $\alpha_1 = \alpha_2 = \alpha_3 = 1$ by default. In addition, we set a threshold for each loss function to tolerate error and improve robustness, because approximate solutions of $\hat{\mu}$, $\hat{\omega}$, $\hat{\Phi}$, and $\hat{\Theta}$ usually have errors with respect to the ground truth, especially in real scan data.

### 3.3. Motion Part Segmentation and Optimization

However, these hypotheses are redundant, such that we must merge similar motion hypotheses. We define that there is no movement if $\Theta$ and $\Phi$ are no more than $\mathcal{E}_\Theta$ (0.1 radians by default) and $\mathcal{E}_\Phi$ (0.05 length by default), respectively, considering noisy data and computation errors. Further, we define rotation as $\Theta > \mathcal{E}_\Theta$ and $\Phi < \mathcal{E}_\Phi$ and translation by contrast. The combination of rotation and translation satisfies the conditions that $\Theta > \mathcal{E}_\Theta$ and $\Phi > \mathcal{E}_\Phi$. We merge the remaining motion axes by comparing the diversity of $\mu$, $\omega$, $\Phi$, and $\Theta$.

We first define the metric distance between trajectory $\tau_i$ and motion hypotheses as the matching degree to merge motion hypotheses. Therefore, we utilize the motion hypotheses to rebuild the trajectory, denoted as $\hat{\tau}_i$. Then, the metric distance has two terms: rebuild loss and direction cosine distance.

$$dis_{match} = \frac{1}{t} \sum_{1 \leq m \leq t} (\|x_{i,m} - \hat{x}_{i,m}\|^2 + \alpha \cdot cosine(d_{i,m}, \hat{d}_{i,m})), \quad (14)$$

where $\alpha$ is 0.2 in our implementation. Then, we obtain refined motion axes and motion parts using an iteration algorithm. Inspired by Non-Maximum Suppression, we first set the refined motion axis set, $\mathcal{S}$, as empty. We assign the trajectories to candidate axes by

computing $dis_{match}$, which lower $dis_{match}$ means that the matching degree is high. Then, we choose the candidate axis that has the most votes. If $\mathcal{S}$ is empty, we add this candidate axis to $\mathcal{S}$. If $\mathcal{S}$ is not empty, we compute whether trajectories belonging to this candidate axis can also belong to an axis in $\mathcal{S}$. If not, we add this candidate axis to $\mathcal{S}$. We iterate these steps until all trajectories are assigned to a motion axis. We optimize coarse segments by examining whether the label of a trajectory is the same as most of the labels of its *k*-nearest neighbors.

### 4. Experiments

In this section, we first introduce the benchmark datasets used in the experiments (Sec. 4.1). We then compare our approach with state-of-the-art methods on the tasks of motion part segmentation, motion attribute estimation, and motion flow prediction (Sec. 4.2). Finally, we conduct several ablation experiments including the influence of different network designs, the effect of different loss function designs, hyperparameter settings, transferability, etc. (Sec. 4.3).
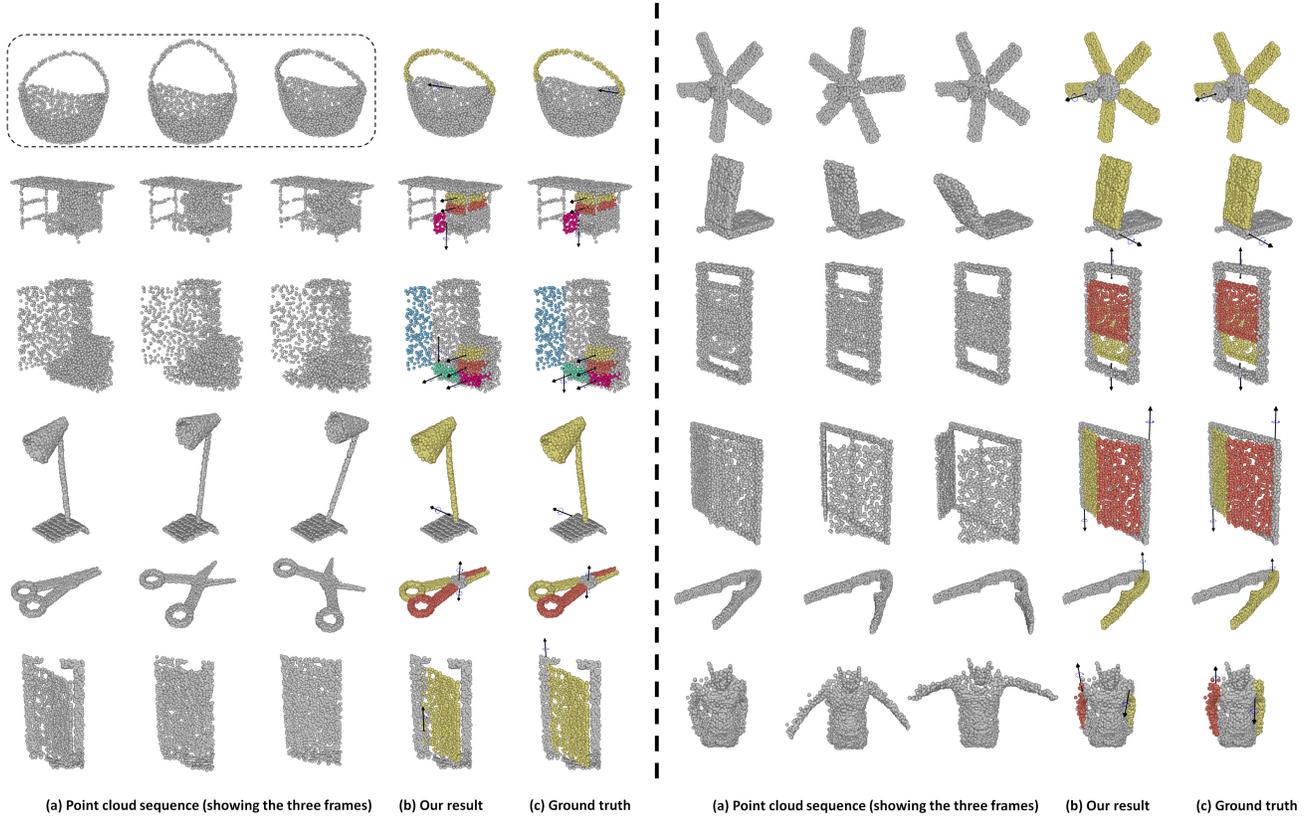
### 4.1. Dataset Settings

