

Efficient MSPSO Sampling for Object Detection and 6D Pose Estimation in 3D Scenes

Xuejun Xing, Jianwei Guo, Liangliang Nan, Qingyi Gu, Xiaopeng Zhang, Dong-Ming Yan

Abstract—The *point pair feature* (PPF) is widely used in manufacturing for estimating 6D poses. The key to the success of PPF matching is to establish correct 3D correspondences between the object and the scene, *i.e.*, finding as many valid similar point pairs as possible. However, efficient sampling of point pairs has been overlooked in existing frameworks. In this paper, we propose a revised PPF matching pipeline to improve the efficiency of 6D pose estimation. Our basic idea is that the valid scene reference points are lying on the object's surface and the previously sampled reference points can provide prior information for locating new reference points. The novelty of our approach is a new sampling algorithm for selecting scene reference points based on the multi-subpopulation particle swarm optimization (MSPSO) guided by a probability map. We also introduce an effective pose clustering and hypotheses verification method to obtain the optimal pose. Moreover, we optimize the progressive sampling for multi-frame point clouds to improve processing efficiency. The experimental results show that our method outperforms previous methods by 6.6%, 3.9% in terms of accuracy on the public DTU and LineMOD datasets, respectively. We further validate our approach by applying it in a real robot grasping task.

Index Terms—6D pose estimation, Multi-subpopulation particle swarm optimization, Point pair features, 3D point cloud

I. INTRODUCTION

Object detection and 6D pose estimation is a key component of various applications in intelligent manufacturing and industrial robotics. The purpose of 6D pose estimation is to determine an object pose that is represented by a 3×3 rotation matrix \mathbf{R} and a 3D translation vector \mathbf{t} . In particular, the object pose is the transformation that converts a 3D point \mathbf{p}_o in the local object coordinate system into the corresponding 3D point \mathbf{p}_s in the scene coordinate system, *i.e.*, $\mathbf{p}_s = \mathbf{R}\mathbf{p}_o + \mathbf{t}$.

Researchers in the machine vision and robotics communities have presented many methods on 6D pose estimation. In general, the key technologies are based on three directions:

Template matching. The template matching method converts the object into a series of templates and searches the position of each template on the entire scene to obtain the object pose [1]–[7]. These methods have the advantages of

strong real-time performance and easy implementation [5], [7]. However, their performance of instances detection and pose estimation in complex scenes is compromised due to changes in light and scene occlusion [4], [8]–[10].

Deep-learning-based methods. In recent years, deep learning has been developed rapidly in the field of image processing and analysis [11]–[15]. Similarly, many pose estimation methods based on deep learning are proposed for point cloud [16]–[19], RGB [7], [10], [20]–[24] or RGB-D images [25]–[32]. These methods show great potential in 6D pose estimation. However, so far there are few methods focusing on unstructured point clouds. Hagelskjaer et al. [16] performed semantic segmentation on the input point cloud based on PointNet, then used traditional methods to estimate the 6D pose of the segmented instances. The methods of [17]–[19] require a segmented point cloud as input, which depends on the segmentation using RGB images. Furthermore, the accuracy of these methods is still lower than that of deep learning models based on RGB-D images and is not yet saturated [18], [19]. On the other hand, the combination of deep learning models and traditional methods for 6D pose estimation has also achieved competitive results [7], [10], [16], [31], [32].

Feature descriptor matching. Another kind of methods mainly used local feature descriptors [33]–[35] to establish similar point pairs between object and scene, which are used to generate the approximate 6D poses. These methods has poor prediction performance when the object has indistinguishable local shape features, such as planar and regular surfaces.

Drost et al. [36] proposed a well-known 6D pose estimation method, which uses point pair feature (PPF) to achieve global model description and local matching. The method firstly describes the 3D object as a PPF Hash table, then samples the reference points in the scene and votes on them to generate hypothetical poses by finding the similar features in the Hash table, and finally generates the optimal pose by a clustering. Since the method was proposed, many variants [9], [37]–[41] have extended it to obtain better performance. In the PPF-based approach, the voting process is the key to the success of finding optimal poses, but this process is computationally expensive with a complexity of $O(n^2)$ [9], [34] where n is the number of scene points. Previous PPF-based methods usually generate a large number of reference points by uniform or random sampling. They do not consider the voting information of the reference points, thus reducing the efficiency of the entire voting process. Besides, invalid scene reference points have a negative impact on the robustness to noise, occlusion, and cluttered scenes, as well as increasing calculation time for pose clustering and hypotheses verification. Second, PPF-

X. Xing is with School of Artificial Intelligence, University of Chinese Academy of Sciences, Beijing 100049, China, and NLPR, Institute of Automation, Chinese Academy of Sciences, Beijing 100190, China.

J. Guo, X. Zhang, D.-M. Yan are with NLPR, Institute of Automation, Chinese Academy of Sciences, Beijing 100190, China, and School of Artificial Intelligence, University of Chinese Academy of Sciences, Beijing 100049, China. (Jianwei Guo is the corresponding author. E-mail: jianwei.guo@nlpr.ia.ac.cn.)

L. Nan is with Delft University of Technology, Delft 2628BL, Netherlands.

Q. Gu is with Institute of Automation, Chinese Academy of Sciences, Beijing 100190, China.

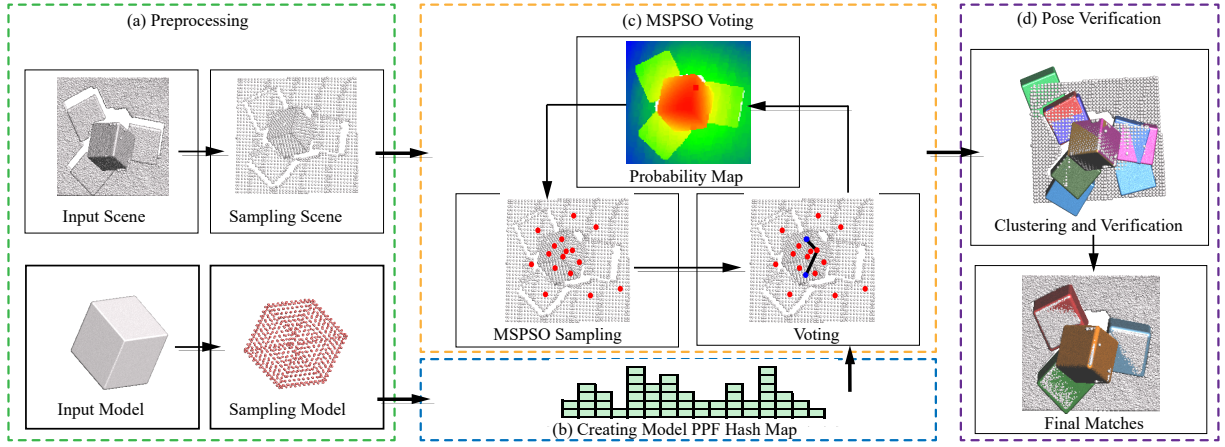


Fig. 1. Illustration of the proposed pipeline: (a) In preprocessing, the input object and scene point cloud are downsampled; (b) The model PPF hash map is generated using the point pairs of the object; (c) We progressively sample scene reference points by the MSPSO algorithm based on a probability cloud map, then perform PPF matching and Hough voting to generate pose hypotheses; (d) An improved pose clustering and verification is conducted to return final poses by filtering out duplicated and invalid poses.

based approaches are imperfect for multi-frame point clouds that contain a sequence of point clouds in a time series and are often encountered in industrial applications. The consistency information between frames is often ignored by the previous 6D pose estimation approaches for point clouds.

In this paper, we propose a new PPF-based framework for 6D pose estimation with an emphasis on intelligent scene reference points sampling. Our observation is based on that a large number of votes indicates that the points around the reference point have a high probability on the instance's surface, while a few votes indicate low probability. Thus the number of votes can be used to construct a probability map of the instance position in the scene. We propose to use the MSPSO algorithm [42], [43] for finding valid reference points which are located on the instance's surface. As a result, we propose the method which can combine MSPSO and probability map to iteratively sample scene reference points (see Fig. 1 (c)). In this way, we can improve the robustness of our algorithm to occlusion and complex scenes and reduce the computational burden. Besides, we further extend this framework to multi-frame point clouds by considering similar object information between adjacent scene frames. It can utilize the prior knowledge in previous frames to improve the matching efficiency in real industrial applications. In summary, the main contributions of this work include:

- A novel MSPSO sampling approach based on a probability map for progressively selecting valid scene reference points, which can simultaneously improve the efficiency and accuracy of 6D pose estimation.
- An effective pose clustering and hypothesis verification strategy to improve the robustness of 6D pose estimation of point clouds.
- A simple probability map-based approach for object detection from multi-frame point clouds, which exploits the correlation between frames and can significantly reduce the number of invalid scene reference points, thus speed up the matching process.

