## CheckerPose: Progressive Dense Keypoint Localization for Object Pose Estimation with Graph Neural Network

Ruyi Lian Haibin Ling

Department of Computer Science, Stony Brook University, Stony Brook, NY 11794-2424, USA

{rulian,hling}@cs.stonybrook.edu

#### Abstract

Estimating the 6-DoF pose of a rigid object from a single RGB image is a crucial yet challenging task. Recent studies have shown the great potential of dense correspondence-based solutions, yet improvements are still needed to reach practical deployment. In this paper, we propose a novel pose estimation algorithm named CheckerPose, which improves on three main aspects. Firstly, CheckerPose densely samples 3D keypoints from the surface of the 3D object and finds their 2D correspondences progressively in the 2D image. Compared to previous solutions that conduct dense sampling in the image space, our strategy enables the correspondence searching in a 2D grid (i.e., pixel coordinate). Secondly, for our 3Dto-2D correspondence, we design a compact binary code representation for 2D image locations. This representation not only allows for progressive correspondence refinement but also converts the correspondence regression to a more efficient classification problem. Thirdly, we adopt a graph neural network to explicitly model the interactions among the sampled 3D keypoints, further boosting the reliability and accuracy of the correspondences. Together, these novel components make CheckerPose a strong pose estimation algorithm. When evaluated on the popular Linemod, Linemod-O, and YCB-V object pose estimation benchmarks, CheckerPose clearly boosts the accuracy of correspondence-based methods and achieves stateof-the-art performances. Code is available at https: //github.com/RuyiLian/CheckerPose.

### 1. Introduction

Object pose estimation from RGB images aims to estimate the rotation and translation of a given rigid object relative to the camera. It is crucial in various applications including robot grasping and manipulation [84, 70, 71], autonomous driving [43, 78, 34], augmented reality [44, 67], *etc.* Most existing methods [56, 68, 46, 22, 51, 81, 48, 38,



Figure 1: **Illustration of CheckerPose.** We evenly sample dense keypoints from the object surface, and predict the 2D locations in the input image. We design a binary code representation to progressively localize each keypoint in the iteratively refined 2D grids. To improve the localization, we also use graph neural networks to explicitly model the interactions between 3D keypoints. Note: we plot 8 keypoints for better visualization, while use 512 keypoints in practice.

61] first estimate an intermediate geometric representation, *i.e.*, the correspondences between 3D object keypoints and 2D image locations, and then recover the object pose using the Perspective-n-Point (PnP) algorithm. Theoretically, for a rigid object, four pairs of 3D-2D correspondences can determine a unique pose [55, 12, 52]. In practice, however, sparse correspondences easily degrade due to occlusion, background clutter, lighting variation, *etc*.

Increasing the number of 3D-2D correspondences is a feasible solution to enhance robustness, especially when combined with outlier removal mechanisms such as RANSAC. Recent methods [81, 48, 38, 18, 21, 75, 10] densely sample 2D image pixels and predict their 3D object coordinates. While these dense predictions improve the robustness of pose estimation, they have several drawbacks. Firstly, the predictions consider only visible pixels and ignore global relations between visible and occluded keypoints, making them unstable when the object is under severe occlusions. Secondly, estimating the corresponding 3D coordinates is nontrivial. Finally, the rich shape prior information is not effectively encoded.

To overcome the above issues, we propose a novel 6D pose estimation algorithm, named *CheckerPose*, which improves dense correspondence with three cooperative components: dense 3D sampling, progressive 2D localization through binary coding, and shape prior encoding with graph neural network, as illustrated in Figure 1.

For dense correspondence, CheckerPose samples 3D keypoints on the object surface and then finds their 2D pixel correspondences in the 3D-to-2D matching way. Compared to previous solutions that conduct dense sampling in the 2D image space, our strategy enables more efficient correspondence searching in a 2D grid (*i.e.*, pixel coordinate) using 2D binary coding, as well as explicit shape prior modeling with graph representation.

Then, to facilitate the localization of dense keypoints, we propose a 2D hierarchical binary coding to represent a 2D image position. Specifically, we superpose a grid on the input image and predict which cells contain the desired keypoints. The precision of the 2D keypoint location is controlled by the resolution of the grid. This novel representation allows us to refine the correspondence progressively. We first localize the keypoints in the  $2 \times 2$  grid, and then iteratively subdivide each cell and localize the keypoints in the refined grid. Inspired by ZebraPose [64], we use binary codes on the *x* and *y* directions to represent each cell, which makes the grids have a checkerboard pattern.

Furthermore, to capture the shape prior of the 3D object, we adopt a graph neural network to explicitly model the interactions among the sampled 3D keypoints and to guide the progressive correspondence estimation. In particular, we construct the k-nearest neighbor (k-NN) graph of the dense keypoints and utilize graph network layers to fuse information from a keypoint and its neighbors. By stacking multiple such layers, we can capture non-local interactions between invisible and visible keypoints, and thus significantly improve the prediction robustness of invisible keypoints.

To summarize, our main contributions are as follows:

- We propose to localize dense 3D keypoints in the input image, to establish dense correspondences for instance-level object pose estimation.
- We design a hierarchical binary coding strategy for 2D projections, which enables progressive localization of dense keypoints.
- We utilize graph neural networks to explicitly model the interactions between 3D keypoints and improve the predictions of invisible keypoints.

Together, these novel contributions make our CheckerPose a strong pose estimation algorithm. We conduct extensive experiments on the popular benchmarks including Linemod [17], Linemod-Occlusion [2], and YCB-V [79], and CheckerPose consistently achieves state-of-the-art performances.

#### 2. Related Work

In this section we review previous studies that are closely related to our work, mainly including different types of pose estimators and graph neural networks.

**Direct Methods.** Given an input RGB image, direct methods estimate the 6D pose of the object in the image without intermediate geometric representations, *e.g.*, 3D-2D correspondences. Traditional direct methods mainly adopt template matching techniques with hand-crafted features [25, 13, 16], and thus can not handle textureless objects well. Recent deep learning based methods utilize features learned by CNNs to directly regress 6D pose [79] or formulate the rotation estimation as a classification task by discretizing the rotation space SO(3) [72, 63, 28, 65].

**Correspondence Guided Methods.** Instead of direct estimation, correspondence guided methods [50, 56, 68, 46, 22, 51, 21, 23, 81, 48, 38, 75, 10, 64] follow a two-stage framework: they first predict a set of correspondences between 3D object frame coordinates and 2D image plane coordinates, and then recover the pose from the 3D-2D correspondences with a PnP algorithm [32, 30, 11, 73, 6]. RANSAC can be used to remove the outliers in the correspondences. Keypoint-localization based methods [50, 56, 68, 46, 22, 51, 21, 23] estimate the 2D coordinates for a sparse set of predefined 3D keypoints, while dense methods [81, 48, 38, 75, 10, 64] predict the 3D object frame coordinate of each 2D image pixel. Compared with sparse correspondences, dense correspondences contain richer context information of the scene and is more robust to occlusion.

Graph Neural Networks for 3D Vision. In 3D vision tasks, point clouds and meshes are important input data formats since they can efficiently represent complex shapes. Compared with convolutional neural networks (CNNs), graph neural networks (GNNs) [62] can handle inputs with irregular structures and effectively model the long-range dependencies, and thus are widely used for processing point clouds and meshes. While meshes can be naturally treated as graphs, a common practice of constructing graphs from point clouds is to treat each 3D point as graph nodes and connect each node to its k nearest neighbors [77, 60, 8]. GNN-based methods have been proposed for representation learning [60, 77, 40, 74], detection [59, 8], segmentation [53, 33], data generation [54, 39], camera pose inference [35, 36], etc. Graph techniques have also been used for learning dense correspondences between 3D shapes [58] using both local and global information. For object pose estimation, GNNs are mainly used for RGB-D inputs [9, 83] to enhance the feature extraction from different modalities. Another recent application is to learn geometric structures of the sparse keypoints for domain adaptation [82].



Figure 2: Framework of our progressive dense keypoint localization with graph neural network, *i.e.*, CheckerPose. Given an RGB image and object detection results, we progressively generate the binary codes representing the 2D locations of N 3D keypoints. (a) Initial graph embedding generation: we use a CNN backbone network to extract feature  $F_I^{(0)}$  from the zoomed-in RoI  $I_O$ , and then transform  $F_I^{(0)}$  to the initial keypoint embeddings  $F_G^{(0)}$  in the k-NN graph  $\mathcal{G}$ . (b) Progressive prediction: we use a graph neural network to generate the binary code representation in a coarse-to-fine manner. We adopt an additional CNN decoder network to generate image features with increased resolutions from  $F_I^{(0)}$ , and fuse the features in the graph neural network based on current predictions. Object segmentation masks M are predicted as an auxiliary learning task. (c) Feature fusion: to fuse the image feature  $F_I$  into the graph embeddings  $F_G$ , for each keypoint, we crop a feature patch from  $F_I$  based on the current localization result, and concatenate the flattened feature with keypoint embedding. We then use a shared MLP to fuse the concatenation and the result is the updated keypoint embedding.

**Our work** follows the two-stage framework and combines the strengths of both keypoint-based methods and dense methods, by localizing a dense set of predefined 3D keypoints to establish dense correspondences. Moreover, it utilizes GNNs to efficiently model the interactions among dense 3D keypoints and thus improve the localization in the input RGB image for monocular object pose estimation.

#### 3. Method

#### 3.1. Problem Formulation and Method Overview

Given an RGB image I and a rigid object O, our goal is to estimate rotation  $\mathbf{R} \in SO(3)$  and translation  $\mathbf{t} \in \mathbb{R}^3$  of O relative to the calibrated camera. We assume the 3D geometry information, *e.g.*, the 3D CAD model, is available, thus we can obtain  $N(N \gg 8)$  keypoints  $\mathcal{P} \subset \mathbb{R}^3$  from the object surface using farthest point sampling (FPS).