We evaluate our methods on three different datasets, including two synthetic sets (**f-data**, **s-data**) and one real set (**r-data**). A fake dataset (**f-data**) is generated randomly and automatically. **f-data** enables the neural network to learn the common motion patterns of random trajectories. We leverage two annotated datasets including the Motion dataset [WZS*19] and the PartNet dataset [MZC*19] to construct the second synthetic dataset (**s-data**). Xiang et al. [XQM*20] enriched the PartNet dataset with motion attributes. We choose 31 categories including cabinet, lamp, and window, to generate motion sequences. For the real data (**r-data**), we select certain reasonable motion sequences in the RBO dataset [MMEB18]. Moreover, Yi et al. [YHL*18] provided a small real scan dataset. Additionally, we scan some real data sequences in-house.

### 4.2. Comparison with State-of-the-Art Methods

We test our method against four alternatives, including both non-learning and learning approaches. Specifically, we compare our method with the traditional method of Yuan et al. [YLX*16], and three network-based approaches including Yi et al. [YHL*18], Wang et al. [WZS*19], and Yan et al. [YHY*19].

#### 4.2.1. Experiment Settings

To the best of our knowledge, we are the first to propose an approach using deep learning to extract part mobility from a point cloud sequence. To ensure a fair comparison, we implement a space-time co-segmentation baseline following Yuan et al. [YLX*16]. Moreover, we employ the trajectory generation algorithm [PB11] because of their method settings. Yi et al. [YHL*18] proposed a neural network architecture with three modules that propose correspondences, estimate 3D deformation flows, and perform segmentation. Their method input a pair of point clouds representing two different articulation states to segment motion parts and

**Figure 5:** *Results of our approach compared with the ground truth on both synthetic and real datasets. Rows 1–5 show the results on the synthetic dataset, and row 6 shows the results on the real dataset. We show three frames, first, middle, and lats, of the point cloud sequence. All results are shown in the first frame.*

| Method | $IoU$ | MD | OE | TA |
|---|---|---|---|---|
| Yuan et al. [YLX*16] | 0.71 | - | - | - |
| Wang et al. [WZS*19] | 0.67 | 0.051 | 0.055 | 0.92 |
| Ours | **0.87** | **0.032** | **0.027** | **0.99** |

**Table 1:** *Comparison between our method and those of Yuan et al. [YLX*16] and Wang et al. [WZS*19] on the tasks of motion part segmentation and motion attribute estimation in terms of IoU, MD, OE, and TA.*

| Method | $IoU$ | RI | EPE |
|---|---|---|---|
| Yi et al. [YHL*18] | 0.71 | 0.80 | 0.029 |
| Ours | **0.83** | **0.87** | **0.024** |

