# KDFNet: Learning Keypoint Distance Field for 6D Object Pose Estimation

Xingyu Liu, Shun Iwase, Kris M. Kitani

*Abstract*— We present KDFNet, a novel method for 6D object pose estimation from RGB images. To handle occlusion, many recent works have proposed to localize 2D keypoints through pixel-wise voting and solve a Perspective-n-Point (PnP) problem for pose estimation, which achieves leading performance. However, such voting process is direction-based and cannot handle long and thin objects where the direction intersections cannot be robustly found. To address this problem, we propose a novel continuous representation called Keypoint Distance Field (KDF) for projected 2D keypoint locations. Formulated as a 2D array, each element of the KDF stores the 2D Euclidean distance between the corresponding image pixel and a specified projected 2D keypoint. We use a fully convolutional neural network to regress the KDF for each keypoint. Using this KDF encoding of projected object keypoint locations, we propose to use a distance-based voting scheme to localize the keypoints by calculating circle intersections in a RANSAC fashion. We validate the design choices of our framework by extensive ablation experiments. Our proposed method achieves state-of-the-art performance on Occlusion LINEMOD dataset with an average ADD(-S) accuracy of 50.3% and TOD dataset mug subset with an average ADD accuracy of 75.72%. Extensive experiments and visualizations demonstrate that the proposed method is able to robustly estimate the 6D pose in challenging scenarios including occlusion.

## I. INTRODUCTION

Estimating the 6D pose of a rigid object, i.e. rotation and translation in 3D space, is a core problem in computer vision and is crucial for applications such as robotic manipulation and augmented reality (AR). We focus on the setting of 6D object pose estimation of a rigid object from RGB images where the 3D textured mesh model of the object is known. Early attempts to this problem include direct pose regression using neural networks in an end-to-end fashion [1, 2]. Recently, 6D pose estimation methods that use object keypoints as intermediate representation have been successful and achieve leading performance on various benchmarks [3–6]. By definition, keypoints are 3D points attached to an object model and are usually a subset of the object surface points. In keypoint-based methods, 2D projections of 3D object keypoints or centers are first located on the image and then the 6D pose can be recovered from such 2D-3D correspondences by solving a Perspective-n-Point (PnP) problem.

There are mainly two types of methods for localizing 2D keypoints: heatmap-based [3, 4] and voting-based [2, 5, 6]. Heatmap-based methods predict probability heatmaps of a keypoint over the image and localize it through an integral operation with the image coordinate map. Though heatmap-based method achieves strong performance on problems such as human pose estimation [7], it is known to be vulnerable to occlusion, because the features of the occluder near the keypoint location can significantly affect the predicted probability map. In voting-based methods, the visible parts of the object hallucinate and vote for the 2D locations of the invisible keypoints [5]. The training objectives of the votings are independent of the occluders. Therefore, compared to heatmap-based methods, voting-based methods are more robust to occlusions and achieve stronger performance on object pose estimation benchmarks where occlusions are important.

The voting schemes of existing methods are direction-based, i.e. every object pixel predicts the 2D direction to the keypoints and the keypoint hypotheses are the intersections of the direction votes [2, 5, 6]. Direction-based voting methods are built upon an important assumption: the angles between the voting directions are large enough so that the keypoint hypotheses can be reliably found by computing the intersections of voted directions. However, this assumption does not hold for long and thin objects where most voting directions concentrate in a small range and the keypoint hypotheses cannot be found or are extremely sensitive to noise in direction.

To address this problem, we propose a novel representation for object keypoint named *Keypoint Distance Field (KDF)*.
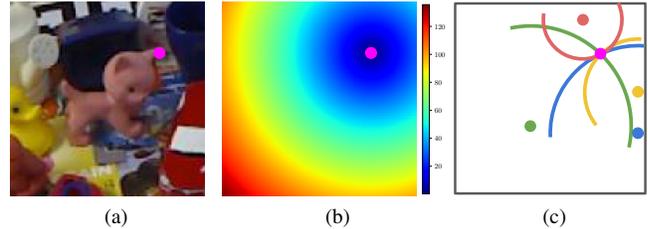


Fig. 1: (a) Our method estimates object 6D pose from RGB images by localizing object keypoints and solving a Perspective-n-Point (PnP) problem. (b) For every keypoint, our deep model regresses **Keypoint Distance Field (KDF)**, a novel 2D map representation for localizing projected keypoints. (c) To recover the 2D locations of the projected keypoints from regressed KDFs, we propose a **distance-based voting scheme** that includes calculating circle intersections in RANSAC fashion.

The authors are with the Robotics Institute, Carnegie Mellon University, 5000 Forbes Ave, Pittsburgh, PA 15213, United States. Contact: {xingyul3, siwase, kkitani}@cs.cmu.edu

Defined as a 2D map of the same spatial size as the RGB image, the KDF is the function of 2D Euclidean distances to a certain projected keypoint. Given a perfect KDF, the 2D location of the projected keypoint can be easily recovered. Note that KDF is able to represent keypoints that are invisible or even outside the image field of view. Given $K$ keypoints, there are $K$ KDFs defined.

With the KDF representation, in this paper, we introduce a novel keypoint-based 6D object pose estimation framework named KDFNet. The core of KDFNet is a fully convolutional neural network (CNN) that predicts the KDFs of the object keypoints through per-pixel regression. To efficiently recover the 2D locations of the projected keypoints from the predicted KDFs, we propose a distance-based voting scheme. The voting hypotheses are generated through circle intersections where the centers of the circles are the pixel voters and the radii of the circles are the predicted keypoint distance values. The projected keypoints are the hypothesis with the maximum consensus of distance prediction among pixel voters. The core idea is illustrated in Figure 1.

We evaluated our framework on one of the most popular 6D pose estimation benchmarks, Occlusion LINEMOD [8], and compare it against related baseline approaches. On Occlusion LINEMOD, our method achieves an accuracy of 50.3%, significantly outperforms related baselines such as [5] and the current state-of-the-art HybridPose [6]. In addition, we evaluate the keypoint estimation of KDFNet on TOD [4], a stereo-RGB dataset for object pose estimation, and compare against previous stereo-RGB keypoint methods including KeyPose [4]. On TOD dataset, our method also achieves state-of-the-art results.