II. METHODOLOGY OVERVIEW

Our input includes an *object* \mathcal{O} represented by either a surface mesh or a point cloud, and a real-scanned *scene* point cloud \mathcal{S} . Each occurrence of the object in the scene is referred to one of its *instances* \mathcal{I} . The goal of our approach is to correctly determine the 6D poses of all instances in the scene. We first shortly explain the main underlying ideas and concepts of the original *Drost-PPFM* method [36].

A. Pose Estimation based on Point Pair Features

Point Pair Feature is a global feature descriptor, which describes the relative relationship between the position and direction of the two points. Given a reference point \mathbf{p}_r and a target point \mathbf{p}_t with normal \mathbf{n}_r and \mathbf{n}_t , respectively, the **PPF** is a 4-d vector which is defined as:

$$\mathbf{PPF}(\mathbf{p}_r, \mathbf{p}_t) = (||\mathbf{d}||_2, \angle(\mathbf{n}_r, \mathbf{d}), \angle(\mathbf{n}_t, \mathbf{d}), \angle(\mathbf{n}_r, \mathbf{n}_t)), \quad (1)$$

where $\mathbf{d} = \mathbf{p}_t - \mathbf{p}_r$, $\angle(\mathbf{a}, \mathbf{b})$ denotes the angle between vectors. **Drost-PPFM** detects instances and estimates their poses via two main phases: offline global model description and online matching. In the offline phase, the object is described by a *model PPF hash table*, which is created from **PPF**s between each points pair. Specifically, this method calculates the **PPF** of each point pair and discretizes the distance and angle elements in **PPF** with a quantization step size of Δ_{dist} and Δ_{angle} , respectively. Then the point pairs with equal discrete feature vectors are recorded in the same hash node.

In the online phase, *Drost-PPFM* mainly performs feature matching and pose clustering. First, the scene point cloud is down-sampled to generate a point set \mathcal{S}' , of which 1/5 points are uniformly sampled as the scene reference points. To generate hypothetical poses, the scene **PPF**s, which are calculated by each scene reference point and all other points in \mathcal{S}' , are matched by the object model **PPF** hash map and the pose votes are cast by performing Hough voting. Finally, the hypothetical poses are clustered to improve robustness and the pose in the cluster with the highest accumulated weight

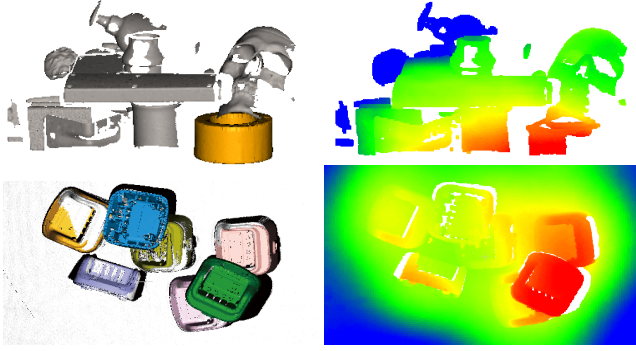


Fig. 2. Visual results of our generated pose estimation guided by the cloud probability map. Left: final pose estimation where the detected objects are highlighted in color. Right: the probability map where the warmer color indicates a higher probability for locating the objects.

is returned as the optimal pose. More details about these procedures can be found in the original paper [36].

B. Overview of Our Approach

We present a new 6D pose estimation algorithm, called *PPFM-MSPSO*, which consists of three main components (see Fig. 1). First, we down-sample the input model and scene point cloud, then generate a model description hash map. Second, we generate the hypothetical poses by progressively sampling scene reference points based on an MSPSO algorithm and a probability map. The sampling method can exploit the voting information of each reference point to guide the search direction of MSPSO and significantly improve the accuracy of pose estimation by increasing the ratio of valid scene reference points. Finally, a pose verification operation is performed to filter out false poses due to the interference of noise and background and predict appropriate 6D poses of the target object.

III. METHOD

A. Preprocessing

Generally, the input 3D object and scene point clouds have a high density with a large number of points. The input is usually downsampled to obtain a sparse set of points to speed up the matching process. However, traditional randomly or uniformly downsampling methods based on voxelization may ignore some important geometrical shape features. To solve this problem, we adopt an adaptive multi-scale voxel downsampling method [9], which takes the normal vectors of the point cloud into account. Specifically, we first discretize the point cloud by creating a multi-resolution grid structure. Then in each voxel cell at each level (in a fine-to-coarse order), the similar points, whose normals angle difference is less than a threshold θ , are merged. As a result, the geometrical features (e.g., edges or large curvature surface) are more preserved to sample discriminative points.

B. MSPSO Voting Module

The selection of scene PPFs has a great influence on the performance of the algorithm. With the unknown number and

locations of the object instances in the scene, a large number of reference scene points need to be sampled to cover as many instances as possible, leading to a high computational cost of voting. In addition, invalid reference points generate many invalid votes that increase the amount of calculation for hypothetical poses clustering and verification, as well as reduce the matching accuracy.

Instead, we present a new voting module where the scene reference points are progressively sampled with the MSPSO algorithm to better exploit the knowledge obtained in the previous voting. A probability map is utilized to guide the MSPSO sampling. Thus, the MSPSO voting process includes three steps (see Fig. 1 (c)): the probability map initializing and updating, MSPSO sampling, and scene PPFs voting. We first initialize the probability map for scene sampling. Then we search the optimal particle position based on the probability map for each subpopulation and select the scene point where the particle is located as a reference point for Hough voting. In the voting process, those reference points generate hypothetical poses by Hough voting [36] and update the probability map with the maximum number of votes to motivate or punish the search path of MSPSO.

We next introduce the probability map and MSPSO sampling in detail.

Probability map. The probability map P_{map} assigns a value ρ_i to each point in the down-sampled point cloud S' , where ρ_i indicates the probability that this point is located on one of the real instances. Initially, each value ρ_i of the probability map is set to a constant. After a selected reference scene point s_r generates votes, the probability of s_r is calculated as following:

$$\rho_{s_r} = \frac{V_{s_r} - V_{min}}{V_{max} - V_{min}}, \quad (2)$$

where V_{s_r} is the maximum number of votes cast by s_r , V_{max} and V_{min} are the number of votes when the probability is 1 and 0 respectively. Afterwards, the probability ρ_i for other point in S' is updated according to ρ_{s_r} :

$$\rho_i = \frac{\sum_{j=1}^{N_s} \omega_{j,s_r} * \rho_{j,s_r}}{\sum_{j=1}^{N_s} \omega_{j,s_r}}, \quad (3)$$

where N_s is the number of scene reference points that have been voted. ρ_{j,s_r} is the probability obtained by voting with j -th scene reference point. ω_{j,s_r} is determined by a Gaussian kernel function, $\omega_{j,s_r} = \exp(-\frac{\|p_{j,s_r} - p_i\|^2}{2 * \sigma^2})$, where σ is the width parameter that is set to $\sigma = 0.25d_{obj}$ (d_{obj} is the diagonal length of the object's bounding box).

However, updating ρ_i for all points is time consuming. To speed it up, we only update the neighboring points around s_r , i.e., the points whose distance to s_r is less than $r_{max} * d_{obj}$, and r_{max} is the update radius coefficient. Several examples of the resulting probability map are shown in Fig. 2.

MSPSO sampling. Different from *Drost-PPFM*, we present a novel MSPSO sampling algorithm to increase the sampling ratio of valid scene reference points and reduce the number of sampling points. We solve this problem by iteratively optimizing the position of the scene reference point in the search space with regard to the above probability map.

Algorithm 1 MSPSO Voting with probability map P_{map}

Input: Model PPF hash table, downsampled scene cloud \mathcal{S}' , MSPSO iteration numbers t_{max} , subpopulation size N , and scene reference point sampling rate τ .