We adopt a two-stage pipeline for object pose estimation: we first predict 2D projection  $\rho \in \mathbb{R}^2$  for each keypoint  $P \in \mathcal{P}$ , and then regress the rotation and translation from the 3D-2D correspondences via a PnP solver. For the input RGB image, we use an off-the-shelf object detector [57, 69] to detect the object bounding box and extract the zoomedin Region of Interest (RoI)  $I_O$ , following the common prac-

tice in instance-level object pose estimation [38, 75, 10, 64]. Figure 2 illustrates our proposed pipeline. We first process the input RoI  $I_O$  by a backbone network to obtain backbone feature  $F_I^{(0)}$  and keypoint embedding  $F_G^{(0)}$  in the k-NN graph  $\mathcal{G}$ . Then we use graph network layers (*i.e.*, Edge-Conv [77]) to progressively localize the keypoints, which are represented as binary codes  $\mathbf{b_v}, \mathbf{b_x}$ , and  $\mathbf{b_y}$ . We also use a standard CNN decoder to transform  $F_I^{(0)}$  to a series of image feature maps, and fuse the features in the graph neural network based on the current predicted locations. The CNN decoder also outputs object segmentation masks Mas an auxiliary learning task. Finally, we convert the binary codes to 2D coordinates and use a PnP solver to recover the pose from the established correspondences. We describe our method, named CheckerPose due to the checkerboardlike binary pattern, in details as follows.

#### 3.2. Hierarchical Representation of 2D Keypoints

Establishing 3D-2D correspondences provides an intermediate representation for object pose estimation. In this work, we focus on localizing a dense set of predefined 3D keypoints  $\mathcal{P}$  in the 2D image plane. For  $N(N \gg 8)$  3D keypoints  $\mathcal{P}$ , we first predict whether their 2D projections



Figure 3: **Keypoint location representation.** (a) We represent the 2D projection coordinate as the center of the cell containing the 2D projection. (b) We iteratively refine the grid and represent the cell as binary codes  $\mathbf{b_x}$ ,  $\mathbf{b_y}$ .

appear in the RoI  $I_O$ , and then localize the keypoints inside  $I_O$ , denoted as  $\mathcal{P}_I$ . In contrast to directly regressing the precise coordinates, we superpose a  $2^d \times 2^d$  grid S on the RoI  $I_O$  and predict which cell  $s \in S$  contains the 2D projection  $\rho$  (Figure 3 (a)). Then we can use the coordinate of the cell center to approximate  $\rho$ , and only need to predict the discrete index  $(i_x, i_y)(0 \le i_x, i_y \le 2^d - 1)$  of the cell s, which is much easier than precise regression. The localization precision is controlled by the resolution of the grid S, and approaches the actual 2D projection as  $d \to \infty$ .

Based on the approximate representation, we can further localize the keypoint  $P \in \mathcal{P}_I$  in a coarse-to-fine manner. As shown in Figure 3 (b), at the beginning, we superpose a  $2 \times 2$  grid  $S^{(1)}$  on the RoI  $I_O$  and predict the index of the cell  $s_P^{(1)}$ . Then at iteration  $j \ (2 \le j \le d)$ , we increase the grid resolution from  $2^{j-1} \times 2^{j-1}$  to  $2^j \times 2^j$  by evenly splitting each cell  $s^{(j-1)} \in S^{(j-1)}$  into halves on both x and y directions. With the prediction of  $s_P^{(j-1)}$  in iteration j-1, we only need to search the corresponding  $2 \times 2$  sub-cells to find  $s_P^{(j)}$  in the refined grid  $S^{(j)}$ .

Inspired by ZebraPose [64], we use binary codes to concisely represent the hierarchical localization. For the cell  $s_P$  in the final  $2^d \times 2^d$  grid S, we use a d-bit binary code  $\mathbf{b_x}$  to represent the index  $i_x$  as

$$i_x = \sum_{k=1}^{a} b_x(k) \times 2^{d-k},$$
 (1)

where  $b_x(k)$  is the k-th bit of  $\mathbf{b}_{\mathbf{x}}$ . We use another d-bit binary code  $\mathbf{b}_{\mathbf{y}}$  to represent the index  $i_y$  in the same way. The first  $j(1 \le j \le d)$  bits of  $\mathbf{b}_{\mathbf{x}}$  and  $\mathbf{b}_{\mathbf{y}}$  also represent the cell  $s_P^{(j)} \in S^{(j)}$ . We use an additional 1-bit binary code  $\mathbf{b}_{\mathbf{v}}$ to indicate the existence of the projection  $\boldsymbol{\rho}$  in the RoI  $I_O$ , where  $\mathbf{b_v} = 1$  means  $\rho \in I_O$  while  $\mathbf{b_v} = 0$  means  $\rho \notin I_O$ . Compared with dense representations (*e.g.*, heatmaps [50, 46] and vector-fields [51, 22]), our representation needs only 2d + 1 binary bits for each keypoint, thus greatly reduces the memory usage for dense keypoint localization. In addition, during inference, we can efficiently convert the binary codes to the 2D coordinates. Furthermore, our representation can be naturally predicted in a progressive way, which allows to gradually improve the localization via iterative refinements.

#### 3.3. Dense Keypoint Localization via Graph Neural Network

Modeling the interactions among the keypoints  $\mathcal{P}$  is crucial for predicting their 2D locations. For the keypoints that are invisible due to occlusions or self-occlusions, the features of the visible ones provide additional clues to infer the 2D locations. However, previous keypoint-based methods mainly use convolutional neural networks (CNNs), which can not handle inputs with irregular structure and thus fail to explicitly capture the interactions among  $\mathcal{P}$ .

We instead utilize graph neural networks (GNNs) to process the features  $F = \{f_1, \dots, f_N\}$  of N keypoints  $\mathcal{P}$ . To construct a graph  $\mathcal{G}$  from  $\mathcal{P}$ , we treat each keypoint  $P_i \in \mathcal{P}(1 \leq i \leq N)$  as a graph node, and connect  $P_i$ to its k nearest neighbors in 3D Euclidean space to generate edges  $\mathcal{E}$ . We adopt the EdgeConv operation [77] as our graph network layer, which directly models local interactions between  $P_i$  and its neighbors. For edge  $(i, j) \in \mathcal{E}$ , we compute the feature  $e_{ij}$  as

$$e_{ijm} = \text{ReLU}(\theta_m \cdot (f_j - f_i) + \phi_m \cdot f_i), \qquad (2)$$

where  $e_{ijm}$  is the *m*-th channel of  $e_{ij}$ , and  $\theta_m$ ,  $\phi_m$  are the weights of the filters. The feature of  $P_i$  is updated by aggregating the edge features as

$$f'_{im} = \max_{j:(i,j)\in\mathcal{E}} e_{ijm},\tag{3}$$

where  $f'_{im}$  is the *m*-th channel of updated feature  $f'_i$ . By stacking multiple EdgeConv operations, our network can gradually learn the non-local interactions in a computationally efficient way for dense keypoints  $\mathcal{P}$ .

As shown in Figure 2 (a), to obtain the initial keypoint embeddings  $F_G^{(0)}$  in  $\mathcal{G}$ , we first use a backbone network to extract a  $C_0 \times 2^{d_0} \times 2^{d_0}$  feature map  $F_I^{(0)}$  from RoI  $I_O$ , where  $C_0$  is the number of the feature channels, and  $2^{d_0} \times 2^{d_0}$  is the spatial size. We then reshape  $F_I^{(0)}$  to  $C_0 \times 2^{2d_0}$  by flattening the spatial dimensions, and use a 1D convolutional network layer to obtain a  $N \times 2^{2d_0}$  feature map, which is regarded as the initial  $2^{2d_0}$ -dimensional embeddings  $F_G^{(0)}$  for N keypoints. After obtaining  $F_G^{(0)}$ , we use a graph neural network to predict the 1-bit indicator code  $\mathbf{b_v}$ , and progressively generate the *d*-bit index codes  $\mathbf{b_x}$ ,  $\mathbf{b_y}$ . Specifically, at stage 0, we apply  $L_0$  EdgeConv [77] operations to  $F_G^{(0)}$  to get the updated embeddings  $F_G^{(1)}$ , and then use shared MLPs to generate  $\mathbf{b_v}$  and the first  $d_0$  bits of  $\mathbf{b_x}$ ,  $\mathbf{b_y}$ , respectively. Then at stage  $j(1 \le j \le d - d_0)$ , we apply  $L_j$  EdgeConv operations to  $F_G^{(j)}$  to obtain  $F_G^{(j+1)}$ , and use shared MLPs to generate new bits  $b_x(d_0 + j)$ ,  $b_y(d_0 + j)$  for  $\mathbf{b_x}$ ,  $\mathbf{b_y}$ , respectively. We regard stage  $j(1 \le j \le d - d_0)$  as refinement stage, since it refines the localization from the low-resolution grid  $S^{(d_0+j-1)}$  to the high-resolution one  $S^{(d_0+j)}$ .