**Table 2:** *Comparison between our method and Yi et al. [YHL*18] on the tasks of motion part segmentation and motion flow prediction in terms of IoU, RI, and EPE.*
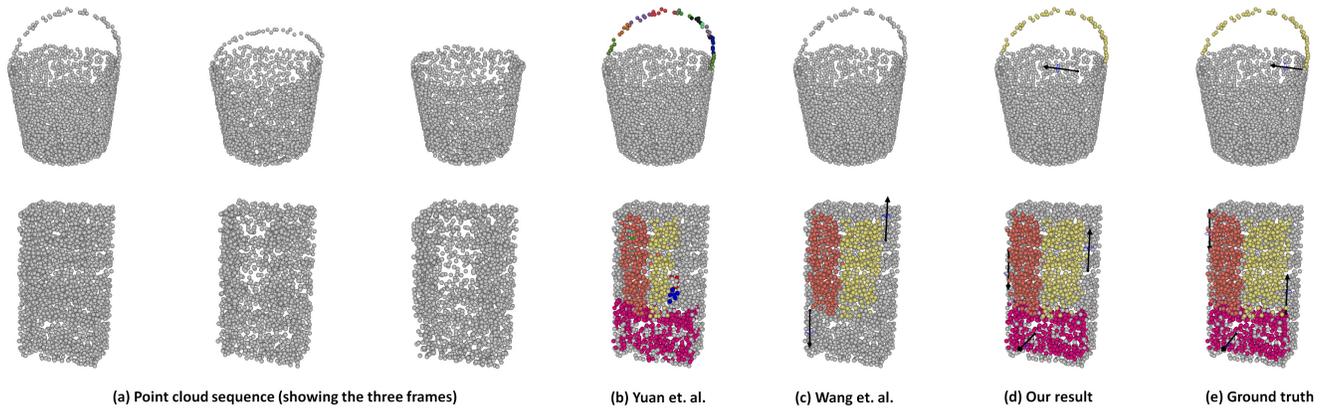
estimate 3D flows. To better provide a contrasting experiment between their approach and ours, we set $t = 2$ in our network settings. We also take a pair of point clouds as input to train and test our network. In terms of the approach of Wang et al. [WZS*19], we train and evaluate their network and ours using the same training and testing data. However, it must be mentioned that their network requires a single point cloud as input. Considering this difference, the data for their network are randomly sampled frames from point cloud sequences. For the comparison between Yan et al. [YHY*19] and our method, the data for their approach are the first frame of the point cloud sequence because their approach segments motion parts by predicting a temporal sequence from a single point cloud.

### 4.2.2. Evaluation Metrics

To evaluate motion part segmentation, we utilize Intersection over Union ($IoU$) defined in [YKC*16]. Moreover, we employ Rand Index (RI) defined in [CGF09] to compare with Yi et al. [YHL*18] and the mean Average Precision (mAP) defined in [YHY*19] to compare with [YHY*19]. To evaluate motion attribute estimation, we measure the Minimum Distance (MD), Orientation Error (OE), and Type Accuracy (TA), as introduced in [WZS*19], which are used for measuring the distance, angle, and motion type accuracy between the predicted motion axis line and the ground truth. To evaluate motion flow prediction, we employ 3D end-point-error defined in [YX16], which is the average $L_2$ distance between the predicted flow and the ground truth flow.

|  |  |  |  |  |
|---|---|---|---|---|
| (a) Point cloud sequence (showing the three frames) | (b) Yuan et. al. | (c) Wang et. al. | (d) Our result | (e) Ground truth |

**Figure 6:** *Our results are compared with those of Yuan et al. [YLX*16], and Wang et al. [WZS*19].*

### 4.2.3. Results and Analyses

Yuan et al. [YLX*16] adopted a clustering method to obtain local segments and propagated these segments to the neighboring frames. Finally, they merged all frames using a space–time segment grouping technique to obtain the final motion part segmentation. However, their method tends to have degraded performance when dealing with tiny motion parts. Moreover, their approach has difficulty in processing 3D shapes with many diverse large motions. Tab. 1 shows that our approach outperforms theirs in terms of *IoU*. They did not exploit the fact that motion parts have local characteristics. They generated a motion hypothesis by randomly choosing a number of trajectory triplets. It resulted in the creation of several fragments that were difficult to group together, especially in tiny parts. Additionally, the parameters of their approaches are often different depending on the case.

Wang et al. [WZS*19] parsed 3D shapes into part mobility by observing a single static point cloud based on a large, well-labeled dataset. Their method can only predict motion parameters of 3D shapes that have been seen in the training data because of the supervised learning based on a large, well-annotated dataset. It seems very easy to get confused in extracting motion parts from a single snapshot. For example, it would be challenging to discover whether a door opens on the right or left when it is closed. More importantly, their network is likely to consider that it is not a motion part when a door is closed. They used a similar matrix to segment motion parts. However, this type of method is also not conducive to dealing with tiny parts. In contrast to their method, which directly predicts an axis as a regression problem, we minimize the cosine distance for $\omega$ and solve linear algebraic equations to compute $\mu$, which makes our result more robust and stable than theirs (See Tab. 1).

Yi et al. [YHL*18] aimed to discover motion parts of objects and estimated 3D flows by analyzing the underlying articulation states and geometry of shapes. Their network architecture alternates between correspondence, deformation flow, and segmentation prediction iteratively in an ICP-like fashion. Theirs is a supervised learning method that estimates 3D point-wise flows and then segments the motion parts. In contrast to their method, we compute the motion axes and motion ranges and then segment motion parts using

| Method | *IoU* | mAP |
|---|---|---|
| Yan et al. [YHY*19] | 0.86 | 0.76 |
| Ours | **0.87** | **0.77** |

**Table 3:** *Comparison between our method and Yan et al. [YHY*19] on the task of motion part segmentation in terms of IoU and mAP.*

these motion parameters. The results are reported in Tab. 2, and they demonstrate that our method achieves higher *IoU* and RI, and lower EPE than theirs.

Yan et al. [YHY*19] predicted a temporal sequence of point-wise displacements from the input shape using LSTM. Then, the RPM-Net used these displacements to learn all the movable points. Next, they obtained a pointwise distance matrix by computing the Euclidean distance between these learned point features. Finally, they selected a clustering method to separate the points in the motion parts according to a distance matrix. They tended to study the part mobility from point cloud sequence using LSTM, but they still used a single point cloud. Tab. 3 shows that the results obtained from these models are close to each other. Although their method deal with motion segmentation well, it has the same defect as that of Wang et al. [WZS*19], which is due to supervised learning. It is not easy to imagine a motion sequence from a single point cloud.