## II. RELATED WORK

**Learning Object Keypoints.** Many previous works have explored deep learning methods for localizing 3D keypoints of an object [3, 4, 9–11] or a human [7, 12] from an RGB image to estimate their poses. The core of these methods is to predict the probability heatmap for the 2D keypoint and localizing it through integral with coordinate map. This idea is also used in 2D object detection where the corners of 2D object bounding boxes are formulated as keypoints and are localized through probability heatmaps [13, 14]. Besides heatmap, 2D direction field representation has also been proposed for localizing keypoints [2, 5]. Our method also uses keypoint 2D location as a bridge to 6D object pose. We propose a novel distance field representation of keypoints and the associated voting scheme.

**Object Pose Estimation** The earliest attempts for object pose estimation used CNN to directly regress 6D pose [1, 2, 15]. Rough 2D bounding boxes of the objects may be predicted first to localize the objects more accurately [2]. However, directly estimating the 6D poses using CNN assumes the neural network can implicitly remember the images of the object in all possible 6D poses which is difficult and prone to occlusion and background clutter. Instead, methods such as [16, 17] leverages object coordinate maps as the dense 2D-3D correspondence representation for

6D pose. Compared to dense methods, keypoints are a more flexible representation and are used in [4–6, 11]. Our pose estimation method is keypoint-based and predicts distance maps to localize the keypoints.

**Voting for Understanding Objectness.** Voting has also been used in 3D object detection where objectness can be inferred from voting consensus [18–21]. Specifically, direction-based voting methods have been adopted by previous works to robustly localize object centers [2] or object keypoints [5, 10] on the images where the voting scores are based on the number of direction inliers. Instead of calculating the mean of predictions, voting-based methods find the maximum consensus among predictions. Therefore, voting-based methods are known to be robust to noise and occlusions. Our method also leverages voting for robustly detecting keypoints. Different from previous methods, our voting scheme is distance-based.

## III. METHOD

Given an RGB image of a rigid object, the goal of our framework is to predict the rotation $\mathbf{R}$ and translation $\mathbf{T}$ of the object. One of the popular approaches used in previous methods is pixel-wise direction voting of object keypoints. In this approach, every pixel votes for the direction to the keypoints and the intersections of the voted directions yield keypoint candidates whose voting scores are then evaluated based on the ratio of voter inliers. The candidate with the highest scores is decided as the keypoint estimate, and the object pose can be recovered from all keypoint estimates by solving PnP. The intuition behind this method is that given an object, the location of the invisible keypoint can be inferred from the visible parts [2, 5, 6]. Though this method is known to be robust to occlusion, it is built upon the assumption that the voting directions can reliably yield keypoint candidates regardless of the geometry of the objects. This assumption is not true. A counterexample is illustrated in Section IV and Figure 4.

To address this problem, we proposed a novel framework for estimating 6D pose of 3D objects from RGB input. The core of our method is a novel representation for object keypoint named **Keypoint Distance Field (KDF)**. Inspired by recent works [4, 5], we first predict KDF to localize 2D keypoints through voting, then compute object 6D pose by solving a PnP problem. Figure 2 illustrates our framework. In this section, we will introduce the representation of KDF and describe the corresponding voting scheme. Then we present the implementation details of our approach.

### A. Representation of Keypoint Distance Field

Keypoint Distance Field (KDF) is specifically defined for each 2D projected keypoint. Suppose that the height and width of the input RGB image are $H$ and $W$ respectively and that the object has $K$ keypoints defined. KDF for the $k$-th projected keypoint $(u_k, v_k)$ is a two-dimensional array $\mathbf{D} \in \mathbb{R}^{H \times W}$ of the same size as the input image. The element of the KDF located at $(u, v)$ stores the 2D Euclidean distance
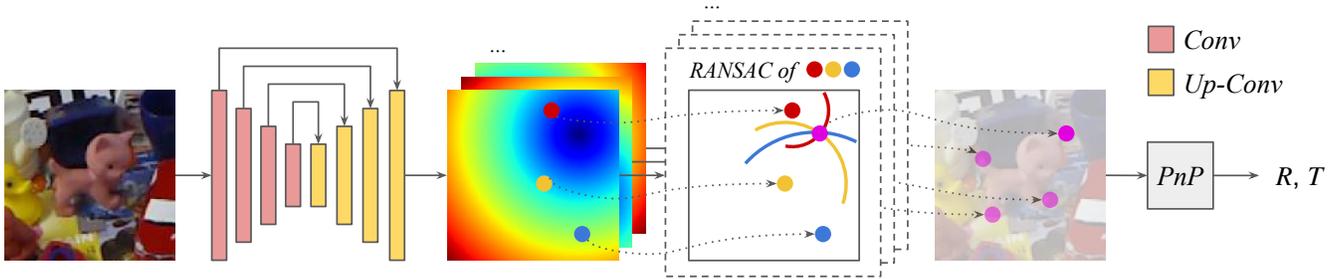
Fig. 2: **Overall Architecture of KDFNet**. We used a fully convolutional neural network to predict the KDFs defined by all keypoints of the objects. For every predicted KDF, the distance-based voting scheme randomly samples three positions $\mathbf{p}_{i_1}$, $\mathbf{p}_{i_2}$, $\mathbf{p}_{i_3}$ for $N$ times to generate keypoint hypotheses. The keypoints location are determined by selecting the hypothesis with the highest voting score defined by Equation (7). Finally, the 6D pose is calculated by solving a PnP problem from the predicted object keypoint locations.

between the element and $(u_k, v_k)$

$$\mathbf{D}_{u,v} = \left|\left|[u, v] - [u_k, v_k]\right|\right|_2 \quad (1)$$

Note that the KDF can still be defined even when $(u_k, v_k)$ is outside of the image. We use a fully convolutional neural network to regress $\mathbf{D}$. Inspired by previous work on 2D object detection [22], we adopt the following parameterization for regression

$$t_{\mathbf{D}} = \log(\mathbf{D}/r) \quad (2)$$

where $r \in \mathbb{R}^+$ is a hyperparameter. The value of $r$ is chosen to be close to the geometric mean of maximum and minimum possible distances so that the lower and upper bounds of parameterized distance are symmetric about zero, which is easier for the neural network to regress. For example, with a $256 \times 256$ image input, $r$ can be set to be 16.