Output: a set of hypothetical poses

```

1: initialize  $P_{map}$ ;
2:  $n_s \leftarrow |\mathcal{S}'| * \tau$ ;  $G \leftarrow \text{int}(n_s / (t_{max} * N)) + 1$ ;
3: for each subpopulation  $s \in \{s_i | i = [1, G]\}$  do
4:    $s \leftarrow$  rejection sampling of  $P_{map}$  on  $\mathcal{S}'$ ;
5:   Initialize  $pbest_{s_r}$  and  $gbest_s$ ;
6:   for  $t \leftarrow 1$  to  $t_{max}$  do
7:     for each particle  $s_r \in s$  do
8:        $s_r$  votes and produces hypothetical poses;
9:       update  $pbest_{s_r}$  and  $gbest_s$ ;
10:      update  $P_{map}$  with Eq. (2) and Eq. (3);
11:    end for
12:    for each particle  $s_r \in s$  do
13:      for  $k \leftarrow 1$  to  $k_u$  do
14:        calculate velocity  $\mathbf{v}_{s_r}$  with Eq. (4);
15:        calculate position  $\mathbf{x}_{s_r}$  with Eq. (5);
16:         $\mathbf{s}_r \leftarrow$  the nearest point of  $\mathbf{x}_{s_r}$ ;
17:         $r \leftarrow \text{rand}(0, 1)$ ;
18:        if  $r < \rho_{s_r}$  then
19:           $\mathbf{s}_r \leftarrow \mathbf{s}_r$ ; break;
20:        end if
21:      end for
22:      if  $k > k_u$  then
23:         $\mathbf{s}_r \leftarrow$  rejection sampling of  $P_{map}$  on  $\mathcal{S}'$ ;
24:      end if
25:    end for
26:  end for
27: end for

```

Here we choose the MSPSO algorithm instead of the traditional particle swarm optimization (PSO) [44], [45], which has better global search ability than PSO and is more suitable for searching the best scene reference point on the probability map. In MSPSO, a scene reference point is regarded as a particle, and each particle has a flying velocity, which is composed of three parts: the upper flying velocity, the optimal historical position of the particle and the subpopulation:

$$v_{s_r}^{d,t+1} = \omega * v_{s_r}^{d,t} + c_1 * r_1 * (pbest_{s_r}^d - x_{s_r}^{d,t}) + c_2 * r_2 * (gbest_s^d - x_{s_r}^{d,t}), \quad (4)$$

where $v_{s_r}^{d,t}$ and $x_{s_r}^{d,t}$ are respectively the velocity and position value of the d -th dimension of the particle s_r in the s subpopulation in the t -th iteration time. ω is the inertia factor between 0 and 1 (set to 0.7 by default), which controls the trade-off between global and local experience. $pbest_{s_r}^d$ is the d -th dimension personal best position of particle s_r in s subpopulation found, and $gbest_s^d$ is the d -th dimension best position found by any particle in s subpopulation; c_1 is the cognitive acceleration constant, c_2 is the social acceleration constant, usually $c_1=c_2=2$; r_1 and r_2 are two random numbers uniformly distributed in the range $[0,1]$. Therefore, the new position of the particle is updated by:

$$x_{s_r}^{d,t+1} = x_{s_r}^{d,t} + v_{s_r}^{d,t+1}. \quad (5)$$

MSPSO voting. The detailed MSPSO voting process is shown in Algo. 1. Specifically, given the downsampled scene cloud \mathcal{S}' , we first calculate the number of subpopulation G :

$$G = \text{int}(n_s / (t_{max} * N)) + 1, \quad (6)$$

where $n_s = |\mathcal{S}'| * \tau$, t_{max} is the number of iterations, N is the size of the subpopulation, and τ is scene reference point sampling rate (line 2). Then we sample and vote on each subpopulation in turn (lines 3-27). In each subpopulation s , we first initialize each particle s_r by the rejection sampling algorithm based on the probability map P_{map} (line 4). Then, we iteratively perform t_{max} times of optimization voting, which includes three steps: s_r as scene reference point exercising PPFs Hough voting to generate hypothetical poses (line 8), updating P_{map} (line 10), and optimizing positions of the particles (lines 12-25).

In addition, Eq. (5) cannot be directly applied to optimize particles' positions, because the search space of traditional MSPSO is continuous, while the point cloud \mathcal{S}' is a discrete point set in the 3D space. In our algorithm, we first calculate the updated position \mathbf{x}_{s_r} of the particle s_r with Eq. (4) and Eq. (5) (lines 14-15). Then we find the nearest scene point \mathbf{s}_r of \mathbf{x}_{s_r} (line 16). Finally, we judge whether \mathbf{s}_r is the new position of the particle s_r (lines 17-20). In general, these steps can be iterated until the particle s_r is updated. To improve the global optimization performance, we stop this iteration and sample s_r on \mathcal{S}' by the rejection sampling of P_{map} if the particle s_r is not updated after two iterations.

C. Pose generation and verification

The above MSPSO voting scheme produces a set of pose votes. Unfortunately, the candidate object poses contain a significant fraction of incorrect poses, thus the hypothetical pose with the highest vote is not necessarily the correct match. Next, we conduct a pose verification step to search for the optimal pose in the hypothetical poses set.

Pose clustering. The candidate object poses are redundant because the poses of the reference points on the same instance are similar. Grouping similar poses can improve the efficiency of pose verification. To measure the similarity of poses, we define a point triplet \mathbf{FP}_{obj}^c , which consists of the center point $\mathbf{c}_{obj} = (c_x, c_y, c_z)$ of the object's bounding box and two auxiliary points $\mathbf{p}_{aux}^1 = (c_x - d_{obj}, c_y, c_z)$ and $\mathbf{p}_{aux}^2 = (c_x, c_y - d_{obj}, c_z)$. Specifically, the similarity of the poses are calculated as following:

$$e_{i,j} = D_{\text{Chebyshev}}((\mathbf{T}_i * \mathbf{FP}_{obj}^c), (\mathbf{T}_j * \mathbf{FP}_{obj}^c)). \quad (7)$$

In the clustering process, we first sort the hypothetical poses set \mathbf{T} in descending order with respect to the number of votes. Then we consider each pose in \mathbf{T} and judge whether they are similar to one of the subsequent poses. If it is true, we add the latter pose to the group of former poses, and the votes are accumulated. Finally, the poses in each cluster are averaged.

Hypotheses verification. After pose clustering, we obtain a set of new pose hypotheses with voting scores. Due to sensor noise and background clutter, the poses with the highest voting score may not be the correct pose, therefore further verification

is still required. Our verification indicator is the degree of overlap between the transformed model and the scene. Specifically, if one transformed model point \mathbf{p}'_m and its closest scene point \mathbf{p}_s meet the condition: $\|\mathbf{p}'_m - \mathbf{p}_s\|_2 < \varepsilon_D$ and $\angle(\mathbf{n}_{\mathbf{p}'_m}, \mathbf{n}_{\mathbf{p}_s}) < \varepsilon_A$, where ε_D is a distance threshold and ε_A is an angle threshold, we call these two points are overlapping. To recalculate the pose score, we first transform the object model into the scene with the pose, and find the nearest point of each model point in the scene and determine whether they are overlapping, and finally denote the pose verification score as the ratio of the number of overlapping points to the number of model points. Among all pose hypotheses, the pose with the highest verification score is the optimal matching pose.

D. Multi-instances, Multi-objects and Multi-frames Detection

To detect multiple instances in the scene, if only the top k cluster center poses are returned, this will easily lead to duplicate or missing poses. We apply a non-maximum suppression (NMS) algorithm to filter out duplicate poses. In our implementation, we return the pose with the highest score as the first instance pose. Then we visit other poses in descending order and calculate the Intersection over Union (IOU) of the visited pose and the already returned instance poses. If the IOU is smaller than a threshold (e.g., 0.4), we accept the pose as a new instance pose, otherwise, we discard it. We iterate this process until k instance poses are retrieved.

To solve the multiple object detection, we first vote to generate hypothetical poses and cluster them for all objects, then merge the hypothetical poses of all objects and sort them by overlapping area size. Finally, we filter out inaccurate poses based on the aforementioned NMS algorithm.