### 4.3. Ablation Experiments

In this section, we first conduct several ablation studies to verify the influence of different backbone network designs and different loss function designs. Then, we evaluate our method on different datasets. Finally, we explore the hyperparameter settings in our method and compare them with their model size and timing.

### 4.3.1. Analysis of the Effects of Different Network Designs

To verify the effectiveness of PointRNN on the tasks of motion part segmentation and motion attribute estimation, we design two baselines to replace the network backbone. To build baseline1, we adopt a slightly modified version of Liu et al. [LYB19]. We merge

| Method | $IoU$ | MD | OE | $\Theta_e$ | $\Phi_e$ | $r_e$ | TA |
|---|---|---|---|---|---|---|---|
| baseline1 | 0.86 | 0.029 | 0.054 | 0.054 | 0.042 | 0.025 | 0.97 |
| baseline2 | 0.84 | 0.034 | 0.062 | 0.064 | 0.052 | 0.029 | 0.96 |
| Ours | **0.88** | **0.027** | **0.028** | **0.046** | **0.028** | **0.023** | **0.98** |

**Table 4:** *Compare our network with two other network designs on the tasks of motion part segmentation and motion attribute estimation in terms of IoU, MD, OE, $\Theta_e$, $\Phi_e$, $r_e$, and TA.*

| Method \ Dataset | **f-data** | **s-data** |
|---|---|---|
| Myronenko and Song [MS09] | **0.018** | 0.027 |
| Ours | 0.021 | **0.023** |

**Table 5:** *Comparison between our method and that of Myronenko and Song [MS09] in terms of $r_e$.*

| Loss Function | | MD | OE | $\Theta_e$ | $\Phi_e$ | $r_e$ | TA |
|---|---|---|---|---|---|---|---|
| $\mathcal{L}_{axis}$ | -**f** | 0.035 | 0.017 | 0.058 | 0.026 | 0.061 | 0.96 |
| $\mathcal{L}_{rebuild} + \mathcal{L}_{pose}$ | -**f** | 0.031 | 0.017 | 0.051 | 0.036 | 0.086 | 0.90 |
| $\mathcal{L}_{pose} + \mathcal{L}_{axis}$ | -**f** | 0.032 | 0.016 | 0.053 | **0.022** | 0.045 | 0.94 |
| $\mathcal{L}_{total}$ | -**f** | **0.027** | **0.014** | **0.047** | 0.023 | **0.020** | **0.97** |
| $\mathcal{L}_{axis}$ | -**s** | 0.043 | 0.030 | 0.056 | 0.030 | 0.045 | 0.97 |
| $\mathcal{L}_{rebuild} + \mathcal{L}_{pose}$ | -**s** | 0.048 | 0.033 | 0.057 | 0.034 | 0.061 | 0.93 |
| $\mathcal{L}_{pose} + \mathcal{L}_{axis}$ | -**s** | 0.037 | 0.029 | 0.054 | 0.030 | 0.034 | **0.99** |
| $\mathcal{L}_{total}$ | -**s** | **0.032** | **0.027** | **0.046** | **0.027** | **0.023** | **0.99** |

**Table 6:** *Effectiveness of different loss function designs on the task of motion attribute estimation. The first four rows show the results on **f-data**, and the last four rows show the results on **s-data**.*

| Dataset | $IoU$ | MD | OE | $\Theta_e$ | $\Phi_e$ | $r_e$ | TA |
|---|---|---|---|---|---|---|---|
| **f-data** | - | **0.027** | **0.014** | 0.047 | **0.023** | **0.020** | 0.97 |
| **s-data** | **0.87** | 0.032 | 0.027 | **0.046** | 0.027 | 0.023 | **0.99** |
| **r-data** | 0.79 | 0.058 | 0.036 | 0.079 | 0.037 | 0.058 | **0.99** |

**Table 7:** *Performance of our method on three datasets on the tasks of motion part segmentation and motion attribute estimation in terms of IoU, MD, OE, $\Theta_e$, $\Phi_e$, $r_e$, and TA.*

the point cloud sequence into a single point cloud and feed it to a PointNet++ to extract features. Then, we minimize the loss function $\mathcal{L}_{total}$ to estimate motion axis and motion range. Simultaneously, we design baseline2 by referring to Yan et al. [YHY*19]. In baseline2, we extract $t$ global features per frame using a shared PointNet++. Then, we adopt an LSTM Network to encode the time information. Similarly, we utilize $\mathcal{L}_{total}$ as the loss function. Different from previous work, our method also estimates the motion ranges including translation distance and rotation angle. Therefore, we three additional novel metrics, $\Theta_e$, $\Phi_e$, and $r_e$. $\Theta_e$ is the rotation error between predicted rotation angle and ground truth rotation angle. Similarly, $\Phi_e$ is the translation error between predicted translation distance and ground truth translation distance. $r_e$ is the error of rebuilding the last frame from the first frame, which is used for a comprehensive evaluation. The results in 4 demonstrate that our network can achieve better performance than the other two baselines because our network can take advantage of detailed spatiotemporal information from motion sequence.