**Loss Function.** To predict a set of continuous values, we can either regress the values directly or convert it to classification problems with multiple discretized values. Since the range of possible distance values is large even after parameterization, we choose to use direct regression for KDF prediction to avoid large discretization error. We use a standard soft $L_1$ loss function for parameterized distance

$$\mathcal{L}(\mathbf{D}, \mathbf{D}^*) = \frac{1}{H \times W} \begin{cases} \frac{1}{2e}(t_{\mathbf{D}^*} - t_{\mathbf{D}})^2, & |t_{\mathbf{D}^*} - t_{\mathbf{D}}| < e \\ |t_{\mathbf{D}^*} - t_{\mathbf{D}}| - \frac{e}{2}, & |t_{\mathbf{D}^*} - t_{\mathbf{D}}| \geq e \end{cases} \quad (3)$$

where $\mathbf{D}^*$ is the ground truth KDF and $e$ is a threshold value. Note the loss function is the mean of all elements.

**Objects with Symmetry.** Symmetric objects can cause ambiguity among mutually equivalent rigid transforms. A group of keypoints placed on the object may thus be indistinguishable. Inspired by [4], to deal with objects with discrete symmetry, we define keypoints such that each keypoint is part of the symmetry permutation group and then apply a permutation loss during training. Suppose keypoints $k_1, k_2, \ldots, k_S$ are mutually equivalent and the symmetric permutation of their indices is $\text{Sym}(S) = \{(i_1, i_2, \ldots, i_S)\}$. Then the KDF regression loss under symmetric permutation is

$$\mathcal{L}_{\text{symm}} = \min_{s \in \text{Sym}(S)} \mathcal{L}(\mathbf{D}, \mathbf{D}_s^*) \quad (4)$$

where $\mathbf{D}_s^*$ is the ground truth KDFs after applied with permutation $s$ on keypoints indices.

For objects with continuous symmetry such as cylinders, the symmetric permutation group is infinite. Therefore, Equation (4) cannot be used. Alternatively, we define keypoints on the rotation axis, and predicting at least two keypoints is sufficient for determining the 6D pose of the object with the ambiguity of the rotation.

### B. Distance-based Voting Scheme

In 2D heatmap-based methods [3, 4], only the small image area that is close to a keypoint and has high probabilities can have an effect on the final predictions. Therefore, the prediction is affected by the pixels in that small area and suffers from occlusion. On the contrary, we apply a RANSAC-based voting scheme to take distance predictions from more pixels into account, so that occluded keypoints or even keypoints outside the image can be robustly handled.

The first step is to generate hypotheses of keypoint locations through sampling. For keypoint $k$, we first determine the set of elements $\mathcal{P} = \{\mathbf{p}_i\}$ on the $k$-th KDF map that will participate in the voting. From $\mathcal{P}$, we randomly sample three elements $\mathbf{p}_{i_1}, \mathbf{p}_{i_2}, \mathbf{p}_{i_3} \sim \mathcal{P}$. For every $j \in \{i_1, i_2, i_3\}$, the set of hypothesized keypoint locations predicted by pixel $\mathbf{p}_j$ is a circle whose center is $\mathbf{p}_j$ and radius equals to $\mathbf{D}_{\mathbf{p}_j}$ — the predicted KDF value at $\mathbf{p}_j$. Given a perfect KDF, all three circles are supposed to intersect at one location. In practice, we find the best possible location(s) agreed by at least two of the three circles using the following procedure. Each pair of the three circles returns two intersections, but at most one intersection is valid as a keypoint hypothesis. The third circle is used to decide the valid hypothesis based on which intersection is closer to the third circle. Mathematically, suppose the two circles predicted by $\mathbf{p}_{i_1}$ and $\mathbf{p}_{i_2}$ intersects at $\mathbf{h}_{i_1}$ and $\mathbf{h}_{i_2}$ which can be obtained by jointly solving the following two quadratic equations

$$||\mathbf{h} - \mathbf{p}_{i_1}||_2 = \mathbf{D}_{\mathbf{p}_{i_1}}, ||\mathbf{h} - \mathbf{p}_{i_2}||_2 = \mathbf{D}_{\mathbf{p}_{i_2}} \quad (5)$$

Then the valid hypothesis generated by $\mathbf{p}_{i_1}$ and $\mathbf{p}_{i_2}$ is given by

$$\mathbf{h}_{i_1, i_2} = \operatorname*{arg\,min}_{\mathbf{h} \in \{\mathbf{h}_{i_1}, \mathbf{h}_{i_2}\}} \left| ||\mathbf{h} - \mathbf{p}_{i_3}||_2 - \mathbf{D}_{\mathbf{p}_{i_3}} \right| \quad (6)$$
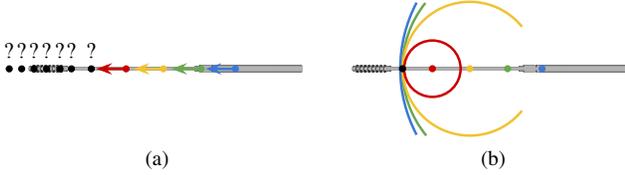
(a)                          (b)

Fig. 3: Voting scheme comparison. (a) **Direction-based voting** cannot reliably localize a point when the angles between directions are too small. (b) **Distance-based voting** can still accurately localize a point in this situation.

The other two valid hypotheses $\mathbf{h}_{i_2,i_3}$ and $\mathbf{h}_{i_3,i_1}$ can be obtained similarly. In total, there are three valid hypotheses generated. The above sampling and hypothesis generation is repeated for $N$ times to generate $3N$ hypotheses $\mathcal{H} = \{\mathbf{h}_l \mid l = 1, 2, \ldots, 3N\}$ for the $k$-th keypoint. Next, all elements from $\mathcal{P}$ vote for these hypotheses. The voting score $v_l$ of $\mathbf{h}_l$ is the number of elements whose distance prediction error is within the threshold $\theta$

$$v_l = \sum_{\mathbf{p} \in \mathcal{P}} \mathbb{I}\Big( \big| \|\mathbf{h}_l - \mathbf{p}\|_2 - \mathbf{D}_{\mathbf{p}} \big| < \theta \Big) \qquad (7)$$

where $\mathbb{I}(\cdot)$ is the indicator function. Then the $k$-th keypoint is determined as the hypothesis with the highest voting score $\mathbf{h}_{l_{\max}}$, where $l_{\max} = \arg\max_l \{v_l\}$. In this way, we find the location agreed by most of the voters within an error range. The above process is repeated for all KDFs to predict all keypoints.