Compared with generating all bits at the network output layer, our progressive prediction enables image feature fusion at each refinement stage. As shown in Figure 2 (b), starting with the image feature map  $F_I^{(0)}$  with low spatial resolution  $2^{d_0} \times 2^{\overline{d}_0}$ , we use an additional CNN-based decoder to progressively generate image feature maps  $F_I^{(1)}, \cdots, F_I^{(d-d_0)}$  with increased spatial resolutions  $2^{d_0+1} \times 2^{d_0+1}, \cdots, 2^d \times 2^d$ , respectively. We also add skip connections between the backbone and the decoder to recover the high-resolution details lost in  $F_I^{(0)}$ . As shown in Figure 2 (c), at the beginning of the refinement stage j, for each keypoint P, we select local image feature from  $F_{I}^{(j)}$ based on the localization result in the previous stage. We then concatenate  $F_l^{(j)}$  with the keypoint embedding in the graph  $\mathcal{G}$ , and use a shared MLP to fuse the concatenation. The fused feature is used as the updated keypoint embedding. Since the initial keypoint embeddings  $F_G^{(0)}$  are obtained from  $F_I^{(0)}$ , fusing the local image features in the refinement stages provides critical high-resolution details for fine-grained localization.

#### 3.4. Training

For the 1-bit indicator code  $\mathbf{b}_{\mathbf{v}}$  of keypoint  $P \in \mathcal{P}$ , our network output  $\hat{\mathbf{b}}_{\mathbf{v}}$  is the probability that  $\mathbf{b}_{\mathbf{v}} = 1$ . We use binary cross-entropy loss for  $\mathbf{b}_{\mathbf{v}}$  as below:

$$\mathcal{L}_{v} = \frac{1}{N} \sum_{P \in \mathcal{P}} \mathbf{b}_{\mathbf{v}} \log \hat{\mathbf{b}}_{\mathbf{v}} + (1 - \mathbf{b}_{\mathbf{v}}) \log(1 - \hat{\mathbf{b}}_{\mathbf{v}}), \quad (4)$$

where N is the number of the keypoints. For d-bit index codes  $\mathbf{b_x}, \mathbf{b_y}$ , since we only localize the keypoints inside the RoI (*i.e.*,  $\mathbf{b_v} = 1$ ), denoted as  $\mathcal{P}_I$ , we compute binary cross-entropy loss for each bit of  $\mathbf{b_x}$  as

$$\mathcal{L}_{x} = \frac{1}{dN_{I}} \sum_{P \in \mathcal{P}_{I}} \sum_{k=1}^{d} b_{x}(k) \log(\hat{b}_{x}(k)) + (1 - b_{x}(k)) \log(1 - \hat{b}_{x}(k)), \quad (5)$$

where  $N_I$  is the number of keypoints inside the RoI,  $\hat{b}_x(k)$  is the network prediction for k-th bit of  $\mathbf{b}_{\mathbf{x}}$ . We compute

the loss  $\mathcal{L}_y$  for  $\mathbf{b}_y$  in the same way as  $\mathcal{L}_x$ .

Besides predicting the 2D projections as binary codes, we also enforce the network to output object segmentation masks. To do this, we apply a single CNN layer to the final image feature map  $F_I^{(d-d_0)}$  and obtain a  $2 \times 2^d \times 2^d$  output, which serves as the full segmentation mask  $M_{\text{full}}$  and the visible one  $M_{\text{vis}}$ . We input the network predictions to the sigmoid function and apply  $L_1$  loss as the mask loss  $\mathcal{L}_{\text{mask}}$ . Generating the masks can be regarded as an auxiliary task to facilitate the learning of image features.

The overall loss function  $\mathcal{L}$  is a combination of  $\mathcal{L}_v$ ,  $\mathcal{L}_x$ ,  $\mathcal{L}_y$ , and  $\mathcal{L}_{\text{mask}}$  as

$$\mathcal{L} = \mathcal{L}_v + \mathcal{L}_x + \mathcal{L}_y + \mathcal{L}_{\text{mask}}.$$
 (6)

Before training the whole network, we pretrain the layers that generate  $\mathbf{b}_{\mathbf{v}}$  and the first  $d_0$  bits of  $\mathbf{b}_{\mathbf{x}}$ ,  $\mathbf{b}_{\mathbf{y}}$ . This encourages the backbone network to quickly adapt to the object keypoints with smaller GPU memory usage, and makes the initial localization to be good for local image feature fusion in the refinement stages.

#### **3.5. Inference**

During inference, we first discard the keypoints with  $\mathbf{b_v} = 0$ . Then we convert the binary codes to the corresponding cells in the final grid *S* (Eq. 1), and use the 2D coordinates of the cell centers as the keypoint projections. In this way, we establish dense 3D-2D correspondences from the network outputs without time-consuming computation operations, *e.g.*, voting for the vector-field representations [51]. Finally we use the RANSAC/PnP [32] or Progressive-X [1] solvers to obtain the object pose from the dense 3D-2D correspondences.

We empirically find that for textureless objects with severe self-occlusions, discarding the correspondences outside  $M_{\rm vis}$  can improve the pose estimation results. To quantify the self-occlusions of a given object O, we uniformly sample 2,562 camera viewpoints on a sphere, and use the Hidden Point Removal (HPR) operator [27] to estimate the visibility of point  $P \in O$  from each viewpoint. We then calculate the proportion of the viewpoints for which P is visible, denoted as V(P). If  $0.2 \leq V(P) < 0.4$ , then Pis considered to be easily self-occluded. Note we ignore the points with V(P) < 0.2, to make our estimation robust to the classification error of the HPR operator. The overall self-occlusion of the object O can be computed by

$$r_{\rm so}(O) = \frac{1}{|O|} \sum_{P \in O} \mathbb{1}(0.2 \le V(P) < 0.4), \qquad (7)$$

where |O| is the number of vertices of the object CAD model, and  $\mathbb{1}(\cdot)$  is the indicator function. If  $r_{so}(O) \ge 0.5$ , *i.e.*, over half part of O is easily to be self-occluded, then we regard O as severely self-occluded.

#### 4. Experiments

#### 4.1. Experimental Setup

**Implementation Details.** Our method is implemented using PyTorch [49] and trained using the Adam optimizer [29] with a batch size of 32. We pretrain our network for 50, 000 steps with learning rate of 2e-4. We use N = 512 keypoints, and utilize k = 20 nearest neighbors to construct the k-NN graph  $\mathcal{G}$ . For the binary code representation, we set d = 6 and  $d_0 = 3$ . We resize the input RoIs to  $256 \times 256$ , and use HRNet [76] as our image feature backbone to extract  $1024 \times 8 \times 8$  feature map  $F_I^{(0)}$ . Then we apply  $L_0 = 2$  EdgeConv operations to get  $\mathbf{b_v}$  and the first  $d_0 = 3$  bits of  $\mathbf{b_x}$ ,  $\mathbf{b_y}$ , and obtain the full binary codes after 3 refinement stages with  $L_j = 3$  (j = 1, 2, 3) EdgeConv operations.

Datasets. We conduct our experiments on three commonly-used datasets for object pose estimation: Linemod (LM) [17], Linemod-Occlusion (LM-O) [2], and YCB-V [79]. LM consists of 13 sequences of real images with ground truth poses for a single object with background clutter and mild occlusion. Each sequence contains around 1, 200 images. Following [3], we utilize about 15% images for training while keeping the rest for testing. We additionally use 1,000 synthetic RGB images for each object during training following [38, 75, 10]. LM-O consists of 1,214 images from a sequence of LM [17], where ground truth poses of eight objects with partial occlusion are annotated for testing. YCB-V is composed of more than 110,000 real images of 21 objects with severe occlusion and clutter. Apart from the real training images, we also utilize the physically-based rendered data following [19] for training on LM-O and YCB-V.

**Evaluation Metrics.** We employ the common evaluation metric ADD(-S) for object pose estimation. ADD(-S) measures whether the average distance between the model points transformed by the predicted pose and the ground truth is less than 10% of the object's diameter (0.1d). For symmetric objects, ADD(-S) metric computes the deviation to the closest model point. On YCB-V, we also compute the AUC (area under curve) of ADD-S and ADD(-S) with a maximum threshold of 10 cm [79]. On LM, we also report the  $n^{\circ}$ , n cm metric, measuring the percentage of predicted 6D poses with rotation error below  $n^{\circ}$  and translation error below n cm. For symmetric objects  $n^{\circ}$ , n cm computes the smallest error for all possible ground truth poses [37, 75].

#### 4.2. Ablation Study on LINEMOD Dataset

We present ablation experiments on LM [17] in Table 1 to verify the effectiveness of each module. We also study the number of keypoint N and the size of neighborhood k in Supplementary. We train a single pose estimator for all objects for 120k steps, with a fixed learning rate of 1e-4 for

Method	A	ADD(-S)			5°5cm
	0.02d	0.05d	0.1d	2 2011	5 5011
GDR-Net [75]	35.5	76.3	93.7	62.1	N/A
SO-Pose [10]	45.9	83.1	96.0	76.9	98.5
EPro-PnP [7]	44.8	82.0	95.8	81.0	98.5
Ours (w/o GNN)	26.4	77.8	95.2	67.7	97.9
Ours (w/o Prog.)	14.1	56.9	85.8	42.3	94.1
Ours (w/o $M_{\rm full}$ )	30.2	82.8	96.7	79.3	98.9
Ours (w/o $M_{\rm vis}$ )	34.1	82.8	96.6	79.1	98.9
Ours (ResNet34)	31.3	80.2	95.6	74.2	98.6
Ours (RANSAC/PnP)	31.1	81.4	96.6	78.4	98.9
CheckerPose (Ours)	35.7	84.5	97.1	79.7	98.9

Table 1: Ablation Study on the LM Dataset.

the first 100k steps and a smaller learning rate of 5e-5 for the remaining steps. During inference, we utilize the detection results from Faster-RCNN [57] by [38]. We do not use any segmentation masks to filter the correspondences for fair comparison. Without specification, we use Progressive-X [1] to compute pose from the dense correspondences.