### 4.3.2. Effectiveness Verification of Self-Supervised Method

Our method is self-supervised because our data are not manually annotated, either in training or testing. Moreover, supervised information is automatically generated from the characteristic distribution of the data. To this end, we define a loss function, $\mathcal{L}_{axis}$. We solve an absolute orientation problem to generate a motion matrix ($\hat{\boldsymbol{M}}$) using the method of Myronenko and Song [MS09]. Then, we parse the supervised information ($\hat{\mu}, \hat{\omega}, \hat{\Phi}, \hat{\Theta}$) from $\hat{\boldsymbol{M}}$.

In this section, we conduct an ablation study to compare with the performance of [MS09], since not improving on [MS09] would have meant the neural network architecture would not be necessary. The results in Tab. 5 show that our method is better than [MS09] on **s-data**, but worse on **f-data**. It demonstrates a hypothesis that our method can handle more complex data. **s-data** usually contains two or more motion parts and a static part. However, **f-data** only has one motion part. Obviously, **s-data** is more complex than **f-data**. The output of [MS09] is an analytical solution. If the grouped trajectories cover different motion parts, it will output the wrong answer. However, our loss function contains three components in-
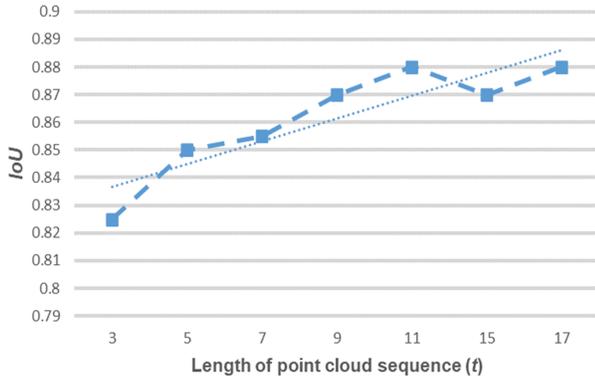
stead of only utilizing $\mathcal{L}_{axis}$, and the deep learning method learns common motion features and outputs an approximation solution. It might slightly alleviate the problem of the wrong solution on **s-data**.

### 4.3.3. Effect of Different Loss Function Designs

To train our network using self-supervised learning, we define a novel self-supervised loss function including three components. To analyze the influence of different loss function designs, we experiment with an ablated version of our network with four combinations of all loss functions on two datasets including **f-data** and **s-data**. We utilize the loss function to train the network on the task of motion hypotheses generation. Therefore, different loss designs mainly influence the network performance with respect to axis generation. Thus, we train and test our network on the task of motion attribute estimation to focus on examining the quality of the predicted motion axes. Tab. 6 lists the results of including $\mathcal{L}_{axis}$, $\mathcal{L}_{rebuild} + \mathcal{L}_{pose}$, $\mathcal{L}_{pose} + \mathcal{L}_{axis}$, and $\mathcal{L}_{total}$. The results demonstrate that our method achieves optimal performance for motion axis estimation, thereby verifying its optimal effect when using $\mathcal{L}_{rebuild}$, $\mathcal{L}_{pose}$, and $\mathcal{L}_{axis}$ together. Moreover, the shown results are applicable to both **f-data** and **s-data**. The reason for this is that it is possible to generate the same $\boldsymbol{M}$ using two different parameter sets of $\mu$, $\omega$, $\Phi$, and $\Theta$. There is an interactive constraint between different components of the loss function.

### 4.3.4. Performance on Three Datasets

We have three datasets used in experiments, and these datasets have their own unique features. The difficulty of **f-data** stems from its randomness. Each element of **f-data** has totally different motion parameters, including motion axes and motion ranges. **s-data** has

**Figure 7:** *Impact of the length of a point cloud sequence, where a higher IoU value is better. The x-axis corresponds to the length of the point cloud sequence (t).*
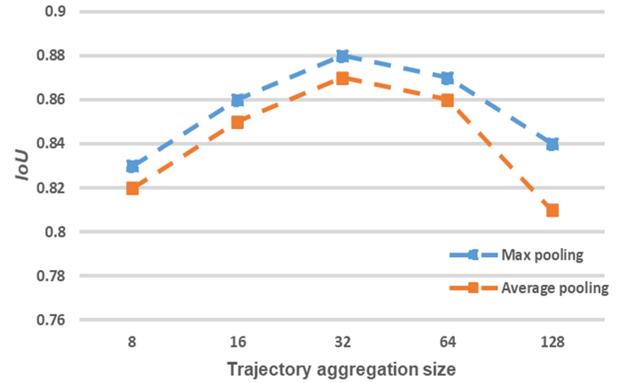


**Figure 8:** *Trajectory aggregation analysis, where a higher IoU value is better. The x-axis corresponds to different hyperparameters of trajectory aggregation size (k).*

| category | $IoU$ | MD | OE | $\Theta_e$ | $r_e$ | TA |
|---|---|---|---|---|---|---|
| Fan     -**f** | 0.89 | 0.014 | 0.014 | 0.032 | 0.028 | 1 |
| Laptop  -**f** | 0.95 | 0.018 | 0.013 | 0.035 | 0.021 | 1 |
| Scissor -**f** | 0.79 | 0.022 | 0.017 | 0.087 | 0.031 | 1 |
| Fan     -**s** | 0.93 | 0.012 | 0.014 | 0.028 | 0.024 | 1 |
| Laptop -**s** | 0.97 | 0.017 | 0.014 | 0.031 | 0.017 | 1 |
| Scissor -**s** | 0.84 | 0.019 | 0.015 | 0.069 | 0.028 | 1 |