### C. Overall Framework and Implementation

The KDF prediction CNN is instantiated by a segmentation network with ResNet [23] as encoder backbone and additional up-convolution layers and skip connections as the decoder. The ResNet backbone is initialized with an ImageNet pre-trained model and then finetuned on 6D pose estimation datasets. During the training of KDF regression network, we randomly generate bounding box crops around the object to introduce more variations. Random photometric data augmentation is used during training. The overall architecture is illustrated in Figure 2.

Though loss in Equation (3) is applied to all elements, to reduce regression error, one can choose to apply the loss only on elements within a certain keypoint distance during training if the image size is too large. In this case, the predicted KDF and the ground truth KDF will be different at large keypoint distances during inference, and a rough initial estimate of the object location from detection or segmentation is needed and $\mathcal{P}$ in Equation (7) only includes elements within a certain keypoint distance. We will show in Section V-E that such training loss strategy does not affect the voting of the 2D keypoints. In practice, within a rough initial range of keypoint 2D location, we randomly sample 4,096 pixels as the set of voters $\mathcal{P}$. Among the voters, we sample 1,024 three tuples of pixels to generate 3,072 keypoint hypotheses whose voting scores are determined by

TABLE I: **Toy dataset results.** Evaluation metric is **2D projection error at 1px threshold**.

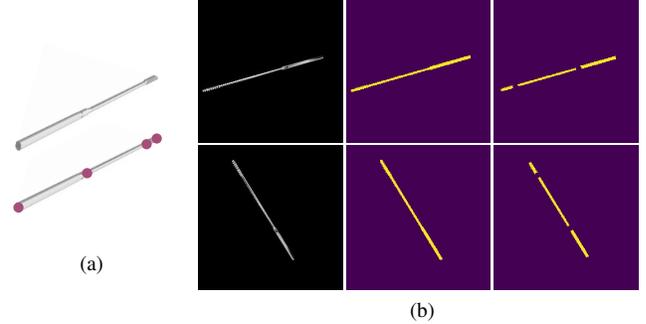| method | PVNet [5] | KDFNet (ours) |
|---|---|---|
| GT mask | 93.2 | **95.5** |
| GT mask + keypoint occlusions | 75.1 | **95.6** |



(a)

(b)

Fig. 4: **Toy dataset.** (a) The 3D model of the object used in the toy dataset: a medical swab stick. There are four keypoints defined which are shown as purple circles. (b) Two data samples drawn from the toy dataset. Columns 1-3 are respectively the RGB image input, ground truth object mask, and ground truth object mask with additional occlusions applied around the keypoints.

all the sampled voters. Given the predicted 2D keypoint locations for each object, the 6D pose can be computed by solving a PnP problem.

### IV. DIRECTION-BASED VS. DISTANCE-BASED VOTING: A TOY EXPERIMENT

The key difference between our method and previous works [2, 5, 6] is the predicted representation used in voting, i.e. keypoint distance vs. keypoint direction. In this section, we construct a toy synthetic dataset where previous direction-based voting methods fail in localizing the keypoints. Through this simple dataset, we show the drawbacks of direction-based representation and the advantage of our proposed KDF representation.

The object we use in the dataset is the 3D mesh model of a medical swab stick for COVID-19 testing, illustrated in Figure 4(a). The dataset consists of synthetically rendered images of the swab stick in various 6D poses. The translation part of the poses $\mathbf{T}$ is the same for every image and the rotation part $\mathbf{R}$ is uniformly randomly sampled. The backgrounds are all black. The dataset has 30,000 training and 3,000 validation images. The first column of Figure 4(b) illustrates two examples of images in our dataset. Semantic masks are provided for every image. We define four keypoints along the medical swab stick.

The baseline we compare our KDFNet against is PVNet [5], a direction-based keypoint method. To fairly compare, we assume the voters of both PVNet and our KDFNet are the pixels that belong to the same masks. We test two types

of masks: 1) ground truth object mask, and 2) ground truth object mask applied with additional small circular occlusions at keypoint locations. For evaluation metrics, we measure average 2D projection error which is described in more detail in Section V-B. The evaluation threshold is 1 pixel. The results are illustrated in Table I. Our model is not affected by occlusions on pixel voters. However, the performance of PVNet drops significantly.

We provide an explanation as follows. The geometry of the medical swab stick is long and thin, in which case most of the keypoint voting directions are close to being on the same line except for the small area around the keypoint. As illustrated in Figure 3, intersections cannot be robustly found from the voting directions if the small area near the keypoints is occluded. On the contrary, our model is based on distance voting and can still reliably locate the keypoints from circle intersection under occlusion. Through this extremely simple dataset, we show the drawbacks of direction-based keypoint voting methods and the advantage of distance-based voting adopted by our KDFNet.

## V. EXPERIMENTS

In this section, we describe the details of our experiment settings in terms of the dataset, implementation details, and evaluation metrics. We then present the experiment results including ablation studies. Additionally, we provide visualizations on KDF and predicted 6D poses of objects.

### A. Dataset and Implementation

**Occlusion LINEMOD [8]** is a standard benchmark for 6D object pose estimation. It contains videos of desktop objects in a cluttered scene. It is a subset of the LINEMOD dataset [25] that mainly focuses on objects under occlusion. Together with the annotated images, high-quality 3D scanned models of the objects are also provided. The test set of Occlusion LINEMOD consists of an image sequence of 1,214 frames, each annotated with the 6D poses of 8 objects from LINEMOD dataset. The data used to train our model are the same as [5]: real images from LINEMOD and synthetically rendered images using the scanned 3D object models. We adopt the same object keypoint set as [5], which are generated by Farthest Point Sampling (FPS) of the object 3D point set.