Furthermore, in most industrial applications (such as the bin-picking task), we usually need to capture the scene and detect objects in successive frames. Such multi-frame data is often overlooked in previous literature. Since the object information is similar between adjacent frames, using the prior knowledge in previous frames can significantly reduce the number of invalid scene reference points, thus speed up the matching process. To make full use of such information, we propose a strategy for initializing and updating the probability map, which is different from the single-frame case. In the initialization phase of the probability map, we first transform the model point cloud into the current scene through:

$$\mathcal{O}_c = \mathbf{T}_{p \rightarrow c} * \mathbf{T}_p * \mathcal{O}, \quad (8)$$

where $\mathbf{T}_{p \rightarrow c}$ is the rigid transformation from the previous frame to the current frame. \mathbf{T}_p is the instance pose estimated in the previous frame. Then, the probability $\rho_{i,m}$ of scene points that overlap with the transformed model \mathcal{O}_c are assigned $\rho_{i,m}^o$ (default as 0.9), and the points without overlap are assigned $\rho_{i,m}^0$ (default as 0.1). Furthermore, we set the scene reference point sampling rate τ to 1/80.

Finally, different from the lack of prior knowledge when we assign the initial value of the probability map in a single frame application, the prior knowledge obtained in the previous frame has high credibility and is helpful for scene reference point sampling. To exploit the prior knowledge, we merge the

initial value into the updating process of the probability map, which is formulated as:

$$\rho_i = \frac{\rho_{i,m} + \sum_{j=1}^{N_s} \omega_{j,s_r} * \rho_{j,s_r}}{1 + \sum_{j=1}^{N_s} \omega_{j,s_r}}. \quad (9)$$

E. Computational complexity

Our approach includes three stages: preprocessing, MSPSO voting, and pose verification. The preprocessing uses the voxel downsampling algorithm, and its time complexity is $O(n)$ where n is the number of points in the input point cloud. In the MSPSO voting phase, the most complicated calculation is the PPF Hough voting. The PPF target point sampling adopts the smart sampling strategy of Hinterstoisser et al. [9], and the PPF reference point is sampled by our MSPSO based on the probability map. The time complexity of this stage is $O(kn_s)$, where k is usually at least one magnitude smaller than $n_{s'}$ (please refer to [9] for details), $n_{s'}$ is the number of points in the down-sampled scene point clouds, and $n_s = n_{s'} * \tau$. Our experiments have suggested that the best accuracy can be achieved when τ is set to 1/10. Finally, the time complexity of the pose verification is $O(n_{T'} n_{\mathcal{O}'} n_{s'}^{\frac{2}{3}})$, where $n_{T'}$ is the number of poses after clustering hypothetical poses, and $n_{\mathcal{O}'}$ is the number of points in the object down-sampled point clouds. The operation at this stage is mainly to find the nearest neighbors by using a KD-tree. In general, the scale of the down-sampled point cloud is around 2000 or less, and the number of hypothetical poses is about 500 or less, so the verification can be done quite efficiently.

F. Discussion

Our work has two differences from the previous studies. First, we present a new voting module where the scene reference points are progressively sampled to better exploit the knowledge obtained in the previous voting. It can increase the inlier rate, retain more instance surface information, and improve computational efficiency and robustness of the algorithm. The sampling rate of reference points is reduced from 20% of Dorst et al. [36] to 10%, while we can still achieve better accuracy. Note that Hinterstoisser et al. [9] proposed a smart sampling algorithm for sampling the second point in each point pair and achieved satisfactory results. However, we have a different purpose, where our probability map-based MSPSO is to effectively sample the scene reference points, while the smart sampling algorithm aims to optimize the sampling for the PPF target points.

Second, previous 6D pose estimation methods based on pure point clouds have not yet considered the correlation between frames. In our method, we novelly transfer the information from the previous frame to the current frame with probability, which significantly reduces the number of invalid scene reference points and thus speeds up the matching process. Our experiments show that under the premise of ensuring the recognition rate, the sampling rate of the scene reference points can be reduced to 1.25% or even lower.

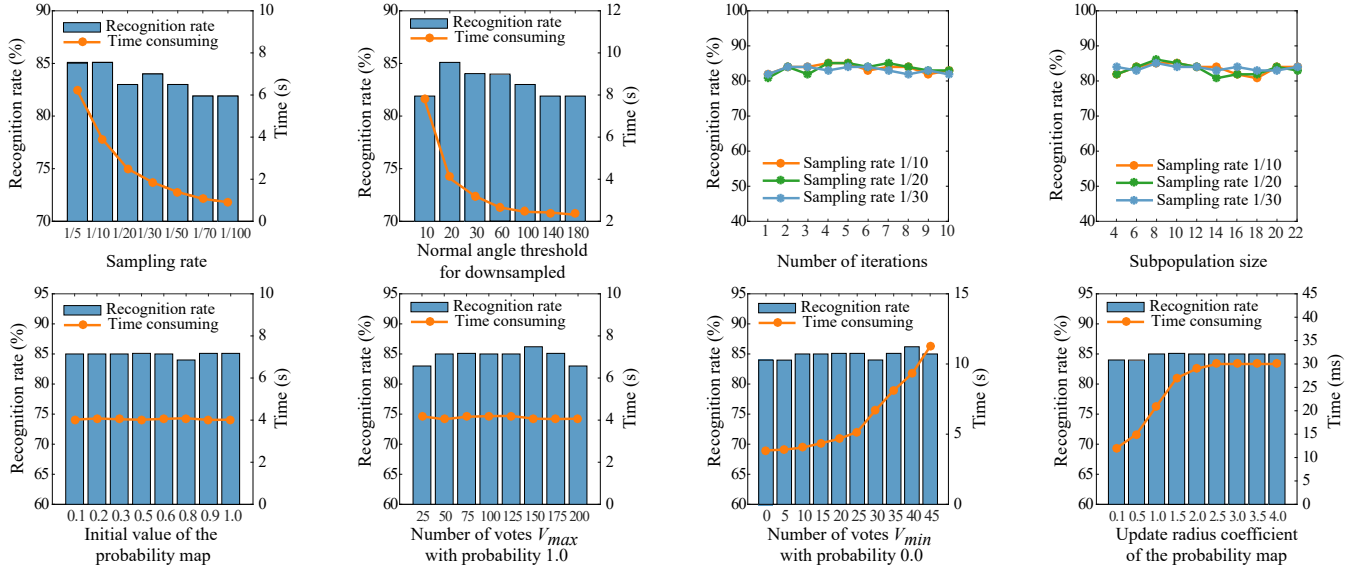


Fig. 3. Parameter analysis using a model from DTU [46]. To analyze each parameter, we set all other parameters as fixed values, which are the optimal values selected through experiments. If not specified, all parameters of our approach are set as default values: the sampling rates $\tau = 1/10$, the normal angle threshold for downsampling $\theta = 20^\circ$, the iteration numbers $t_{max} = 4$, the subpopulation size $N = 8$, the initial value of the probability map $\rho_i = 0.1$, the number of votes with probability 1.0 ($V_{max} = 50$), the number of votes with probability 0.0 ($V_{min} = 10$), and the update radius coefficient $r_{max} = 1.0$.

IV. EXPERIMENTAL RESULTS

In this section, we first evaluate the proposed algorithm and verify its effectiveness by conducting ablation studies. Then we demonstrate the performance qualitatively and quantitatively through visually inspecting our results and conducting a comparison with state-of-the-art approaches.

A. Experimental Setup

Datasets. We first carry out experiments on a real-scanning dataset, DTU [46], in which there are many cylindrical and flat 3D object models that are more challenging for pose estimation. To show our application on real industrial data, we also built a small-scale dataset, called RIPC, which consists of 9 objects and 100 scenes. This dataset also contains multiple instances and multiple objects. In addition, although our algorithm focuses on pose estimation of unorganized 3D point clouds, we also compare with some state-of-the-art deep learning-based methods on two RGBD datasets, including LineMOD (LM) [2] and LineMOD Occlusion (LM-O) [47].