**Comparison with State of the Art.** As shown in Table 1, our method outperforms the state-of-the-art methods [75, 10, 7] w.r.t. ADD(-S) 0.05d, ADD(-S) 0.1d, and  $5^{\circ}5$ cm, and achieves comparable results w.r.t. ADD(-S) 0.02d and  $2^{\circ}2$ cm. The improvement of ADD(-S) 0.1d indicates that our method can facilitate the estimation of hard cases and serve as a good initialization for refinement methods [37, 26, 80]. Since the 2D coordinates of our estimated correspondences are approximated by the cell centers (Sec. 3.2), our pose estimation results in terms of ADD(-S) 0.02d may be further improved by increasing the grid resolution.

Effectiveness of Graph Neural Networks. Our network utilizes GNN layers, *e.g.*, EdgeConv [77], to explicitly model the interactions between different keypoints. We also report the result of removing all GNN layers in Table 1. Without GNN layers, the keypoints still interact indirectly via local image feature fusion modules, since the keypoints with close 2D locations share the similar local image features. However, the performance of pose estimation degrades significantly, demonstrating that it is important to directly model the keypoint interactions with GNN layers.

**Effectiveness of Progressive Prediction.** Progressively generating the binary codes enforces our network to gradually refine the localization in the iteratively subdivided grids. It also enables image feature fusion based on the intermediate estimations, which can provide crucial highresolution details for fine-grained localization. As shown in Table 1, the accuracy decreases significantly without

Method	PVNet [51]	S. Stage [21]	Hybrid [61]	RePose [26]	GDR-Net [75]	SO-Pose [10]	Zebra [64]	Ours
ape	15.8	19.2	20.9	31.1	46.8	48.4	57.9	58.3
can	63.3	65.1	75.3	80.0	90.8	85.8	95.0	95.7
cat	16.7	18.9	24.9	25.6	40.5	32.7	60.6	62.3
driller	65.7	69.0	70.2	73.1	82.6	77.4	94.8	93.7
duck	25.2	25.3	27.9	43.0	46.9	48.9	64.5	69.9
eggbox*	50.2	52.0	52.4	51.7	54.2	52.4	70.9	70.0
glue*	49.6	51.4	53.8	54.3	75.8	78.3	88.7	86.4
holep.	36.1	45.6	54.2	53.6	60.1	75.3	83.0	83.8
mean	40.8	43.3	47.5	51.6	62.2	62.3	76.9	77.5

Table 2: Comparison with State-of-the-art Methods on the LM-O Dataset. We report the Average Recall (%) of ADD(-S). (\*) denotes symmetric objects. We highlight the best result in red color, and the second best result in blue color.

progressively generating the binary codes, which clearly demonstrates the importance of progressive prediction.

Effectiveness of Object Segmentation Masks. Our network outputs the full segmentation mask  $M_{\rm full}$  and the visible one  $M_{\rm vis}$  as auxiliary tasks. As shown in Table 1, the performance degrades without either  $M_{\rm full}$  or  $M_{\rm vis}$ . The ADD(-S) 0.02d metric drops significantly without  $M_{\rm full}$ , indicating that predicting  $M_{\rm full}$  facilitates image feature extraction for keypoint localization, since all the keypoints should be located within  $M_{\rm full}$ . The degraded performance without  $M_{\rm vis}$  also implies that predicting  $M_{\rm vis}$  provides important context information including occlusions.

**Impact of Backbone Networks.** We report the results of our method with different backbone networks in Table 1. After replacing HRNet [76] by ResNet34 [15], our method still achieves comparable results with state of the art, which demonstrates the efficacy of our method regardless of the backbone networks.

**Influence of PnP Solvers.** We show the results with different PnP solvers during inference in Table 1. Since our correspondences are established from the binary codes, a small perturbation of our network prediction can result in flipped bit values, which may correspond to dramatically different locations in the input RoI. Compared with RANSAC/PnP [32], Progressive-X [1] contains a spatial coherence filter to efficiently remove such outliers, and thus achieves better performance w.r.t. to all the metrics, especially ADD(-S) 0.02d.

#### 4.3. Comparison to State of the Art

In this section we present the quantitative results of our method on LM-O and YCB-V datasets. We train a single CheckerPose for each object for 380,000 steps with a fixed learning rate of 1e-4. During inference, we utilize the detections from FCOS [69] provided by CDPNv2 [38].

Experiments on the LM-O dataset. We report the recall of ADD(-S) metric for the LM-O dataset in Table 2. Based on the criterion discussed in Sec. 3.5, we filter out the correspondences outside the visible segmentation masks  $M_{\rm vis}$  for textureless objects with severe self-occlusions, including can, cat, driller, and eggbox. Without the filtering operation, the average recall of ADD(-S) of our method is 77.1, which surpasses previous methods. The detailed results of each object without filtering are provided in supplementary material. The additional filtering operation further improves the performance of our method. The intuition is that it is infrequent to observe an easily self-occluded keypoint P in the training images. Besides, due to the lack of texture, it is also hard to infer the location of P from other keypoints with distinguishable features. Such objects may require much more training steps to achieve stable estimations for easily self-occluded keypoints. Simply discarding correspondences outside  $M_{\rm vis}$  reduces unstable localization results when our network is trained for limited steps, and enhances the robustness of pose estimation.

**Experiments on the YCB-Video dataset.** We report the averaged metrics of 21 objects in Table 3, and provide detailed results in the suppl.. Based on the criterion discussed in Sec. 3.5, we use visible segmentation masks to filter correspondences for foam\_brick. We also apply the filtering operation to pudding\_box because it is severely occluded by gelatin\_box, which is a distraction object with similar texture. As shown in Table 3, CheckerPose achieves the best performance w.r.t. ADD(-S) and AUC of ADD(-S), and is comparable with state of the art w.r.t. AUC of ADD-S.

#### 4.4. Qualitative Results

In Figure 4, we provide localization results of eight keypoints for the occluded and flipped bowl. While our network directly outputs the 2D locations, the results of other dense methods [64, 75] are computed by projecting the keypoints using the estimated poses. Figure 4 (a) visualizes the reprojections of ZebraPose [64], where the keypoints

Method	ADD(-S)	AUC ADD-S	AUC ADD(-S)
SegDriven [22]	39.0	_	_
SingleStage [21]	53.9	-	—
CosyPose [31]	_	89.8	84.5
RePose [26]	62.1	88.5	82.0
GDR-Net [75]	60.1	91.6	84.4
SO-Pose [10]	56.8	90.9	83.9
ZebraPose [64]	80.5	90.1	85.3
DProST [47]	65.1	_	77.4
CheckerPose (Ours)	81.4	91.3	86.4

Table 3: Comparison on the YCB-Video Dataset. We report the ADD(-S), and AUC of ADD-S and ADD(-S). Following [79], the symmetric metric is used for all objects in ADD-S while only for symmetric objects in ADD(-S). We highlight the best result in red color, and the second best result in blue color. "-" denotes unavailable results.





(c) CheckerPose (Ours)



Figure 4: Keypoint localization. (a) Keypoint locations based on the predicted pose of ZebraPose [64]. (b) Keypoint locations based on the pose estimated by GDR-Net [75]. (c) Keypoint locations output by our network. (d) The ground truth keypoint locations. Considering the symmetry of the bowl, we use the equivalent rotations closest to our prediction to project the keypoints in (a), (b), and (d).

concentrate on the visible pixels. Since ZebraPose generates pixel-wise 3D coordinates from the visible regions, it predicts a drastically wrong pose for the severely occluded bowl. As shown in Figure 4 (b), the reprojections of GDR-Net [75] cover the region similar to the ground truth (Figure 4 (d)). However, the order of the blue keypoint and the red one changes from clockwise to counterclockwise, indicating the bowl is actually faced up. Since GDR-Net is an end-to-end method, it may memorize poses that frequently appear in the training samples. As shown in Figure 4 (c), our network is capable of localizing the keypoints for the upside-down object with severe occlusion. More qualitative results can be found in the Supplementary Material.

#### 4.5. Runtime Analysis

We test the running speed on the LM-O dataset. Given a  $640 \times 480$  RGB image, we evaluate the speed on a desktop with an Intel 3.30GHz CPU and an NVIDIA GeForce GTX 1080 GPU (8G), which is reasonable in real-world application. The FCOS detector [69] takes 87 ms for each image. The runtime of establishing the dense 3D-2D correspondences by our network is 78 ms. RANSAC/PnP [32] takes only 1 ms to recover pose from the correspondences, while Progressive-X [1] takes 32 ms. Under the same testing environment, ZebraPose [64] requires 10ms for generating 3D-2D correspondences by CNN and around 350ms to estimate pose using Progressive-X. The overall running time of our method is greatly reduced, because we establish at most 512 candidate 3D-2D correspondences while ZebraPose outputs  $128^2$  candidates in the worst case.

#### 5. Conclusion

In this work, we propose a novel way to establish dense correspondences for object pose estimation, by progressively localizing dense 3D keypoints in the input image. With dense keypoints including occluded and self-occluded ones, we comprehensively explore the available geometry information and enhance the robustness of pose estimation under severe occlusion. We adopt graph neural networks to explicitly model the keypoint interactions, and design a hierarchical binary code representation for the 2D locations. The experiments on LM, LM-O and YCB-V datasets demonstrate that our method achieves state-of-the-art performance of instance-level object pose estimation.

Acknowledgement. We thank all reviewers for valuable comments and suggestions. The work was supported in part by US National Science Foundation Grants 2006665 and 2128350, and by the Defense Advanced Research Projects Agency (DARPA) under Agreement No. HR0011-22-9-0077. This work is also supported in part by the SBU/BNL Seed Grant Award. Any opinions, findings, and conclusions, or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of funding agencies.