**TOD Dataset [4]** is a dataset for keypoint estimation of transparent objects. It consists of 48,000 stereo images from 600 stereo videos of 15 transparent objects placed on simple textured tabletops. Every stereo image is annotated with 3D keypoints. Together with the images, high-quality 3D scanned models of the objects are also provided. In TOD dataset, there are 2,880 training images and 320 test images for each object. Since the objects are transparent, we did not render additional synthetic images during training but adopt geometric and photometric data augmentations during training. Though TOD dataset was originally created for keypoint estimation only, it can still be used to train and evaluate the 6D poses of the objects. To this end, we select a subset of the objects in TOD that have at least three keypoints

defined, i.e. the 7 mug categories, so that there are enough keypoints for recovering 6D poses.

### B. Evaluation Metrics

**ADD(-S) Metrics.** We use ADD [25] and ADD-S [2] in our evaluation. When computing ADD distance, we transform the model point set by the predicted and the ground truth poses respectively, and compute the mean 3D Euclidean distance between the two point sets. Given an object with 3D model point set of $\mathcal{M} = \{\mathbf{x}_i \in \mathbb{R}^3 \mid i = 1, 2, \ldots, M\}$, the ADD distance is calculated as

$$\text{ADD} = \frac{1}{|\mathcal{M}|} \sum_{\mathbf{x} \in \mathcal{M}} ||(\mathbf{R}\mathbf{x} + \mathbf{T}) - (\mathbf{R}^*\mathbf{x} + \mathbf{T}^*)|| \quad (8)$$

where $\mathbf{R}^*$ and $\mathbf{R}$ are the ground truth and estimated rotation, and $\mathbf{T}^*$ and $\mathbf{T}$ are the ground truth and estimated translation. For symmetric objects, ADD-S [2] is used instead. When computing ADD-S distance, the 3D distances are calculated between the closest points

$$\text{ADD-S} = \frac{1}{|\mathcal{M}|} \sum_{\mathbf{x}_1 \in \mathcal{M}} \min_{\mathbf{x}_2 \in \mathcal{M}} ||(\mathbf{R}\mathbf{x}_1 + \mathbf{T}) - (\mathbf{R}^*\mathbf{x}_2 + \mathbf{T}^*)||$$
$$(9)$$

We use the following two evaluation metrics. (1) **ADD(-S) accuracy**: ADD(-S) accuracy measures the proportion of correct pose predictions. A pose prediction is considered correct if the ADD(-S) distance is less than the threshold of 10% of the model's diameter. (2) **ADD(-S) AUC**: the area under ADD(-S) accuracy-threshold curve where the maximum threshold is set to 10cm. To compare with previous works, we use ADD(-S) accuracy on Occlusion LINEMOD datasets, and ADD accuracy and AUC on TOD dataset.

**2D Projection Metric.** When computing 2D projection error, we transform the model point set by the predicted and the ground truth poses respectively, and compute the mean 2D distance between the image projections of model points. Given the camera projection function $\Pi : \mathbb{R}^3 \to \mathbb{R}^2$, the 2D projection error is calculated as

$$\text{2D-Proj} = \frac{1}{|\mathcal{M}|} \sum_{\mathbf{x} \in \mathcal{M}} ||\Pi(\mathbf{R}\mathbf{x} + \mathbf{T}) - \Pi(\mathbf{R}^*\mathbf{x} + \mathbf{T}^*)|| \quad (10)$$

A pose prediction is considered as correct if the distance is less than the threshold of 5 pixels. We use the **2D projection accuracy** to evaluate Occlusion LINEMOD dataset.

### C. Experiment Results and Comparison

We train the KDFNet model on the LINEMOD dataset and rendered synthetic images. We evaluate and compare our model against previous baselines [2, 5, 6, 17, 24] on Occlusion LINEMOD dataset. The results are illustrated in Table II. Among these baselines, the most relevant is PVNet [5], a direction-based keypoint voting method which was also compared against in the toy experiment in Section IV. Our method achieves the best average ADD(-S) accuracy of 50.3% among all baselines while being the best on 5 of the 8 objects. In terms of 2D projection accuracy, our model also achieves the best among all baselines with an average of 66.5% while being the best on 6 of the 8 objects. In

TABLE II: **ADD(-S) accuracy** (left) and **2D projection error accuracy** (right) results on Occlusion LINEMOD [8].

| methods | PoseCNN [2] | Oberweger [24] | Pix2Pose [17] | PVNet [5] | HybridPose [6] | KDFNet (ours) |
|---|---|---|---|---|---|---|
| ape | 9.6 | 17.6 | **22.0** | 15.8 | 20.9 | 19.5 |
| can | 45.2 | 53.9 | 44.7 | 63.3 | 75.3 | **78.4** |
| cat | 0.9 | 3.31 | 22.7 | 16.7 | 24.9 | **28.2** |
| driller | 41.4 | 45.2 | 44.7 | 65.7 | 70.2 | **75.1** |
| duck | 19.6 | 17.2 | 15.0 | 25.2 | 27.9 | **38.6** |
| eggbox | 22.0 | 25.9 | 25.2 | 50.1 | **52.4** | 51.2 |
| glue | 38.5 | 39.6 | 32.4 | 49.6 | **53.8** | 52.1 |
| holepuncher | 22.1 | 21.3 | 49.5 | 12.4 | 54.2 | **59.0** |
| average | 24.9 | 30.4 | 32.0 | 40.8 | 47.5 | **50.3** |

| methods | PoseCNN [2] | Oberweger [24] | PVNet [5] | KDFNet (ours) |
|---|---|---|---|---|
| ape | 34.6 | **69.6** | 69.1 | 66.6 |
| can | 15.1 | 82.6 | 86.1 | **91.1** |
| cat | 10.4 | 65.1 | 65.1 | **72.5** |
| driller | 31.8 | 73.8 | 73.1 | **79.7** |
| duck | 7.4 | 61.4 | 61.4 | **71.3** |
| eggbox | 1.9 | **13.1** | 8.4 | 6.1 |
| glue | 13.8 | 54.9 | 55.4 | **59.6** |
| holepuncher | 23.1 | 66.4 | 69.8 | **85.3** |
| average | 17.2 | 60.9 | 61.1 | **66.5** |

TABLE III: Ablation Study on the number of keypoints $K$. Evaluation metrics are **ADD(-S) accuracy** and **2D projection**.