**Table 8:** *Experiment on the motion feature learning. The first three rows show the results, which train our network on **f-data**, and the last three rows show the results, which train our network on **s-data**. Our method can handle unseen categories, because motion patterns and type are learned by the network.*

various categories. Most 3D shapes in **s-data** have two or more motion parts. **r-data** generally has a single partial view of 3D shapes with more noise. Therefore, we need to verify the performance on these datasets, respectively. We train and test our network using three different datasets. The results are reported in Tab. 7, and results demonstrate that not only can our approach learn motion parameters from synthetic data, but it can also be applied to real noisy data.

### 4.3.5. Analysis on Motion Feature Learning

Supervised learning methods tend to require that object categories in test data must have been seen in train data. However, our approach is a self-supervised deep learning algorithm, and our network can estimate the motion axes by learning the feature representation of trajectories other than only by learning geometrical characteristics. Therefore, we aim to verify that our network has the ability to learn motion features. We designed an experiment in which we train our network on **f-data** and **s-data**, and then test it on **s-data**. We present three categories with different levels of complexity results in Tab. 8. Two sets of the result that train our network on the different datasets are similar. The results demonstrate that the network can produce correct motion axes, even though it

did not see any real objects. It implies that our network can learn motion patterns from trajectories.
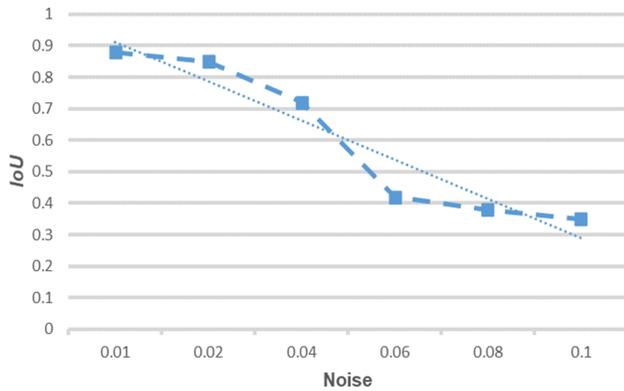
### 4.3.6. Analysis on the Length of Point Cloud Sequence

Considering that our motivation is to learn part mobility from a point cloud sequence, the number of frames ($t$) is an important hyperparameter in our approach. Thus, we must determine how many frames are sampled from one motion sequence. To achieve this, we conduct an ablation experiment on $t$ to discuss the impact of the length of the point cloud sequence. We reduce the number of frames from 17 to 3 for sufficient verification. The results are depicted in Fig. 7. As the number of frames increases, the performance gain is consistent but slower and slower. Hence, we adopt 11 frames as input to balance performance and efficiency.

### 4.3.7. Effect of Trajectory Aggregation

The aggregation of trajectories is of great importance in PointRNN, as it describes a motion part consisting of a group of similar motion trajectories. Thus, it is important to analyze how different aggregation designs influence performance. We report five different sizes with two types of aggregations: max and average, as shown in Fig. 8. As the aggregation radius ($k$) increases, PointRNN improves until it peaks at approximately 32. The results demonstrate that it is challenging to learn the features when clustering too few trajectories in a local region. In contrast, a larger region usually contains more than one motion part and results in decreased performance. Moreover, the max-pooling operation achieves better results than average by aggregating the features in local regions.

### 4.3.8. Analysis on the Quality of Trajectory

We also discuss the impact of different qualities of trajectory. Given a point cloud sequence, we first generate a set of trajectories. Then, we add different levels of random noise to these trajectories. Fig. 9 shows the relationship between $IoU$ and noise. Our method maintained more than 0.7 ($IoU$) after adding 0.04 random noise. As the

**Figure 9:** *Effect of different qualities of trajectory. The result illustrates that the performance decreases with increasing of noise.*

| Method | Model size | Timing |
|---|---|---|
| Yuan et al. [YLX*16] | - | -/180s |
| Wang et al. [WZS*19] | 86.7MB | 38h/7.32s |
| Yi et al. [YHL*18] | 78MB | 35.5h/4.88s |
| Yan et al. [YHY*19] | 153.3 MB | 13.2h/ 0.55s |
| Ours | **25.2MB** | **11.6h/0.43s** |

**Table 9:** *Comparison of Model size and processing time. Our model is smaller and faster.*

noise degree increases, the performance declines sharply. The results demonstrate that our method can tolerate noise to a certain extent, but it will fail when under the addition of excessive noise.

### 4.3.9. Model Size and Speed

Our proposed model is highly efficient (See Tab. 9), as it leverages sparsity in point clouds by using clusters and does not require many stages like those of Wang et al. [WZS*19] and Yi et al. [YHL*18]. Compared to previous methods, our model is more than $3\times$ smaller in size and more than $10\times$ times faster than those of Wang et al. [WZS*19], and Yi et al. [YHL*18]. While the method in Yuan et al. [YLX*16] is a non-learning method, it costs approximately 3 min to segment a 3D shape without giving motion parameters. Furthermore, our model is more than $6\times$ smaller in size than that of Yan et al. [YHY*19]. Moreover, we employ the model size to compare the size of the training parameters of different models to evaluate the memory consumption of the algorithm.