#### References

- [1] Daniel Barath and Jiri Matas. Progressive-x: Efficient, anytime, multi-model fitting algorithm. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (*ICCV*), pages 3780–3788, 2019. 5, 6, 7, 8, 13
- [2] Eric Brachmann, Alexander Krull, Frank Michel, Stefan Gumhold, Jamie Shotton, and Carsten Rother. Learning 6d object pose estimation using 3d object coordinates. In *Proceedings of the European Conference on Computer Vision* (ECCV), pages 536–551. Springer, 2014. 2, 6, 13, 14, 15
- [3] Eric Brachmann, Frank Michel, Alexander Krull, Michael Ying Yang, Stefan Gumhold, et al. Uncertaintydriven 6d pose estimation of objects and scenes from a single rgb image. In *Proceedings of the IEEE Conference* on Computer Vision and Pattern Recognition (CVPR), pages 3364–3372, 2016. 6
- [4] Dingding Cai, Janne Heikkilä, and Esa Rahtu. SC6D: symmetry-agnostic and correspondence-free 6d object pose estimation. In 2022 International Conference on 3D Vision (3DV). IEEE, 2022. 15
- [5] Pedro Castro and Tae-Kyun Kim. CRT-6D: fast 6d object pose estimation with cascaded refinement transformers. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 5746–5755, 2023. 15
- [6] Bo Chen, Alvaro Parra, Jiewei Cao, Nan Li, and Tat-Jun Chin. End-to-end learnable geometric vision by backpropagating pnp optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (CVPR), pages 8100–8109, 2020. 2
- [7] Hansheng Chen, Pichao Wang, Fan Wang, Wei Tian, Lu Xiong, and Hao Li. EPro-PnP: generalized end-to-end probabilistic perspective-n-points for monocular object pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2781–2790, 2022. 6
- [8] Jintai Chen, Biwen Lei, Qingyu Song, Haochao Ying, Danny Z Chen, and Jian Wu. A hierarchical graph network for 3d object detection on point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 392–401, 2020. 2
- [9] Wei Chen, Xi Jia, Hyung Jin Chang, Jinming Duan, Linlin Shen, and Ales Leonardis. FS-Net: fast shape-based network for category-level 6d object pose estimation with decoupled rotation mechanism. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (CVPR), pages 1581–1590, 2021. 2
- [10] Yan Di, Fabian Manhardt, Gu Wang, Xiangyang Ji, Nassir Navab, and Federico Tombari. SO-Pose: exploiting selfocclusion for direct 6d pose estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (*ICCV*), pages 12396–12405, 2021. 1, 2, 3, 6, 7, 8
- [11] Luis Ferraz, Xavier Binefa, and Francesc Moreno-Noguer. Very fast solution to the pnp problem with algebraic outlier rejection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 501– 508, 2014. 2

- [12] Xiao-Shan Gao, Xiao-Rong Hou, Jianliang Tang, and Hang-Fei Cheng. Complete solution classification for the perspective-three-point problem. *IEEE transactions on pattern analysis and machine intelligence*, 25(8):930–943, 2003. 1
- [13] Chunhui Gu and Xiaofeng Ren. Discriminative mixture-oftemplates for viewpoint classification. In *Proceedings of European Conference on Computer Vision (ECCV)*, pages 408– 421. Springer, 2010. 2
- [14] Rasmus Laurvig Haugaard and Anders Glent Buch. Surfemb: Dense and continuous correspondence distributions for object pose estimation with learnt surface embeddings. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6749–6758, 2022. 15
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 770–778, 2016. 7
- [16] Stefan Hinterstoisser, Cedric Cagniart, Slobodan Ilic, Peter Sturm, Nassir Navab, Pascal Fua, and Vincent Lepetit. Gradient response maps for real-time detection of textureless objects. *IEEE transactions on pattern analysis and machine intelligence*, 34(5):876–888, 2011. 2
- [17] Stefan Hinterstoisser, Vincent Lepetit, Slobodan Ilic, Stefan Holzer, Gary Bradski, Kurt Konolige, and Nassir Navab. Model based training, detection and pose estimation of texture-less 3d objects in heavily cluttered scenes. In Asian Conference on Computer Vision (ACCV), pages 548–562. Springer, 2012. 2, 6, 13
- [18] Tomas Hodan, Daniel Barath, and Jiri Matas. EPOS: estimating 6d pose of objects with symmetries. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 11703–11712, 2020. 1
- [19] Tomáš Hodaň, Martin Sundermeyer, Bertram Drost, Yann Labbé, Eric Brachmann, Frank Michel, Carsten Rother, and Jiří Matas. BOP challenge 2020 on 6D object localization. *European Conference on Computer Vision Workshops (EC-CVW)*, 2020. 6, 15
- [20] Yinlin Hu, Pascal Fua, and Mathieu Salzmann. Perspective flow aggregation for data-limited 6d object pose estimation. In Proceedings of the European Conference on Computer Vision (ECCV), pages 89–106. Springer, 2022. 15
- [21] Yinlin Hu, Pascal Fua, Wei Wang, and Mathieu Salzmann. Single-stage 6d object pose estimation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 2930–2939, 2020. 1, 2, 7, 8, 16
- [22] Yinlin Hu, Joachim Hugonot, Pascal Fua, and Mathieu Salzmann. Segmentation-driven 6d object pose estimation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 3385–3394, 2019. 1, 2, 4, 8, 16
- [23] Yinlin Hu, Sebastien Speierer, Wenzel Jakob, Pascal Fua, and Mathieu Salzmann. Wide-depth-range 6d object pose estimation in space. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 15870–15879, 2021. 2

- [24] Lin Huang, Tomas Hodan, Lingni Ma, Linguang Zhang, Luan Tran, Christopher Twigg, Po-Chen Wu, Junsong Yuan, Cem Keskin, and Robert Wang. Neural correspondence field for object pose estimation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 585–603. Springer, 2022. 15
- [25] Daniel P Huttenlocher, Gregory A. Klanderman, and William J Rucklidge. Comparing images using the hausdorff distance. *IEEE Transactions on pattern analysis and machine intelligence*, 15(9):850–863, 1993. 2
- [26] Shun Iwase, Xingyu Liu, Rawal Khirodkar, Rio Yokota, and Kris M Kitani. RePOSE: fast 6d object pose refinement via deep texture rendering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3303–3312, 2021. 6, 7, 8, 16
- [27] Sagi Katz, Ayellet Tal, and Ronen Basri. Direct visibility of point sets. ACM Transactions On Graphics (TOG), 26(3):24, 2007. 5, 13
- [28] Wadim Kehl, Fabian Manhardt, Federico Tombari, Slobodan Ilic, and Nassir Navab. SSD-6D: making rgb-based 3d detection and 6d pose estimation great again. In *Proceedings* of the IEEE International Conference on Computer Vision (ICCV), pages 1521–1529, 2017. 2
- [29] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015. 6
- [30] Laurent Kneip, Hongdong Li, and Yongduek Seo. Upnp: An optimal o (n) solution to the absolute pose problem with universal applicability. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 127–142. Springer, 2014. 2
- [31] Yann Labbé, Justin Carpentier, Mathieu Aubry, and Josef Sivic. CosyPose: consistent multi-view multi-object 6d pose estimation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 574–591. Springer, 2020. 8, 17
- [32] Vincent Lepetit, Francesc Moreno-Noguer, and Pascal Fua.
  Epnp: An accurate o (n) solution to the pnp problem. *International journal of computer vision*, 81(2):155, 2009. 2, 5, 7, 8, 13
- [33] Guohao Li, Matthias Muller, Ali Thabet, and Bernard Ghanem. DeepGCNs: can GCNs go as deep as CNNs? In Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), pages 9267–9276, 2019. 2
- [34] Shichao Li, Zengqiang Yan, Hongyang Li, and Kwang-Ting Cheng. Exploring intermediate representation for monocular vehicle pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (CVPR), pages 1873–1883, 2021. 1
- [35] Xinyi Li and Haibin Ling. PoGO-Net: Pose graph optimization with graph neural networks. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5875–5885, 2021. 2
- [36] Xinyi Li and Haibin Ling. GTCaR: Graph transformer for camera re-localization. In *Proceedings of the European Conference on Computer Vision (ECCV)*, volume 13670 of *Lecture Notes in Computer Science*, pages 229–246. Springer, 2022. 2

- [37] Yi Li, Gu Wang, Xiangyang Ji, Yu Xiang, and Dieter Fox. DeepIM: deep iterative matching for 6d pose estimation. In Proceedings of the European Conference on Computer Vision (ECCV), pages 683–698, 2018. 6
- [38] Zhigang Li, Gu Wang, and Xiangyang Ji. CDPN: coordinates-based disentangled pose network for real-time rgb-based 6-dof object pose estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (*ICCV*), pages 7678–7687, 2019. 1, 2, 3, 6, 7
- [39] Kevin Lin, Lijuan Wang, and Zicheng Liu. Mesh graphormer. In Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), pages 12939– 12948, 2021. 2
- [40] Zhi-Hao Lin, Sheng-Yu Huang, and Yu-Chiang Frank Wang. Convolution in the cloud: Learning deformable kernels in 3d graph convolution networks for point cloud analysis. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 1800–1809, 2020. 2
- [41] Lahav Lipson, Zachary Teed, Ankit Goyal, and Jia Deng. Coupled iterative refinement for 6d multi-object pose estimation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 6728– 6737, 2022. 15
- [42] Xingyu Liu, Ruida Zhang, Chenyangguang Zhang, Bowen Fu, Jiwen Tang, Xiquan Liang, Jingyi Tang, Xiaotian Cheng, Yukang Zhang, Gu Wang, and Xiangyang Ji. GDRNPP. https://github.com/shanice-1/ gdrnpp\_bop2022, 2022. 15
- [43] Fabian Manhardt, Wadim Kehl, and Adrien Gaidon. ROI-10D: monocular lifting of 2d detection to 6d pose and metric shape. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 2069– 2078, 2019. 1
- [44] Eric Marchand, Hideaki Uchiyama, and Fabien Spindler. Pose estimation for augmented reality: a hands-on survey. *IEEE transactions on visualization and computer graphics*, 22(12):2633–2651, 2015.
- [45] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020. 16
- [46] Markus Oberweger, Mahdi Rad, and Vincent Lepetit. Making deep heatmaps robust to partial occlusions for 3d object pose estimation. In *Proceedings of the European Conference* on Computer Vision (ECCV), pages 119–134, 2018. 1, 2, 4
- [47] Jaewoo Park and Nam Ik Cho. DProST: 6-dof object pose estimation using space carving and dynamic projective spatial transformer. In *Proceedings of the European Conference* on Computer Vision (ECCV), 2022. 8, 16, 17
- [48] Kiru Park, Timothy Patten, and Markus Vincze. Pix2Pose: pixel-wise coordinate regression of objects for 6d pose estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 7668–7677, 2019. 1, 2
- [49] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming

Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32:8026– 8037, 2019. 6

- [50] Georgios Pavlakos, Xiaowei Zhou, Aaron Chan, Konstantinos G Derpanis, and Kostas Daniilidis. 6-DoF object pose from semantic keypoints. In 2017 IEEE international conference on robotics and automation (ICRA), pages 2011–2018. IEEE, 2017. 2, 4
- [51] Sida Peng, Yuan Liu, Qixing Huang, Xiaowei Zhou, and Hujun Bao. PVNet: pixel-wise voting network for 6dof pose estimation. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition (CVPR), pages 4561–4570, 2019. 1, 2, 4, 5, 7
- [52] Mikael Persson and Klas Nordberg. Lambda twist: An accurate fast robust perspective three point (p3p) solver. In *Proceedings of the European conference on computer vision* (ECCV), pages 318–332, 2018. 1
- [53] Xiaojuan Qi, Renjie Liao, Jiaya Jia, Sanja Fidler, and Raquel Urtasun. 3D graph neural networks for rgbd semantic segmentation. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 5199–5208, 2017. 2
- [54] Yue Qian, Junhui Hou, Sam Kwong, and Ying He. PUGeo-Net: a geometry-centric network for 3d point cloud upsampling. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 752–769. Springer, 2020. 2
- [55] Long Quan and Zhongdan Lan. Linear n-point camera pose determination. *IEEE Transactions on pattern analysis and machine intelligence*, 21(8):774–780, 1999. 1
- [56] Mahdi Rad and Vincent Lepetit. BB8: a scalable, accurate, robust to partial occlusion method for predicting the 3d poses of challenging objects without using depth. In *Proceedings* of the IEEE International Conference on Computer Vision (ICCV), pages 3828–3836, 2017. 1, 2
- [57] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28:91–99, 2015. 3, 6
- [58] Mahdi Saleh, Shun-Cheng Wu, Luca Cosmo, Nassir Navab, Benjamin Busam, and Federico Tombari. Bending graphs: Hierarchical shape matching using gated optimal transport. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11757– 11767, 2022. 2
- [59] Weijing Shi and Raj Rajkumar. Point-GNN: graph neural network for 3d object detection in a point cloud. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 1711–1719, 2020. 2
- [60] Martin Simonovsky and Nikos Komodakis. Dynamic edgeconditioned filters in convolutional neural networks on graphs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3693–3702, 2017. 2
- [61] Chen Song, Jiaru Song, and Qixing Huang. HybridPose: 6d object pose estimation under hybrid representations. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 431–440, 2020. 1, 7

- [62] Alessandro Sperduti and Antonina Starita. Supervised neural networks for the classification of structures. *IEEE Transactions on Neural Networks*, 8(3):714–735, 1997. 2
- [63] Hao Su, Charles R Qi, Yangyan Li, and Leonidas J Guibas. Render for CNN: viewpoint estimation in images using cnns trained with rendered 3d model views. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 2686–2694, 2015. 2
- [64] Yongzhi Su, Mahdi Saleh, Torben Fetzer, Jason Rambach, Nassir Navab, Benjamin Busam, Didier Stricker, and Federico Tombari. ZebraPose: coarse to fine surface encoding for 6dof object pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6738–6748, 2022. 2, 3, 4, 7, 8, 15, 16, 17, 20
- [65] Martin Sundermeyer, Zoltan-Csaba Marton, Maximilian Durner, Manuel Brucker, and Rudolph Triebel. Implicit 3d orientation learning for 6d object detection from rgb images. In Proceedings of the European Conference on Computer Vision (ECCV), pages 699–715, 2018. 2
- [66] Matthew Tancik, Pratul P. Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan T. Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. Advances in Neural Information Processing Systems, 2020. 16
- [67] Fulin Tang, Yihong Wu, Xiaohui Hou, and Haibin Ling. 3d mapping and 6d pose computation for real time augmented reality on cylindrical objects. *IEEE Transactions on Circuits* and Systems for Video Technology, 30(9):2887–2899, 2019.
- [68] Bugra Tekin, Sudipta N Sinha, and Pascal Fua. Real-time seamless single shot 6d object pose prediction. In *Proceed*ings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 292–301, 2018. 1, 2
- [69] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. FCOS: fully convolutional one-stage object detection. In Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), pages 9627–9636, 2019. 3, 7, 8
- [70] Jonathan Tremblay, Thang To, Balakumar Sundaralingam, Yu Xiang, Dieter Fox, and Stan Birchfield. Deep object pose estimation for semantic robotic grasping of household objects. In *Conference on Robot Learning*, pages 306–316. PMLR, 2018. 1
- [71] Jonathan Tremblay, Stephen Tyree, Terry Mosier, and Stan Birchfield. Indirect object-to-robot pose estimation from an external monocular rgb camera. In 2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pages 4227–4234. IEEE, 2020. 1
- [72] Shubham Tulsiani and Jitendra Malik. Viewpoints and keypoints. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 1510–1519, 2015. 2
- [73] Steffen Urban, Jens Leitloff, and Stefan Hinz. Mlpnp a realtime maximum likelihood solution to the perspective-n-point problem. arXiv preprint arXiv:1607.08112, 2016. 2
- [74] Nitika Verma, Edmond Boyer, and Jakob Verbeek. FeaSt-Net: feature-steered graph convolutions for 3d shape anal-

ysis. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 2598–2606, 2018. 2

- [75] Gu Wang, Fabian Manhardt, Federico Tombari, and Xiangyang Ji. GDR-Net: geometry-guided direct regression network for monocular 6d object pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16611–16621, 2021. 1, 2, 3, 6, 7, 8, 15, 16, 17, 20
- [76] Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Mingkui Tan, Xinggang Wang, et al. Deep high-resolution representation learning for visual recognition. *IEEE transactions* on pattern analysis and machine intelligence, 43(10):3349– 3364, 2020. 6, 7
- [77] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon. Dynamic graph cnn for learning on point clouds. ACM Transactions On Graphics (TOG), 38(5):1–12, 2019. 2, 3, 4, 5, 6, 13
- [78] Di Wu, Zhaoyong Zhuang, Canqun Xiang, Wenbin Zou, and Xia Li. 6D-VNet: end-to-end 6-dof vehicle pose estimation from monocular RGB images. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, pages 0–0, 2019. 1
- [79] Yu Xiang, Tanner Schmidt, Venkatraman Narayanan, and Dieter Fox. PoseCNN: a convolutional neural network for 6d object pose estimation in cluttered scenes. In *Robotics: Science and Systems*, 2018. 2, 6, 8, 13, 14, 15, 16, 17
- [80] Yan Xu, Kwan-Yee Lin, Guofeng Zhang, Xiaogang Wang, and Hongsheng Li. RNNPose: recurrent 6-dof object pose refinement with robust correspondence field estimation and pose optimization. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 14880–14890, 2022. 6
- [81] Sergey Zakharov, Ivan Shugurov, and Slobodan Ilic. DPOD: 6d pose object detector and refiner. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (*ICCV*), pages 1941–1950, 2019. 1, 2
- [82] Shaobo Zhang, Wanqing Zhao, Ziyu Guan, Xianlin Peng, and Jinye Peng. Keypoint-graph-driven learning framework for object pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (CVPR), pages 1065–1073, 2021. 2
- [83] Guangyuan Zhou, Huiqun Wang, Jiaxin Chen, and Di Huang. PR-GCN: a deep graph convolutional network with point refinement for 6d pose estimation. In *Proceedings of* the IEEE/CVF International Conference on Computer Vision (ICCV), pages 2793–2802, 2021. 2
- [84] Menglong Zhu, Konstantinos G Derpanis, Yinfei Yang, Samarth Brahmbhatt, Mabel Zhang, Cody Phillips, Matthieu Lecce, and Kostas Daniilidis. Single image 3d object detection and pose estimation for grasping. In 2014 IEEE International Conference on Robotics and Automation (ICRA), pages 3936–3943. IEEE, 2014. 1

#### 6. Supplemental Material

#### 6.1. Hyper-parameters in the Pose Solver

We use both RANSAC/PnP [32] and Progressive-X [1] when evaluating the results on the LM dataset [17], and we use Progressive-X for LM-O [2] and YCB-V [79] datasets. For both pose solvers, we set the threshold of reprojection error as 2 pixels. We run 150 iterations when using RANSAC/PnP and run 400 iterations when using Progressive-X.