| metrics | ADD(-S) accuracy | | | | 2D projection | | | |
|---|---|---|---|---|---|---|---|---|
| # of kps $K$ | 8 | 12 | 16 | 20 | 8 | 12 | 16 | 20 |
| ape | 18.5 | 19.0 | 19.3 | **19.5** | **67.3** | 67.0 | 67.1 | 66.6 |
| can | 77.7 | **80.0** | 79.4 | 78.4 | 90.6 | **91.4** | 91.2 | 91.1 |
| cat | 26.9 | 27.2 | **28.3** | 28.2 | **74.9** | 73.8 | 73.1 | 72.5 |
| driller | 72.3 | 74.6 | **75.5** | 75.1 | 77.0 | 78.4 | 79.4 | **79.7** |
| duck | 36.9 | 37.4 | 38.0 | **38.7** | **72.1** | 71.8 | 71.8 | 71.3 |
| eggbox | 50.6 | 50.5 | 51.7 | 51.3 | 6.1 | 5.7 | 5.6 | 6.1 |
| glue | 49.3 | 50.3 | 51.8 | **52.1** | 57.7 | 59 | 59.3 | **59.6** |
| holepuncher | 57.3 | 58.5 | 58.8 | **59.0** | 85.5 | 85.4 | **85.5** | 85.3 |
| average | 48.7 | 49.7 | 50.3 | 50.3 | 66.4 | 66.5 | **66.6** | 66.5 |

TABLE IV: Ablation Study on the threshold value. Evaluation metrics are **ADD(-S) accuracy** and **2D projection**.

| metrics | ADD(-S) accuracy | | | | 2D projection | | | |
|---|---|---|---|---|---|---|---|---|
| threshold $\theta$ | 0.2 | 0.4 | 0.8 | 1.6 | 0.2 | 0.4 | 0.8 | 1.6 |
| ape | 18.6 | 19.5 | 20.0 | **20.2** | 66.6 | 66.6 | **66.9** | 66.6 |
| can | **79.0** | 78.4 | 77.1 | 77.2 | 91.1 | 91.1 | 91.0 | **91.3** |
| cat | **28.6** | 28.2 | 26.4 | 25.3 | 72.6 | 72.5 | **72.9** | 72.9 |
| driller | **75.9** | 75.1 | **75.9** | 75.0 | **80.2** | 79.7 | 79.5 | 79.0 |
| duck | 37.7 | 38.0 | 38.7 | **39.4** | 71.6 | 71.3 | **72.1** | 72.0 |
| eggbox | 49.5 | 51.3 | 52.5 | **58.0** | **6.3** | 6.1 | 6.1 | 6.2 |
| glue | 51.4 | **52.1** | 50.8 | 50.1 | **59.6** | **59.6** | 59.2 | 58.8 |
| holepuncher | **59.1** | 59.0 | 58.5 | 57.1 | **85.3** | **85.3** | 84.9 | **85.3** |
| average | 50.0 | **50.3** | 50.1 | 50.1 | **66.6** | 66.5 | **66.6** | 66.5 |

TABLE V: Ablation Study on the number of keypoint hypotheses $3N$. Evaluation metrics are **ADD(-S) accuracy** and **2D projection**.

| metrics | ADD(-S) accuracy | | | | 2D projection | | | |
|---|---|---|---|---|---|---|---|---|
| # of hypotheses | 48 | 192 | 768 | 3072 | 48 | 192 | 768 | 3072 |
| ape | 20.0 | 19.9 | 19.3 | 19.5 | 66.9 | 66.8 | 66.7 | 66.6 |
| can | 78.3 | 78.8 | 78.7 | 78.4 | 91.3 | 91.2 | 91.1 | 91.1 |
| cat | 28.1 | 28.4 | 28.3 | 28.2 | 72.7 | 72.2 | 72.5 | 72.5 |
| driller | 75.2 | 75.5 | 75.7 | 75.1 | 79.9 | 79.1 | 79.5 | 79.7 |
| duck | 38.9 | 38.5 | 38.4 | 38.7 | 71.7 | 71.5 | 71.6 | 71.3 |
| eggbox | 50.2 | 51.0 | 51.2 | 51.3 | 6.1 | 5.9 | 5.9 | 6.1 |
| glue | 51.1 | 51.2 | 51.6 | 52.1 | 59.6 | 59.8 | 60.4 | 59.6 |
| holepuncher | 59.2 | 58.5 | 59.0 | 59.0 | 85.2 | 85.3 | 85.4 | 85.3 |
| average | 50.1 | 50.2 | 50.3 | 50.3 | 66.7 | 66.5 | 66.7 | 66.5 |

particular, our method outperforms PVNet [5] by a margin of 9.5% in ADD(-S) accuracy and 5.6% in 2D projection accuracy, while also outperforming previous state-of-the-art HybridPose [6] by 2.8% in ADD(-S) accuracy.

Additionally, we train our model on TOD dataset [4] to predict the object keypoints in both stereo images and compare it against KeyPose [4]. KeyPose predicts object keypoints using heatmaps in both stereo images and the 6D object poses are computed by keypoint triangulation from stereo and pose fitting by solve an Orthogonal Procrustes problem. To evaluate KDFNet on TOD dataset, we follow the same procedure in [4] to recover 6D object pose from the predicted 2D keypoints in both stereo images. The results are illustrated in Table VI. Our method achieves state-of-the-art performance and surpasses [4]. Note that our leading margin is not as significant as on Occlusion LINEMOD dataset. An explanation is that TOD dataset did not include occlusion, therefore the proposed KDFNet cannot fully show its ability in dealing with occlusion compared to a heatmap-based method.

*D. Ablation Studies*

We conduct ablation studies on Occlusion LINEMOD dataset in the following three aspects to help and validate the design choices of our architecture.

**Number of Keypoints $K$** defined for the object. Intuitively, increase the number of keypoints will provide more information for the 2D-3D correspondences and yield better 6D pose estimates. In the ablation study, the value of $K$

varies from 8 to 20. The ablation results are illustrated in Table III. For most objects, the performance generally increases as the number of keypoints increases. The performance saturates at $K = 16$. which means localizing additional keypoints does not further improve the estimation of 2D-3D correspondence.