Evaluation metric. Given an estimated 6D object pose $\hat{\mathbf{T}}$ and the ground-truth pose $\hat{\mathbf{T}}$, we use the average distance metric (ADM) to measure the L_2 distance between the model points transformed by $\hat{\mathbf{T}}$ and $\hat{\mathbf{T}}$. We define the pose error e_{ADM} as:

$$e_{ADM} = \max(\text{avg}_{\mathbf{x}_1 \in \mathcal{M}} \min_{\mathbf{x}_2 \in \mathcal{M}} \|\hat{\mathbf{T}}\mathbf{x}_1 - \hat{\mathbf{T}}\mathbf{x}_2\|_2, \|\hat{\mathbf{T}}\mathbf{c}_{obj} - \hat{\mathbf{T}}\mathbf{c}_{obj}\|_2) \quad (10)$$

where \mathcal{M} is the template model of the object, and we consider the distance between transformed object centers \mathbf{c}_{obj} .

We regard an estimated pose as positive if the pose error is less than a threshold ξ_e . Then the performance for detecting each object is quantitatively measured by using recognition rate (RR) that is the ratio of the number of true positive poses compared to the number of all its ground-truth poses. We also calculate the Mean Recall (MR) as the mean of the per-object

recognition rates to evaluate the overall performance on one dataset:

$$MR = \text{avg}_{o \in O} \frac{\sum_{s \in S} |P(o, s)|}{\sum_{s \in S} |G(o, s)|}, \quad (11)$$

where O and S are sets of templates and test scenes. $|P(o, s)|$ is the number of correctly detected poses and $|G(o, s)|$ is the number of ground-truth poses of object o in scene s .

Competitors. We select a commercial machine vision software MVTEC HALCON¹ as a competitor because it contains the optimized and significantly improved implementation of the original Drost-PPFM. We denote it Drost-PPFM*. We also compare to an open-source method Buch-17 [34], which presents a new pose voting and clustering method by integrating a local feature-based recognition pipeline. Here we test several representative 3D local feature descriptors, including PPF [35], spin images (SI) [33], FPFH [48] and SHOT [49].

B. Evaluation

Parameter settings. We conduct experiments with different parameter settings to evaluate their influence on recognition accuracy. In this section, we mainly analyze the eight parameters of the MSPSO sampling algorithm. The results of detailed parameter analysis are illustrated in Fig. 3, where we vary one parameter while all other parameters are fixed. We also show the running times with different sampling rates τ . From the figure, we can observe that the parameter τ affects the performance of our method, where the recognition rate and running efficiency are reduced when the number of sampled scene reference points decreases. We find that when $\tau = 1/10$, the recognition accuracy and running time can be well balanced. Next, we can observe that the recognition rate has reached the best value when $r_{max} = 1.0$, and

¹<https://www.mvtec.com/products/halcon/>

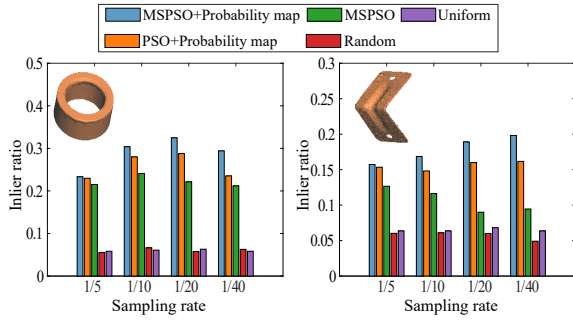


Fig. 4. Ablation studies of five sampling configurations.

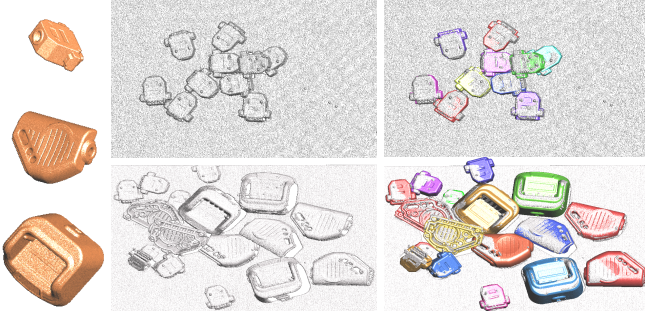


Fig. 5. Given one or more objects and a 3D point cloud scene, our method detects all instances in the scene and estimates the 6D pose for each instance.

the running time increases with the increase of the update radius coefficient in the early stage and remains stable in the later stage. Similarly, parameters θ and V_{min} affect our performance, and we can achieve the best performance when $\theta = 20^\circ$ and $V_{min} = 10$. In addition, our method is robust to t_{max} , N , ρ_i , and V_{max} .

Ablation studies. We now investigate the influence of different components of our method. We test five configurations:

- *Random sampling.* We randomly sample $|\mathcal{S}'| * \tau$ scene reference points in the down-sampled scene cloud \mathcal{S}' .
- *Uniform sampling.* We uniformly sample the scene reference points, i.e., every i -th ($i = 1/\tau$ in default) scene point is selected as a reference scene point.
- *MSPSO.* The scene reference points are progressively sampled with MSPSO but without the guidance of the probability map.
- *PSO+Probability map.* We replace the MSPSO with PSO.
- *MSPSO+Probability map.* This is our full method.

First, Fig. 4 shows the inlier ratio of these algorithms at different sampling rates. Here the inlier ratio represents the ratio of valid reference points on the instance's surface to the total scene reference points. We observe that using particle swarm optimization (MSPSO or PSO) increases the inlier ratio, while *MSPSO+Probability map* and *PSO+Probability map* guided by the probability map could further improve the performance, which is approximately 3 to 5 times that of random and uniform sampling. In comparison, our MSPSO is more effective than PSO.

Further, we analyze the recognition rate of these methods at different sampling rates. The experimental results are shown in Table I. We see that we obtain the best performance with *MSPSO+Probability map*, while the performance of conven-

tional random and uniform sampling is the worst.

Multi-instance and multi-object detection. Fig. 5 verifies the feasibility of our method for solving multi-instances and multi-objects detection and pose estimation problem. In Fig. 5, the input object point cloud is sampled from CAD models, while the input scene point cloud is captured by a structured-light scanner. If a scene consists of more than one object, we perform the steps of MSPSO sampling, PPF voting, pose clustering, and verification in parallel for all objects (this way we can also distinguish different objects). Then in the final step of returning appropriate poses, we put the pose hypotheses of all objects together and rank them according to the overlapping information of each pose in the scene. From Fig. 5 we see that our method detects all of the instances for each kind of object and returns their correct 6D poses.

Multi-frame Scenes. Finally, to demonstrate that we can speed up the matching process using multi-frame data, we conduct tests on the robotic-arm grasping system. For comparison, we also report the running time without multi-frame data as well as the time of Drost-PPFM*. The experimental results are shown in Table II. In a single-threaded environment, our method can obtain similar efficiency to Drost-PPFM*. Furthermore, we achieve the best performance through simple multi-core parallel processing using the OpenMP library.

C. Comparisons

Comparison on DTU dataset. Table III shows the quantitative comparison of the performance score and timing on DTU, where the recognition rate of each object and each method is reported. The DTU dataset is quite challenging because of its high occlusion and clutter. We further avoid testing the simple feature-rich objects and select flat, cylindrical, and thin-edge objects without many local distinguishable features. The results show that our algorithm outperforms other competitors for most of the test objects. In terms of running time, the commercial software HALCON (Drost-PPFM*) is the fastest because it takes advantage of the hardware and each step is also fully optimized. Comparing to Buch-17, our method is reasonably faster, but each step can still be further accelerated on GPU to meet the needs of industrial applications.

Comparison on real industrial scenes. Next, we conduct a comparison on our proposed RIPC dataset which is collected from real industrial scenes. Fig. 6 shows a qualitative comparison using several complex scenes for a robotic gripping task. The grasping objects involve common housings, polyhedrons, and threaded industrial devices. Similarly, we see that Drost-PPFM* and our approach achieve good performance on these scenes, but ours is still slightly better. Table IV reports the recognition rate on all of 9 objects in 100 scenes. This quantitative comparison verifies that our algorithm achieves the best-performing results on this dataset.

The surface geometric elements of industrial objects are simple. The local features of the point cloud are less discriminative. It is difficult to accurately estimate the pose for methods based on local features. The key to the success of Drost-PPFM* and our method is the excellent global description performance of PPF. Our method optimizes the

TABLE I
RESULTS OF ABLATION EXPERIMENTS USING TWO MODELS FROM DTU DATASET [46]. THE BEST RESULT OF EACH MEASUREMENT IS MARKED IN **BOLD**.