# 6.2. Additional Ablation Experiments on LINEMOD Dataset

Theoretically, increasing the number of keypoint N leads to more candidate 3D-2D correspondences and enhances the robustness of pose estimation. In our current implementation, we adopt k = 20 in EdgeConv following [77], and N = 512 based on our available computation resources. We also conduct ablation studies of N and k on the LM dataset in Table 4, showing that larger N and k help improve the performance.

N	k	A 0.02d	ADD(-S) 0.05d	0.1d	2°2cm	5°5cm
128	10 15 20	29.4 29.2 29.8	81.3 81.0 82.0	96.4 96.1	75.6	98.8 98.6 98.7
256	20 10 15 20	<b>36.0</b> 32.0 33.6	84.2 83.4 82.9	96.8 96.8 96.4	79.1 78.1 75.8	98.8 98.8 98.8 98.8
512	10 15 20	30.4 29.9 35.7	82.0 82.8 84.5	96.6 96.3 <b>97.1</b>	76.3 76.4 <b>79.7</b>	98.7 98.5 <b>98.9</b>

Table 4: Ablation Study of N and k on the LM Dataset.

#### 6.3. Filtering Operation on LM-O and YCB-V

As discussed in the main paper, we empirically find that for a textureless object O with severe self-occlusions, filtering out the correspondences outside the visible segmentation masks  $M_{\text{vis}}$  can improve the pose estimation results. We quantify the self-occlusions of O using  $r_{\text{so}}(O)$ . As a common practice, the visibility of point  $P \in O$  from each viewpoint can be determined by checking the intersections between the camera rays and the object mesh. However, this may produce undesired results for our task. For example, the mesh of the bowl in the YCB-V dataset can be treated as a half sphere with very small thickness. When sampling the dense keypoints from the surface, we get keypoints from both outer side and inner side. For the keypoint on the inner side of the bowl, it is considered as easily self-occluded when we use ray intersections to determine

Object	$r_{\rm so}$	filtering
ape	0.356	×
can	0.650	1
cat	0.584	1
driller	0.657	1
duck	0.483	×
eggbox	0.529	1
glue	0.362	X
holep.	0.354	×

Table 5: Quantitative measure  $r_{so}$  of the self-occlusions of the objects on LM-O [2]. Since the objects do not have strong textures, for the objects with  $r_{so} \ge 0.5$ , we apply the filtering operation during inference, *i.e.*, discarding the correspondences outside the visible segmentation masks.

the visibility. However, since the bowl is textureless and the thickness of the mesh can be ignored, the keypoint is equivalent to the nearest surface point on the outer side, and should not be considered as easily self-occluded. Considering this issue and the slow computation speed, we instead use Hidden Point Removal (HPR) operator [27] to estimate the proportion V(P) of the viewpoints for which P is visible. For a keypoint with high V(P), it may be consistently misclassified as invisible by the HPR operator, so we ignore the points with V(P) < 0.2 estimated by the HPR operator.

We report the value of  $r_{so}(O)$  for each object O of the LM-O dataset in Table 5. Since these objects do not have strong textures, we apply the filtering operation during the inference for the objects with  $r_{so}(O) \ge 0.5$ .

For the objects that requires filtering operation, we report the ADD(-S) metric without filtering in Table 6. We also report the results of using different segmentation masks to filter the correspondences in Table 6. Without the filtering operation, the ADD(-S) values decreases for all the objects. Since all the 2D projections should be located within the full segmentation mask  $M_{\rm full}$ , using  $M_{\rm full}$  to filter the correspondences aims to discard the wrong predictions outside the object area. However, it does not improve the final estimations consistently, which indicates that we still need to discard more unstable correspondences within the object area.

We report the values of  $r_{so}$  for the textureless objects in the YCB-V dataset in Table 7. According to Table 7, only one textureless object, *i.e.*, 061\_foam\_brick, requires filtering operation due to severe self-occlusions.

We further report the ADD(-S) metric w.r.t. the filtering operation for 061\_foam\_brick in Table 8. The ADD(-S) of 061\_foam\_brick remains the same without filtering operation or using  $M_{\rm full}$  rather than  $M_{\rm vis}$  in the filtering operation. This observation suggests that the localization of the easily self-occluded regions may become stable after

Object	w/o Filter	w/ Filter ( $M_{\rm full}$ )	w/ Filter ( $M_{\rm vis}$ )
can	95.2	95.1	95.7
cat	62.0	61.3	62.3
driller	92.6	92.6	93.7
eggbox	68.8	69.6	70.0

Table 6: **ADD(-S) metrics on LM-O [2] w.r.t. the filtering operation.** "w/o Filter" denotes using all predicted correspondences to compute the pose. "w/ Filter  $(M_{\rm full})$ " denotes discarding the correspondences outside the full segmentation mask  $M_{\rm full}$ , while "w/ Filter  $(M_{\rm vis})$ " denotes discarding the correspondences outside the full segmentation mask  $M_{\rm vis}$ .

Object	$r_{\rm so}$	filtering
011_banana	0.240	×
019_pitcher_base	0.221	×
024_bowl	0.498	×
025_mug	0.108	X
036_wood_block	0.438	×
037_scissors	0.365	×
051_large_clamp	0.163	×
052_extra_large_clamp	0.138	×
061_foam_brick	0.542	1

Table 7: Quantitative measure  $r_{so}$  of the self-occlusions of the textureless objects on YCB-V [79]. For the object with  $r_{so} \ge 0.5$ , we apply the filtering operation during inference, *i.e.*, discarding the correspondences outside the visible segmentation masks.

Object	w/o Filter	$ $ w/ Filter ( $M_{\rm full}$ )	w/ Filter ( $M_{\rm vis}$ )
008_pudding_box	66.4	71.0	86.5
061_foam_brick	<b>87.2</b>	87.2	87.2

Table 8: **ADD(-S) metrics on YCB-V** [79] w.r.t. the filtering operation. "w/o Filter" denotes using all predicted correspondences to compute the pose. "w/ Filter  $(M_{\rm full})$ " denotes discarding the correspondences outside the full segmentation mask  $M_{\rm full}$ , while "w/ Filter  $(M_{\rm vis})$ " denotes discarding the correspondences outside the full segmentation mask  $M_{\rm vis}$ .

380,000 training steps. We further investigate the results of 061\_foam\_brick after different training steps in Table 9. After 200,000 steps, the ADD(-S) without filtering is inferior to the result of discarding correspondences outside  $M_{\rm vis}$ . This observation implies that the localization of the easily self-occluded regions are unstable with fewer training steps.

Besides textureless objects with severe self-occlusions, we also apply filtering operation on 008\_pudding\_box

Steps	w/o Filter	w/ Filter ( $M_{\rm full}$ )	w/ Filter ( $M_{\rm vis}$ )
200k	86.1	85.4	86.8
380k	87.2	87.2	87.2

Table 9: **ADD(-S) metrics of 061\_foam\_brick with differ**ent training steps. "w/o Filter" denotes using all predicted correspondences to compute the pose. "w/ Filter  $(M_{\rm full})$ " denotes discarding the correspondences outside the full segmentation mask  $M_{\rm full}$ , while "w/ Filter  $(M_{\rm vis})$ " denotes discarding the correspondences outside the full segmentation mask  $M_{\rm vis}$ .

from the YCB-V dataset. As shown in Figure 5, 008\_pudding\_box is severely occluded by 009\_gelatin\_box. We regard 009\_gelatin\_box as a distraction object for the keypoint localization task of 008\_pudding\_box, since these objects share similar appearances, especially the texts (i.e., "JELL-O"). Such severe occlusions by the same distraction object exist in all the test images of 008\_pudding\_box, and can be automatically detected by checking the object detection results. Thus we discard the correspondences outside  $M_{\rm vis}$  to remove the unstable localization results due to the occlusions by the distraction object. We also report the ADD(-S) metric without filtering and using  $M_{\rm full}$  in filtering in Table 8. Using either  $M_{\rm full}$  or  $M_{\rm vis}$  to filter the correspondences improve the pose estimation results compared with using all predicted correspondences. This indicates that the filtering operation can remove extreme outliers that are far from 008\_pudding\_box to improve the pose estimation. Using  $M_{\rm vis}$  in the filtering operations obtains better results than  $M_{\rm full}$ , which demonstrates that the localization results of the keypoints occluded by the distraction object are not accurate enough for recovering the pose.

#### 6.4. Evaluation of 2D-3D Correspondences

The evaluation results in the main paper focus on the final estimated poses. We additionally evaluate the quality of the established dense correspondences before RANSAC. Specifically, for each test sample, we reproject the 3D keypoints by the ground truth pose and compute the mean distance between the reprojection results and predicted 2D locations. For symmetric objects, we use the equivalent rotation closest to our final estimated pose. To obtain the inlier ratio of the estimated correspondences, we regard a keypoint as an inlier if its reprojection error is less than 5 pixels. We compute the average reprojection error and inlier ratio for each object and report the average values over the whole dataset in Table 10.

#### 6.5. BOP Results on LM-O and YCB-V

We report the performance of our method on LM-O and YCB-Video using the evaluation metrics from BOP chal-



Figure 5: Example of test images for 008\_pudding\_box from the YCB-V dataset. We visualize the zoomed-in RoI based on the detection results. For all test images, 008\_pudding\_box (the brown box) is severely occluded by 009\_gelatin\_box (the red box).

Dataset	LM	LM-O	YCB-V
reprojection error (pixel)	3.4	14.4	10.9
inlier ratio (%)	88.4	67.8	39.6

Table 10: Evaluation results of predicted dense correspondences.

lenge [19] in Table 11 and Table 12, respectively. We mainly select baselines from officially published work. We also include the results of GDRNPP [42] for reference, which improves upon GDR-Net [75] with implementation skills including stronger domain randomization, more powerful detectors, *etc.*, to compensate for the domain gap between training and test images. Without these implementation skills, our method still achieves comparable performance with the state-of-the-art methods, including the refinement based method [41].