**Pixel Threshold Value $\theta$** in Equation (7). We are interested in whether the pose estimation results are sensitive to the choice of $\theta$. The ablation results are illustrated in Table IV. Though the threshold value varies in a wide range from 0.2 to 1.6, both metrics of ADD(-S) and 2D projection accuracy stay almost the same. This means the choice of the hyperparameter $\theta$ does not exert a significant effect on the performance as long as it is within a reasonable range, which shows the stability with respect to $\theta$ of our model during inference.

**Number of Keypoint Hypothesis** $3N$. We are interested in the number of keypoint hypotheses that are sufficient for finding good keypoint estimates. The ablation results are illustrated in Table V. Starting at $3N = 48$, both metrics stay almost the same. This means with the good regression of KDF, a small number of keypoint hypotheses can cover keypoints that result in near-optimal predictions.

*E. Visualization*

We visualize the KDF, voted 2D locations of the keypoints, and the estimated 6D poses on Occlusion LINEMOD dataset in Figure 5. The KDFs are visualized as a heatmap superimposed on the RGB images. Our framework can accurately localize keypoints even under occlusion. The predicted KDF and the ground truth KDF are different at large keypoint distances because we followed the training loss strategy in III-C and only train the model on elements that are within the keypoint distance of 64 pixels. As illustrated in Figure 5(d)-(e), such difference does not affect the voting process and our framework can still localize keypoints on the image and estimate 6D poses accurately.

*F. Run Time Analysis.*

We test KDFNet on a machine with an Intel i7-6850K 3.7GHz CPU, a GTX 1080 Ti GPU and Tensorflow 1.15. Given the input image of resolution of $256 \times 256$ and 192 keypoint hypotheses, the network inference and voting take 116ms and 24ms respectively, resulting in KDFNet running at 7 frames/sec.

## VI. CONCLUSIONS

In this work, we propose a novel method named KDFNet for 6D pose estimation from RGB images. Our method is based on the novel representation of Keypoint Distance Field (KDF). We also proposed a distance-based voting scheme to recover the 2D locations of keypoints from predicted distance fields in a RANSAC fashion. Experiment results show that KDFNet achieves state-of-the-art performance on Occlusion LINEMOD and TOD dataset.

As future work, we will investigate the extension of the proposed idea of distance field and voting to robotic perception problems in other scenarios or modalities, such as object detection [26] in temporal [27] and/or 3D data [28, 29].

## VII. ACKNOWLEDGEMENT

TABLE VI: **TOD dataset** results. Evaluation metrics are **ADD accuracy** (%) and **AUC** with range of 0 to 10cm.

| method | KeyPose [4] | | **KDFNet (ours)** | |
|---|---|---|---|---|
| metrics | accuracy | AUC | accuracy | AUC |
| mug_0 | 70.57 | **90.06** | **71.55** | 89.72 |
| mug_1 | 48.09 | 81.73 | **53.43** | **85.28** |
| mug_2 | 67.72 | **88.95** | **73.42** | 86.82 |
| mug_3 | 72.50 | **88.59** | **76.88** | 88.14 |
| mug_4 | **75.00** | 85.69 | 74.69 | **86.21** |
| mug_5 | 91.48 | 91.14 | **92.43** | **91.81** |
| mug_6 | 87.66 | **89.83** | 87.66 | 89.73 |
| average | 73.29 | 88.00 | **75.72** | **88.24** |

## REFERENCES

[1] A. Kendall, M. Grimes, and R. Cipolla, "Posenet: A convolutional network for real-time 6-dof camera relocalization," in *ICCV*, 2015. 1, 2

[2] Y. Xiang, T. Schmidt, V. Narayanan, and D. Fox, "Posecnn: A convolutional neural network for 6d object pose estimation in cluttered scenes," *RSS*, 2018. 1, 2, 4, 5, 6

[3] G. Pavlakos, X. Zhou, A. Chan, K. G. Derpanis, and K. Daniilidis, "6-dof object pose from semantic keypoints," in *ICRA*, 2017. 1, 2, 3

[4] X. Liu, R. Jonschkowski, A. Angelova, and K. Konolige, "Keypose: Multi-view 3d labeling and keypoint estimation for transparent objects," in *CVPR*, 2020. 1, 2, 3, 5, 6, 7

[5] S. Peng, Y. Liu, Q. Huang, X. Zhou, and H. Bao, "Pvnet: Pixel-wise voting network for 6dof pose estimation," in *CVPR*, 2019. 1, 2, 4, 5, 6

[6] C. Song, J. Song, and Q. Huang, "Hybridpose: 6d object pose estimation under hybrid representations," in *CVPR*, 2020. 1, 2, 4, 5, 6

[7] X. Sun, B. Xiao, F. Wei, S. Liang, and Y. Wei, "Integral human pose regression," in *ECCV*, 2018. 1, 2

[8] E. Brachmann, A. Krull, F. Michel, S. Gumhold, J. Shotton, and C. Rother, "Learning 6d object pose estimation using 3d object coordinates," in *ECCV*, 2014. 2, 5, 6

[9] S. Suwajanakorn, N. Snavely, J. J. Tompson, and M. Norouzi, "Discovery of latent 3d keypoints via end-to-end geometric reasoning," in *NIPS*, 2018. 2

[10] 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 *CVPR*, 2020. 2

[11] J. N. Kundu, M. Rahul, A. Ganeshan, and R. V. Babu, "Object pose estimation from monocular image using multi-view keypoint correspondence," in *ECCV*, 2018. 2

[12] D. Mehta, H. Rhodin, D. Casas, P. Fua, O. Sotnychenko, W. Xu, and C. Theobalt, "Monocular 3d human pose estimation in the wild using improved cnn supervision," in *3DV*, 2017. 2

[13] K. Duan, S. Bai, L. Xie, H. Qi, Q. Huang, and Q. Tian, "Centernet: Keypoint triplets for object detection," in *ICCV*, 2019. 2

[14] H. Law and J. Deng, "Cornernet: Detecting objects as paired keypoints," in *ECCV*, 2018. 2

[15] 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 *ICCV*, 2017. 2

[16] S. Zakharov, I. Shugurov, and S. Ilic, "Dpod: 6d pose object detector and refiner," in *ICCV*, 2019. 2

[17] K. Park, T. Patten, and M. Vincze, "Pix2pose: Pixel-wise coordinate regression of objects for 6d pose estimation," in *ICCV*, 2019. 2, 5, 6

[18] C. R. Qi, O. Litany, K. He, and L. J. Guibas, "Deep hough voting for 3d object detection in point clouds," in *ICCV*, 2019. 2

[19] A. Velizhev, R. Shapovalov, and K. Schindler, "Implicit shape models for object detection in 3d point clouds," in *International Society of Photogrammetry and Remote Sensing Congress*, 2012.