Methods	DTU-13-Tape							DTU-28-Angle_bar						
Sampling Rate	1/5	1/10	1/20	1/40	1/50	1/100	Mean	1/5	1/10	1/20	1/40	1/50	1/100	Mean
MSPSO + Probability map	85.1	85.1	83.0	85.1	83.0	81.9	83.9	58.0	55.7	53.4	52.3	51.1	35.2	51.0
PSO + Probability map	83.0	83.0	81.9	81.9	84.0	78.7	82.1	52.3	54.5	52.3	47.7	44.3	33.0	47.4
MSPSO	83.0	84.0	83.0	84.0	84.0	78.7	82.8	56.8	48.9	51.1	50.0	44.3	31.8	47.2
Random	84.0	81.9	81.9	79.7	78.7	78.7	80.8	50.0	48.9	48.9	46.6	43.2	30.7	44.7
Uniform	81.9	83.0	79.7	80.8	77.7	75.5	79.8	51.1	50.0	47.7	45.5	44.3.2	31.8	45.2

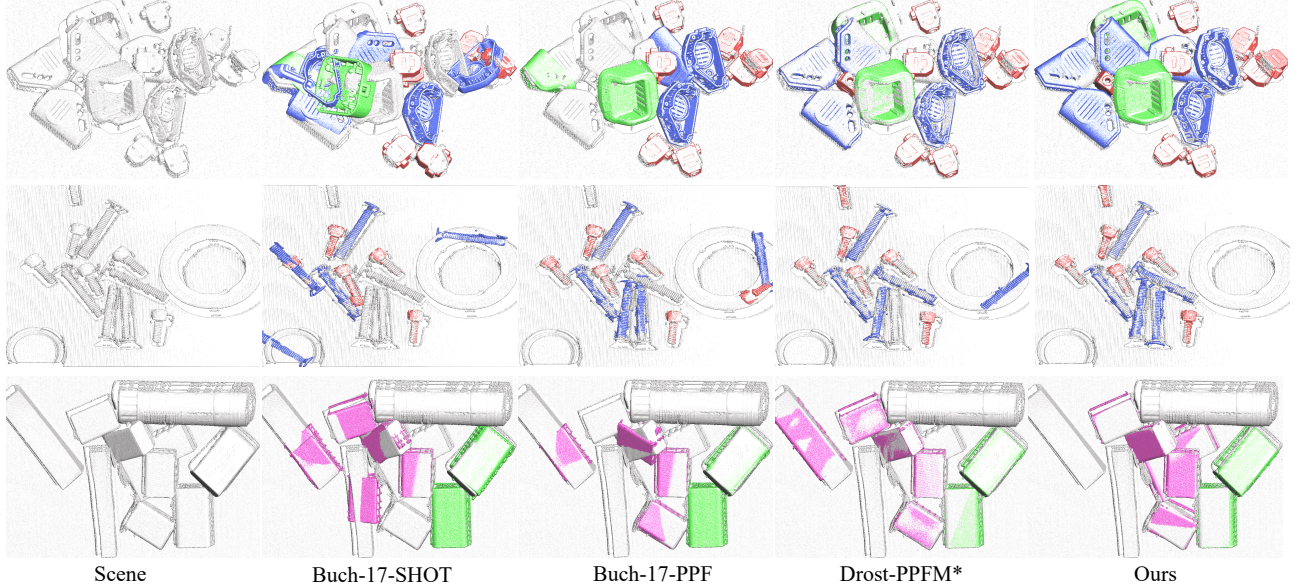






Fig. 6. Qualitative comparison using five scenes from our proposed RIPCD dataset.

TABLE II
AVERAGE RUNNING TIME(S) OF GRASPING 20 OBJECTS IN THE ROBOTIC ARM GRASPING SYSTEM. **MF** IS THE ABBREVIATION FOR OUR MULTI-FRAME POINT CLOUD DETECTION ALGORITHM.

Methods				
Drost-PPFM*	2.8	2.8	2.9	3.5
without MF	5.4	21.6	19.2	26.4
with MF	2.1	2.8	2.2	3.9
with MF and parallel	1.1	1.5	1.3	1.9

algorithm for sampled scene reference points and hypothetical poses verification, which is more robust than Drost-PPFM*.

Comparison on LM and LM-O. Even though our method is not specially designed for RGB-D data, we compare to a series of state-of-the-art deep learning-based methods, including depth-only, RGB-only, and RGB-D methods. Table V shows the comparison results. Our method, which does not use multi-frame information, obtains the best performance among the methods using only depth information. It outperforms CloudAEE+ICP [18] by 3.9% and 13.0% on LM and LM-O, respectively. Moreover, we achieve the same performance as PVN3D [28] on LM, which is one of state-of-the-art works using both RGB and depth information.

The objects in LM-O [47] are heavily occluded, making

pose estimation more difficult. As expected, we observe a significant performance drop for each method. However, our algorithm still outperforms other competitors.

V. CONCLUSION AND FUTURE WORK

We have presented an improved approach based on Drost-PPFM for 6D pose estimation of disordered objects in 3D point cloud scenes. We achieved major improvements in terms of efficiency and robustness by improving the sampling strategy for scene reference points. A new MSPSO algorithm based on the probability map can effectively improve the sampling of scene reference points. At the same time, the method of initializing the probability map with multi-frame data was also proposed to improve the efficiency of object pose estimation. We demonstrated our advantages by comparing it to the SOTA methods on several datasets.

Future work will be to explore more efficient object detection and pose estimation methods. A possible direction is to design an end-to-end deep neural network based on the middle-level structural analysis or high-level semantic information. For example, instead of using point correspondences, we could extract middle-level geometric primitives (such as plane, cylinder, sphere, cone) from objects and scenes to build feature correspondences. Besides, using object detection or segmentation maybe also improve the estimation performance.

TABLE III
QUANTITATIVE COMPARISON ON DTU DATASET. THE POSE ERROR THRESHOLD IS SET TO 0.1.

Methods	RR_2	RR_4	RR_{10}	RR_{12}	RR_{13}	RR_{14}	RR_{15}	RR_{16}	RR_{17}	RR_{20}	RR_{28}	RR_{31}	RR_{37}	RR_{41}	MR	Time(s)
Buch-17-PPF	67.1	71.2	62.4	59.4	11.7	48.7	73.3	14.7	85.0	15.3	19.3	76.8	62.4	29.6	55.7	1.71
Buch-17-SHOT	46.4	51.1	21.2	42.2	15.9	32.4	36.3	14.8	50.2	20.2	40.9	44.5	36.0	42.0	36.8	3.34
Buch-17-SI	64.3	72.3	46.5	48.4	3.2	39.9	62.3	21.3	75.4	18.0	45.5	61.9	61.4	43.2	51.4	1.49
Buch-17-FPFH	35.0	51.1	36.3	35.9	14.9	45.6	56.8	16.3	71.7	10.4	52.3	45.2	41.1	43.2	43.1	22.1
Drost-PPFM*	67.1	69.0	57.9	68.0	77.7	63.5	62.3	42.6	82.8	55.8	29.5	54.8	66.0	86.4	65.8	0.40
Ours	70.0	69.9	68.6	71.1	85.1	71.4	60.9	50.8	86.6	71.4	58.0	67.1	59.4	98.1	72.4	1.53

TABLE IV
QUANTITATIVE COMPARISON ON RIPCD DATASET. THE POSE ERROR THRESHOLD IS SET TO 0.1.