### 5. Limitations and Future Work

There are several avenues for future research. First, our method is based on trajectories such that if the quality of trajectories is extremely poor with many fragments, the performance decreases. Second, if two parts have exactly the same motion, e.g., two drawers in a cabinet have the same motion, our method considers that they have the same motion and should be one motion part. Third,

for hierarchical mobility extraction, we can segment the 3D object into different motion parts, but the motion axis describes a compound movement. For example, the motion of a bulb holder is a composite movement with that of a lamp post. In the future, it would be interesting to infer parts and motions and discover common articulation patterns from various sensor data. Moreover, our method relies on a pre-registration of point clouds using an ICP-like optimization. There are also many deep learning methods for point cloud registration, such as [PMR*20]. It is deserved to design a network encompassing both trajectory generation and motion attribute estimation in an end-to-end fashion in future work. To facilitate future research and reproduce our method more easily, the source code and dataset are available on Github: https://github.com/AGithubRepository/PartMobility.

### 6. Conclusion

In this paper, we introduce a self-supervised method for parsing 3D shapes into several motion parts with motion parameters from point cloud sequences. We first transform point cloud sequences into trajectories. Then, we propose a novel neural network architecture called PointRNN to extract the feature representations from trajectories and produce motion hypotheses. Finally, we propose an iterative algorithm for segmenting motion parts. In contrast to previous studies, our approach first predicts part motion parameters including motion axes and motion ranges and then segments motion parts guided by these motion hypotheses. We experimentally demonstrated that our approach yields significantly better results compared with both the traditional and network-based methods.

### Acknowledgement

### References

[Bel]    BELONGIE S.: Rodrigues' rotation formula. From MathWorld–A Wolfram Web Resource, created by Eric W. Weisstein. https://mathworld.wolfram.com/RodriguesRotationFormula.html. 5

[BM92]   BESL P. J., MCKAY N. D.: A method for registration of 3-d shapes. *IEEE Trans. Pattern Anal. Mach. Intell. 14*, 2 (Feb. 1992), 239–256. 2

[BPDG19]    BEHL A., PASCHALIDOU D., DONNE S., GEIGER A.: Point-flownet: Learning representations for rigid motion estimation from point clouds. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (June 2019). 3

[CGF09]   CHEN X., GOLOVINSKIY A., FUNKHOUSER T.: A benchmark for 3d mesh segmentation. *Acm Transactions on Graphics 28*, 3 (2009), 1–12. 7

[CZ08]   CHANG W., ZWICKER M.: Automatic registration for articulated shapes. *Computer Graphics Forum 27*, 5 (2008), 1459–1468. 2

[FRA11]   FAYAD J., RUSSELL C., AGAPITO L.: Automated articulated structure and 3d shape recovery from point correspondences. In *2011 International Conference on Computer Vision* (Nov 2011), pp. 431–438. 2

[HLRB13] HERMANS T., LI F., REHG J. M., BOBICK A. F.: Learning contact locations for pushing and orienting unknown objects. In *2013 13th IEEE-RAS International Conference on Humanoid Robots (Humanoids)* (Oct 2013), pp. 435–442. 1

[HLVK*17] HU R., LI W., VAN KAICK O., SHAMIR A., ZHANG H., HUANG H.: Learning to predict part mobility from a single static snapshot. *ACM Trans. Graph. 36*, 6 (Nov. 2017), 227:1–227:13. 3

[HS97] HOCHREITER S., SCHMIDHUBER J.: Long short-term memory. *Neural Computation 9*, 8 (1997), 1735–1780. 3

[HSvK18] HU R., SAVVA M., VAN KAICK O.: Functionality representations and applications for shape analysis. *Computer Graphics Forum 37*, 2 (2018), 603–624. 1

[JSGC15] JAIMEZ M., SOUIAI M., GONZALEZ-JIMENEZ J., CREMERS D.: A primal-dual framework for real-time dense rgb-d scene flow. In *2015 IEEE International Conference on Robotics and Automation (ICRA)* (May 2015), pp. 98–104. 2

[KL17] KLOKOV R., LEMPITSKY V.: Escape from cells: Deep kd-networks for the recognition of 3d point cloud models. In *The IEEE International Conference on Computer Vision (ICCV)* (Oct 2017). 3

[KLAK16] KIM Y., LIM H., AHN S. C., KIM A.: Simultaneous segmentation, estimation and analysis of articulated motion from dense point cloud sequence. In *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* (Oct 2016), pp. 1085–1092. 2

[LBS*18] LI Y., BU R., SUN M., WU W., DI X., CHEN B.: Pointcnn: Convolution on x-transformed points. In *Advances in Neural Information Processing Systems 31*. Curran Associates, Inc., 2018, pp. 820–830. 1

[LQG19] LIU X., QI C. R., GUIBAS L. J.: Flownet3d: Learning scene flow in 3d point clouds. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (June 2019). 3