#### 6.6. Detailed Results of YCB-V

We report the detailed evaluation metrics of each object on YCB-V dataset [79] in Table 13 and Table 14. Our method outperforms previous methods w.r.t. ADD(-S) and AUC of ADD(-S), and achieves comparable performance with state of the art w.r.t. AUC of ADD-S.

#### 6.7. Qualitative Results

We provide additional qualitative results for LM-O [2] and YCB-V [79] in Figure 6 and Figure 7, respectively. We render the 3D CAD model based on the predictions of CheckerPose, and highlight the contour in green. We also highlight the ground truth contour in blue. For better visualization, we crop the images and we also show the original

Method	$AR_{MSPD}$	AR <sub>MSSD</sub>	$AR_{VSD}$	AR
SurfEmb [14]	85.1	64.0	49.7	66.3
Coupled [41]	83.1	63.3	50.1	65.5
Zebra [64]	88.0	72.1	55.2	71.8
NCF [24]	_	_	_	63.2
PFA [20]	83.7	66.1	52.3	67.4
CRT-6D [5]	83.7	64.0	50.4	66.0
GDRNPP [42]	88.7	70.1	54.9	71.3
Ours	87.3	72.3	53.7	71.1

Table 11: **Results on LM-O dataset under BOP setup** [19]. The results of Coupled [41] and NCF [24] are obtained from the original paper, and the results of other methods are obtained from https://bop.felk. cvut.cz/leaderboards/. We highlight the best result in red color, and the second best result in blue color. "-" denotes unavailable results.

Method	$AR_{MSPD}$	AR <sub>MSSD</sub>	AR <sub>VSD</sub>	AR
SurfEmb [14]	77.3	62.0	54.8	64.7
Coupled [41]	85.2	83.5	78.3	82.4
Zebra [64]	86.4	83.0	75.1	81.5
NCF [24]	_	_	_	77.5
PFA [20]	84.9	81.4	75.8	80.7
SC6D [4]	80.4	79.6	69.5	76.5
CRT-6D [5]	77.4	77.6	70.6	75.2
GDRNPP [42]	86.9	84.6	76.0	82.5
Ours	85.3	84.4	70.7	80.1

Table 12: **Results on YCB-Video dataset under BOP setup [19].** The results of Coupled [41] and NCF [24] are obtained from the original paper, and the results of other methods are obtained from https://bop.felk. cvut.cz/leaderboards/. We highlight the best result in red color, and the second best result in blue color. "-" denotes unavailable results.

input image on the left for LM-O and YCB-V.

Furthermore, we provide more keypoint localization results of duck, bowl, and banana in Figure 8. For better visualization we only plot eight keypoints that are evenly distributed over the object surface. While our network directly outputs the 2D locations, the results of other dense methods [64, 75] are computed by projecting the keypoints using the estimated poses. Considering the symmetry of the bowl, we use the equivalent rotations closest to our prediction to project the keypoints of bowl.

#### 6.8. Failure Cases and Future Work

We visualize typical failure cases in Figure 9. As shown

Method	SegDriven[22]	S.Stage[21]	RePose [26]	GDR [75]	Zebra [64]	DProST [47]	Ours
002_master_chef_can	33.0	-	-	41.5	62.6	-	45.9
003_cracker_box	44.6	-	-	83.2	98.5	-	94.2
004_sugar_box	75.6	-	-	91.5	96.3	-	98.3
005_tomato_soup_can	40.8	-	-	65.9	80.5	-	83.2
006_mustard_bottle	70.6	-	-	90.2	100.0	-	99.2
007_tuna_fish_can	18.1	-	-	44.2	70.5	-	88.9
008_pudding_box	12.2	-	-	2.8	99.5	-	86.5
009_gelatin_box	59.4	-	-	61.7	97.2	-	86.0
010_potted_meat_can	33.3	-	-	64.9	76.9	-	70.0
011_banana	16.6	-	-	64.1	71.2	-	96.0
019_pitcher_base	90.0	-	-	99.0	100.0	-	100.0
021_bleach_cleanser	70.9	-	-	73.8	75.9	-	<b>89.8</b>
024_bowl*	30.5	-	-	37.7	18.5	-	68.0
025_mug	40.7	-	-	61.5	77.5	-	89.0
035_power_drill	63.5	-	-	78.5	97.4	-	95.9
036_wood_block*	27.7	-	-	59.5	87.6	-	58.7
037_scissors	17.1	-	-	3.9	71.8	-	62.4
040_large_marker	4.8	-	-	7.4	23.3	-	18.8
051_large_clamp*	25.6	-	-	69.8	87.6	-	95.4
052_extra_large_clamp*	8.8	-	-	90.0	98.0	-	95.6
061_foam_brick*	34.7	-	-	71.9	99.3	-	87.2
MEAN	39.0	53.9	62.1	60.1	80.5	65.1	81.4

Table 13: Detailed results on YCB-V [79] w.r.t. ADD(-S). (\*) denotes symmetric objects and "-" denotes unavailable results.

in Figure 9 (a) and (b), the textureless object eggbox from LM-O dataset is severely occluded by a toy car, and a distraction object with similar color also partially appears in the input RoI. As a result, the estimated 2D projections are shifted towards the distraction object. We also present a failure case of objects with textures in Figure 9 (c) and (d). The object in interest is 002\_master\_chef\_can from YCB-V dataset, which is geometrically symmetric. Though the texture is almost symmetric as well, the barcode only appears on one side of the object, which causes the asymmetry. For the given input RoI, the keypoints are localized in the opposite directions, w.r.t. the central axis.

To improve the localization results, one future direction is the selection of 3D keypoints. Since we adopt farthest point sampling algorithm to obtain evenly distributed keypoints, we ignore other factors to make the keypoints more representative. For example, the issue of 002\_master\_chef\_can may be solved by sampling more keypoints in the barcode area. Besides, no positional encoding [45, 66] is leveraged in graph feature aggregation and image feature fusion operations. Such encoding can provide additional cues for textureless regions. In future, we will explore the positional encoding to enhance the keypoint localization process.

Method	CosyPose [31]		GDR-Net[75]		ZebraPose[64]		DProST [47]	Ours	
Metric	AUC of	AUC of	AUC of	AUC of	AUC of	AUC of	AUC of	AUC of	AUC of
	ADD-S	ADD(-S)	ADD-S	ADD(-S)	ADD-S	ADD(-S)	ADD(-S)	ADD-S	ADD(-S)
002_master_chef_can	-	-	96.3	65.2	93.7	75.4	-	87.5	67.7
003_cracker_box	-	-	97.0	88.8	93.0	87.8	-	93.2	86.7
004_sugar_box	-	-	98.9	95.0	95.1	90.9	-	95.9	91.7
005_tomato_soup_can	-	-	96.5	91.9	94.4	90.1	-	94.0	89.9
006_mustard_bottle	-	-	100.0	92.8	96.0	92.6	-	95.7	90.9
007_tuna_fish_can	-	-	99.4	94.2	96.9	92.6	-	97.5	94.4
008_pudding_box	-	-	64.6	44.7	97.2	95.3	-	94.9	91.5
009_gelatin_box	-	-	97.1	92.5	96.8	94.8	-	96.1	93.4
010_potted_meat_can	-	-	86.0	80.2	91.7	83.6	-	86.4	80.4
011_banana	-	-	96.3	85.8	92.6	84.6	-	95.7	90.1
019_pitcher_base	-	-	99.9	98.5	96.4	93.4	-	95.8	91.9
021_bleach_cleanser	-	-	94.2	84.3	89.5	80.0	-	90.6	83.2
024_bowl*	-	-	85.7	85.7	37.1	37.1	-	82.5	82.5
025_mug	-	-	99.6	94.0	96.1	90.8	-	96.9	92.7
035_power_drill	-	-	97.5	90.1	95.0	89.7	-	94.7	88.8
036_wood_block*	-	-	82.5	82.5	84.5	84.5	-	68.3	68.3
037_scissors	-	-	63.8	49.5	92.5	84.5	-	91.7	81.6
040_large_marker	-	-	88.0	76.1	80.4	69.5	-	83.3	72.3
051_large_clamp*	-	-	89.3	89.3	85.6	85.6	-	90.0	90.0
052_extra_large_clamp*	-	-	93.5	93.5	92.5	92.5	-	91.6	91.6
061_foam_brick*	-	-	96.9	96.9	95.3	95.3	-	94.1	94.1
MEAN	89.8	84.5	91.6	84.3	90.1	85.3	77.4	91.3	86.4

Table 14: **Detailed results on YCB-V** [79] w.r.t. AUC of ADD-S and ADD(-S). As in [79], symmetric metric is used for all objects in ADD-S while only for symmetric objects in ADD(-S). (\*) denotes symmetric objects.



Figure 6: **Qualitative results on the LM-O dataset.** For each image on the left, we visualize the 6D pose by rendering the 3D CAD models and highlighting the contours on the right. Blue color denotes ground truth and green color denotes the prediction from CheckerPose.































Dominc













TELLC

CHEE7AT











Figure 7: Qualitative results on the YCB-V dataset. For each image on the left, we visualize the 6D pose by rendering the 3D CAD models and highlighting the contours on the right. Blue color denotes ground truth and green color denotes the prediction from CheckerPose.



Figure 8: **Visualization of keypoint localization.** Each column visualizes the keypoint location results of ZebraPose [64], GDR-Net [75], our method, and ground truth. While our network directly outputs the 2D locations, the results of other dense methods [64, 75] are computed by projecting the keypoints using the estimated poses.



Figure 9: Failure cases. We provide the localization results of eight keypoints that are inliers of the estimated poses.