Methods	RR_{case1}	RR_{case2}	$RR_{connector}$	RR_{cube}	$RR_{eraser-big}$	$RR_{eraser-small}$	$RR_{nameplate}$	$RR_{screw-long}$	$RR_{screw-short}$	MR	Time(s)
Buch-17-PPF	24.2	56.1	87.4	66.8	35.0	24.5	100.0	20.0	76.1	65.2	5.36
Buch-17-SHOT	8.1	40.4	18.0	53.1	12.5	21.3	73.3	24.0	94.4	31.8	10.40
Buch-17-SI	0.0	31.6	62.6	66.4	57.5	12.8	100.0	20.0	93.0	52.4	6.65
Buch-17-FPFH	4.0	19.3	23.3	38.9	42.5	40.4	86.7	20.0	67.6	30.7	39.83
Drost-PPFM*	90.9	100.0	79.9	68.7	82.5	84.0	100.0	28.0	94.4	80.3	3.25
Ours	97.0	100.0	83.7	76.8	95.0	86.2	100.0	40.0	94.4	85.1	4.13

TABLE V
QUANTITATIVE EVALUATION OF 6D POSE ON ADD(S) [50] METRIC ON THE LINEMOD [2] AND LINEMOD OCCLUSION DATASET [47]. THE BEST RESULT OF EACH MEASUREMENT IS MARKED IN **BOLD**. THE NAMES OF SYMMETRIC OBJECTS ARE ALSO IN **BOLD**.

(a) LineMOD											(b) LineMOD Occlusion		
Modality	RGB			RGBD				D					
Method	PoseCNN+DeepIM [20], [21]	PVNet	CDPN [22]	SSD-6D+ICP [25]	PointFusion [26]	DF [27]	PVN3D [28]	CloudPose+ICP [17]	CloudAEE+ICP [18]	Ours	Modality	Method	average
ape	77.0	43.6	64.4	65.0	70.4	92.3	97.3	58.3	92.5	97.5	RGB	PoseCNN [20]	24.9
bvise	97.5	99.9	97.8	80.0	80.7	93.2	99.7	65.6	91.8	100		PVNet [10]	40.8
cam	93.5	86.9	91.7	78.0	60.8	94.4	99.6	43.0	88.9	99.7		RGBD	PoseCNN+ICP [20]
can	96.5	95.5	95.9	86.0	61.1	93.1	99.5	84.7	96.4	99.7	PVN3D [28]		63.2
cat	82.1	79.3	83.8	70.0	79.1	96.5	99.8	84.6	97.5	99.7	D		CloudPose+ICP [17]
driller	95.0	96.4	96.2	73.0	47.3	87.0	99.3	83.3	99.0	99.8		CloudAEE [18]	57.1
duck	77.7	52.6	66.8	66.0	63.0	92.3	98.2	43.2	92.7	98.0		CloudAEE+ICP [18]	66.1
e.bbox	97.1	99.2	99.7	100	99.9	99.8	99.8	99.5	99.8	100		Ours	79.1
glue	99.4	95.7	99.6	100	99.3	100	100	98.8	99.0	100			
holep	52.8	82.0	85.8	49.0	71.8	92.1	99.9	72.1	93.7	98.7			
iron	98.3	98.9	97.9	78.0	83.2	97	99.7	70.3	95.9	99.4			
lamp	97.5	99.3	97.9	73.0	62.3	95.3	99.8	93.2	96.6	99.3			
phone	87.7	92.4	90.8	79.0	78.8	92.8	99.5	81.0	97.4	99.8			
average	88.6	86.3	89.9	79.0	73.7	94.3	99.4	75.2	95.5	99.4			

In addition, we will further improve our algorithm to adapt it to structured RGB-D scenes.

REFERENCES

- [1] S. Hinterstoisser, C. Cagniat, S. Ilic, P. Sturm, N. Navab, P. Fua, and V. Lepetit, "Gradient response maps for real-time detection of textureless objects," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 5, pp. 876–888, 2011.
- [2] S. Hinterstoisser, S. Holzer, C. Cagniat, S. Ilic, K. Konolige, N. Navab, and V. Lepetit, "Multimodal templates for real-time detection of textureless objects in heavily cluttered scenes," in *IEEE International Conference on Computer Vision (ICCV)*, 2011, pp. 858–865.
- [3] T. Hodaň, X. Zabulis, M. Lourakis, Š. Obdržálek, and J. Matas, "Detection and fine 3d pose estimation of texture-less objects in rgb-d images," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2015, pp. 4421–4428.
- [4] A. Tejani, R. Kouskouridas, A. Doumanoglou, D. Tang, and T.-K. Kim, "Latent-class hough forests for 6 dof object pose estimation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 1, pp. 119–132, 2017.
- [5] Y. Konishi, K. Hattori, and M. Hashimoto, "Real-time 6d object pose estimation on cpu," in *2019 IEEE International Conference on Intelligent Robots and Systems (IROS)*, 2019, pp. 3451–3458.
- [6] Z. He, Z. Jiang, X. Zhao, S. Zhang, and C. Wu, "Sparse template-based 6-d pose estimation of metal parts using a monocular camera," *IEEE Trans. on Industrial Electronics*, vol. 67, no. 1, pp. 390–401, 2020.
- [7] S. Daftry, B. Ridge, W. Seto, T.-H. Pham, P. Ilhardt, G. Maggolino, M. Van der Merwe, A. Brinkman, J. Mayo, E. Kulczynski *et al.*, "Machine vision based sample-tube localization for mars sample return," in *2021 IEEE Aerospace Conference*, 2021, pp. 1–12.
- [8] T. Hodan, F. Michel, E. Brachmann, W. Kehl, A. GlentBuch, D. Kraft, B. Drost, J. Vidal, S. Ihrke, X. Zabulis *et al.*, "Bop: Benchmark for 6d object pose estimation," in *European Conference on Computer Vision (ECCV)*, 2018, pp. 19–34.
- [9] S. Hinterstoisser, V. Lepetit, N. Rajkumar, and K. Konolige, "Going further with point pair features," in *European Conference on Computer Vision (ECCV)*, 2016, pp. 834–848.
- [10] S. Peng, Y. Liu, Q. Huang, X. Zhou, and H. Bao, "Pvnet: Pixel-wise voting network for 6dof pose estimation," in *IEEE Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 4561–4570.
- [11] D. Hong, N. Yokoya, J. Chanussot, and X. X. Zhu, "An augmented linear mixing model to address spectral variability for hyperspectral unmixing," *IEEE Transactions on Image Processing*, vol. 28, no. 4, pp. 1923–1938, 2018.
- [12] D. Hong, L. Gao, J. Yao, B. Zhang, A. Plaza, and J. Chanussot, "Graph convolutional networks for hyperspectral image classification," *IEEE Transactions on Geoscience and Remote Sensing*, 2020.
- [13] D. Hong, N. Yokoya, J. Chanussot, J. Xu, and X. X. Zhu, "Joint and progressive subspace analysis (jpsa) with spatial-spectral manifold alignment for semisupervised hyperspectral dimensionality reduction," *IEEE Transactions on Cybernetics*, 2020.
- [14] D. Hong, L. Gao, N. Yokoya, J. Yao, J. Chanussot, Q. Du, and B. Zhang, "More diverse means better: Multimodal deep learning meets remote-sensing imagery classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 59, no. 5, pp. 4340–4354, 2020.
- [15] D. Hong, W. He, N. Yokoya, J. Yao, L. Gao, L. Zhang, J. Chanussot, and X. Zhu, "Interpretable hyperspectral artificial intelligence: When non-convex modeling meets hyperspectral remote sensing," *IEEE Geoscience and Remote Sensing Magazine*, vol. 9, no. 2, pp. 52–87, 2021.
- [16] F. Hagelskjær and A. G. Buch, "Pointvotenet: Accurate object detection and 6 dof pose estimation in point clouds," in *IEEE International Conference on Image Processing (ICIP)*. IEEE, 2020, pp. 2641–2645.
- [17] G. Gao, M. Lauri, Y. Wang, X. Hu, J. Zhang, and S. Frintrop, "6d object pose regression via supervised learning on point clouds," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*, 2020, pp. 3643–3649.
- [18] G. Gao, M. Lauri, X. Hu, J. Zhang, and S. Frintrop, "Cloudaae: Learning 6d object pose regression with on-line data synthesis on point clouds,"