[20] W. Kehl, F. Milletari, F. Tombari, S. Ilic, and N. Navab, "Deep learning of local rgb-d patches for 3d object detection and 6d pose estimation," in *ECCV*, 2016.

[21] M. Sun, G. Bradski, B.-X. Xu, and S. Savarese, "Depth-encoded hough voting for joint object detection and shape
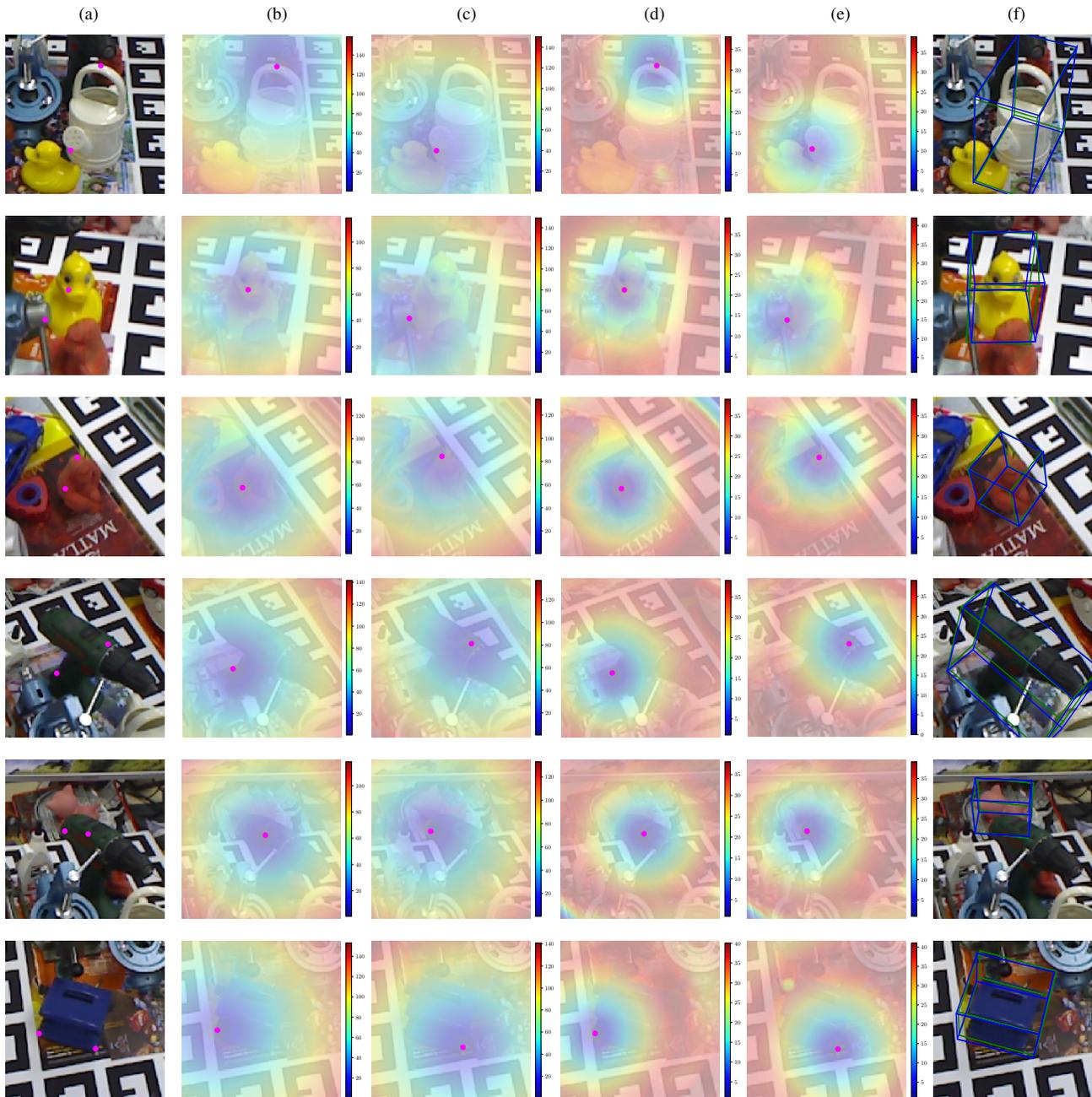
Fig. 5: **Visualization of predicted keypoints, distance fields and 6D poses** on Occlusion LINEMOD dataset. From left to right in each row: (a) input image with two of the ground truth keypoints; (b)(c) two ground truth distance fields and locations of the keypoints; (d)(e) two predicted distance fields and voted location of the keypoints; (f) ground truth 6D pose (green 3D bounding box) and predicted 6D pose (blue 3D bounding box).

recovery," in *ECCV*, 2010. 2

[22] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *NIPS*, 2015. 3

[23] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016. 4

[24] M. Oberweger, M. Rad, and V. Lepetit, "Making deep heatmaps robust to partial occlusions for 3d object pose estimation," in *ECCV*, 2018. 5, 6

[25] S. Hinterstoisser, S. Holzer, C. Cagniart, S. Ilic, K. Konolige, N. Navab, and V. Lepetit, "Multimodal templates for real-time detection of texture-less objects in heavily cluttered scenes," in *ICCV*, 2011. 5

[26] C. R. Qi, Y. Zhou, M. Najibi, P. Sun, K. Vo, B. Deng, and D. Anguelov, "Offboard 3d object detection from point cloud sequences," in *CVPR*, 2021. 7

[27] X. Liu, J.-Y. Lee, and H. Jin, "Learning video representations from correspondence proposals," in *CVPR*, 2019. 7

[28] X. Liu, M. Yan, and J. Bohg, "Meteornet: Deep learning on dynamic 3d point cloud sequences," in *ICCV*, 2019. 7

[29] X. Liu, C. R. Qi, and L. J. Guibas, "Flownet3d: Learning scene flow in 3d point clouds," in *CVPR*, 2019. 7