[LSZ*19] LIU M., SHI Y., ZHENG L., XU K., HUANG H., MANOCHA D.: Recurrent 3d attentional networks for end-to-end active object recognition. *Computational Visual Media 5*, 01 (2019), 92–104. 1

[LWL*16] LI H., WAN G., LI H., SHARF A., XU K., CHEN B.: Mobility fitting using 4d ransac. *Computer Graphics Forum 35*, 5 (2016), 79–88. 1

[LYB19] LIU X., YAN M., BOHG J.: Meteornet: Deep learning on dynamic 3d point cloud sequences. In *The IEEE International Conference on Computer Vision (ICCV)* (October 2019). 8

[MMEB18] MARTLN-MARTLN R., EPPNER C., BROCK O.: The rbo dataset of articulated objects and interactions, 2018. 6

[MS09] MYRONENKO A., SONG X. B.: On the closed-form solution of the rotation matrix arising in computer vision problems. *CoRR abs/0904.1613* (2009). 5, 9

[MS15] MATURANA D., SCHERER S.: 3d convolutional neural networks for landing zone detection from lidar. In *2015 IEEE International Conference on Robotics and Automation (ICRA)* (May 2015), pp. 3471–3478. 3

[MTFA15] MYERS A., TEO C. L., FERMLLER C., ALOIMONOS Y.: Affordance detection of tool parts from geometric features. In *2015 IEEE International Conference on Robotics and Automation (ICRA)* (May 2015), pp. 1374–1381. 1

[MYY*13] MITRA N. J., YANG Y.-L., YAN D.-M., LI W., AGRAWALA M.: Illustrating how mechanical assemblies work. *Commun. ACM 56*, 1 (Jan. 2013), 106–114. 3

[MZC*19] MO K., ZHU S., CHANG A. X., YI L., TRIPATHI S., GUIBAS L. J., SU H.: PartNet: A large-scale benchmark for fine-grained and hierarchical part-level 3D object understanding. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (June 2019). 6

[PB11] PAPAZOV C., BURSCHKA D.: Deformable 3d shape registration based on local similarity transforms. *Computer Graphics Forum 30*, 5 (2011), 1493–1502. 2, 3, 6

[PMR*20] PAIS G., MIRALDO P., RAMALINGAM S., GOVINDU V., NASCIMENTO J., CHELLAPPA R.: 3dregnet: A deep neural network for 3d point registration. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2020), 7191–7201. 11

[QYSG17] 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 30*. Curran Associates, Inc., 2017, pp. 5099–5108. 1, 3, 4

[RLA*19] ROBERTS R., LEWIS J. P., ANJYO K., SEO J., SEOL Y.: Optimal and interactive keyframe selection for motion capture. *Computational Visual Media 5*, 02 (2019), 172–191. 1

[SA07] SORKINE O., ALEXA M.: As-rigid-as-possible surface modeling, 2007. 2

[SMW06] SCHAEFER S., MCPHAIL T., WARREN J.: Image deformation using moving least squares. *ACM Trans. Graph. 25* (07 2006), 533–540. 3

[SSTN18] 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 31*. Curran Associates, Inc., 2018, pp. 2059–2070. 5

[VRS14] VOGEL C., ROTH S., SCHINDLER K.: View-consistent 3d scene flow estimation over multiple frames. In *Computer Vision – ECCV 2014* (Cham, 2014), Springer International Publishing, pp. 263–278. 2

[WZS*19] WANG X., ZHOU B., SHI Y., CHEN X., ZHAO Q., XU K.: Shape2motion: Joint analysis of motion parts and attributes from 3d shapes. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (June 2019). 2, 3, 6, 7, 8, 11

[XQM*20] XIANG F., QIN Y., MO K., XIA Y., ZHU H., LIU F., LIU M., JIANG H., YUAN Y., WANG H., YI L., CHANG A. X., GUIBAS L. J., SU H.: SAPIEN: A simulated part-based interactive environment. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (June 2020). 6

[YHL*18] YI L., HUANG H., LIU D., KALOGERAKIS E., SU H., GUIBAS L.: Deep part induction from articulated object pairs. *ACM Trans. Graph. 37*, 6 (Dec. 2018), 209:1–209:15. 2, 3, 6, 7, 8, 11

[YHY*19] YAN Z., HU R., YAN X., CHEN L., VAN KAICK O., ZHANG H., HUANG H.: Rpm-net: Recurrent prediction of motion and parts from point cloud. *ACM Transactions on Graphics (Proceedings of SIGGRAPH ASIA 2019) 38*, 6 (2019), 240:1–240:15. 2, 3, 4, 6, 7, 8, 9, 11

[YKC*16] 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. 7

[YLX*16] YUAN Q., LI G., XU K., CHEN X., HUANG H.: Space-time co-segmentation of articulated point cloud sequences. *Computer Graphics Forum 35*, 2 (2016), 419–429. 2, 3, 6, 7, 8, 11

[YP06] YAN J., POLLEFEYS M.: A general framework for motion segmentation: Independent, articulated, rigid, non-rigid, degenerate and non-degenerate. In *Computer Vision – ECCV 2006* (Berlin, Heidelberg, 2006), Springer Berlin Heidelberg, pp. 94–106. 2, 4

[YX16] YAN Z., XIANG X.: Scene flow estimation: A survey. *arXiv preprint arXiv:1612.02590* (2016). 7