- in *2020 IEEE International Conference on Robotics and Automation (ICRA)*, 2021.
- [19] Y. Shi, J. Huang, X. Xu, Y. Zhang, and K. Xu, "Stablepose: Learning 6d object poses from geometrically stable patches," in *IEEE Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 15 222–15 231.
- [20] Y. Xiang, T. Schmidt, V. Narayanan, and D. Fox, "Posecnn: A convolutional neural network for 6d object pose estimation in cluttered scenes," *arXiv preprint arXiv:1711.00199*, 2017.
- [21] Y. Li, G. Wang, X. Ji, Y. Xiang, and D. Fox, "Deepim: Deep iterative matching for 6d pose estimation," in *European Conference on Computer Vision (ECCV)*, 2018, pp. 683–698.
- [22] Z. Li, G. Wang, and X. Ji, "Cdpn: Coordinates-based disentangled pose network for real-time rgb-based 6-dof object pose estimation," in *IEEE International Conference on Computer Vision (ICCV)*, 2019, pp. 7678–7687.
- [23] C. Wu, L. Chen, Z. He, and J. Jiang, "Pseudo-siamese graph matching network for textureless objects' 6d pose estimation," *IEEE Trans. on Industrial Electronics*, pp. 1–1, 2021.
- [24] Z. Yang, X. Yu, and Y. Yang, "Dsc-poseNet: Learning 6dof object pose estimation via dual-scale consistency," in *IEEE Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 3907–3916.
- [25] W. Kehl, F. Manhardt, F. Tombari, S. Ilic, and N. Navab, "Ssd-6d: Making rgb-based 3d detection and 6d pose estimation great again," in *IEEE International Conference on Computer Vision (ICCV)*, Oct 2017, pp. 1521–1529.
- [26] D. Xu, D. Anguelov, and A. Jain, "Pointfusion: Deep sensor fusion for 3d bounding box estimation," in *IEEE Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 244–253.
- [27] C. Wang, D. Xu, Y. Zhu, R. Martín-Martín, C. Lu, L. Fei-Fei, and S. Savarese, "Densefusion: 6d object pose estimation by iterative dense fusion," in *IEEE Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 3343–3352.
- [28] Y. He, W. Sun, H. Huang, J. Liu, H. Fan, and J. Sun, "Pvn3d: A deep point-wise 3d keypoints voting network for 6dof pose estimation," in *IEEE Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 11 632–11 641.
- [29] Y. He, H. Huang, H. Fan, Q. Chen, and J. Sun, "Ffb6d: A full flow bidirectional fusion network for 6d pose estimation," in *IEEE Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 3003–3013.
- [30] W. Chen, X. Jia, H. J. Chang, J. Duan, L. Shen, and A. Leonardis, "Fs-net: Fast shape-based network for category-level 6d object pose estimation with decoupled rotation mechanism," in *IEEE Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 1581–1590.
- [31] C. Zhuang, Z. Wang, H. Zhao, and H. Ding, "Semantic part segmentation method based 3d object pose estimation with rgb-d images for bin-picking," *Robotics and Computer-Integrated Manufacturing*, vol. 68, p. 102086, 2021.
- [32] R. König and B. Drost, "A hybrid approach for 6dof pose estimation," in *European Conference on Computer Vision Workshops*. Springer, 2020, pp. 700–706.
- [33] A. E. Johnson and M. Hebert, "Using spin images for efficient object recognition in cluttered 3d scenes," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 21, no. 5, pp. 433–449, 1999.
- [34] A. G. Buch, L. Kiforenko, and D. Kraft, "Rotational subgroup voting and pose clustering for robust 3d object recognition," in *IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 4137–4145.
- [35] A. Buch and D. Kraft, "Local point pair feature histogram for accurate 3d matching," in *British Machine Vision Conference (BMVC)*, 2018.
- [36] B. Drost, M. Ulrich, N. Navab, and S. Ilic, "Model globally, match locally: Efficient and robust 3d object recognition," in *IEEE Computer Vision and Pattern Recognition (CVPR)*, 2010, pp. 998–1005.
- [37] C. Choi, Y. Taguchi, O. Tuzel, M. Liu, and S. Ramalingam, "Voting-based pose estimation for robotic assembly using a 3d sensor," in *IEEE International Conference on Robotics and Automation*, 2012, pp. 1724–1731.
- [38] B. Drost and S. Ilic, "3d object detection and localization using multi-modal point pair features," in *International Conference on 3D Imaging, Modeling, Processing, Visualization & Transmission*, 2012, pp. 9–16.
- [39] T. Birdal and S. Ilic, "Point pair features based object detection and pose estimation revisited," in *International Conference on 3D Vision (3DV)*, 2015, pp. 527–535.
- [40] J. Vidal, C.-Y. Lin, and R. Martí, "6d pose estimation using an improved method based on point pair features," in *International Conference on Control, Automation and Robotics (ICCAR)*, 2018, pp. 405–409.
- [41] J. Vidal, C.-Y. Lin, X. Lladó, and R. Martí, "A method for 6d pose estimation of free-form rigid objects using point pair features on range data," *Sensors*, vol. 18, no. 8, p. 2678, 2018.
- [42] I. Montri and S. Siriporn, "A multi-subpopulation particle swarm optimization: A hybrid intelligent computing for function optimization," in *International Conference on Natural Computation*, 2007, pp. 679–684.
- [43] L. Bi, W. Cao, W. Hu, and M. Wu, "Intelligent tuning of microwave cavity filters using granular multi-swarm particle swarm optimization," *IEEE Transactions on Industrial Electronics*, pp. 1–1, 2020.
- [44] R. Eberhart and J. Kennedy, "A new optimizer using particle swarm theory," in *International Symposium on Micro Machine and Human Science*, 1995, pp. 39–43.
- [45] J. Kennedy and R. Eberhart, "Particle swarm optimization," in *International Conference on Neural Networks*, 1995, pp. 1942–1948.
- [46] T. Sølund, A. G. Buch, N. Krüger, and H. Aanæs, "A large-scale 3d object recognition dataset," in *International Conference on 3D Vision (3DV)*, 2016, pp. 73–82.
- [47] E. Brachmann, A. Krull, F. Michel, S. Gumhold, J. Shotton, and C. Rother, "Learning 6d object pose estimation using 3d object coordinates," in *European Conference on Computer Vision (ECCV)*, 2014, pp. 536–551.
- [48] R. B. Rusu, N. Blodow, and M. Beetz, "Fast point feature histograms (fpfh) for 3d registration," in *IEEE International Conference on Robotics and Automation*, 2009, pp. 3212–3217.
- [49] S. Salti, F. Tombari, and L. Di Stefano, "Shot: Unique signatures of histograms for surface and texture description," *Computer Vision and Image Understanding*, vol. 125, pp. 251–264, 2014.
- [50] S. Hinterstoisser, V. Lepetit, S. Ilic, S. Holzer, G. Bradski, K. Konolige, and N. Navab, "Model based training, detection and pose estimation of texture-less 3d objects in heavily cluttered scenes," in *Asian conference on computer vision*, 2012, pp. 548–562.

Xuejun Xing received his B.S. degree from China University of Petroleum, Dongying, China, in 2010. He is currently working towards the M.S. degree in computer engineering at the School of Artificial Intelligence, University of Chinese Academy of Sciences. His research interests include 6D pose estimation, 3D reconstruction and machine vision.



Jianwei Guo is an associate professor in National Laboratory of Pattern Recognition (NLPR), Institute of Automation, Chinese Academy of Sciences (CASIA). He received his Ph.D. degree in computer science from CASIA in 2016, and bachelor degree from Shandong University in 2011. His research interests include 3D vision, computer graphics and image processing.



Liangliang Nan received his Ph.D. degree from the Graduate University of Chinese Academy of Sciences in 2009. He is currently an assistant professor at the Delft University of Technology (TU Delft). Prior to joining TU Delft, he was a research scientist at KAUST (2013-2018), and associate professor at SIAT (2009-2013). His research interests are in the fields of computer graphics, computer vision, and 3D Geoinformation.



Qingyi Gu received the M.E. and Ph.D. degree in Engineering, Hiroshima University, Japan, in 2010, and 2013 respectively. He is currently a professor in Institute of Automation, Chinese Academy of Sciences, China. His research interest is high-speed image processing, real-time 3D reconstruction, and applications in industry and biomedicine.



Xiaopeng Zhang received the PhD degree in computer science from Institute of Software, Chinese Academic of Sciences in 1999. He is a professor in National Laboratory of Pattern Recognition at Institute of Automation, Chinese Academy of Sciences. His main research interests include image processing, computer graphics and computer vision.





Dong-Ming Yan is a professor in National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences. He received his Ph.D. degree in computer science from Hong Kong University in 2010, and his master and bachelor degrees in computer science and technology from Tsinghua University in 2005 and 2002, respectively. His research interests include computer graphics, computer vision, and pattern recognition.