

Single Shot 6D Object Pose Estimation

Kilian Kleeberger¹ and Marco F. Huber^{1,2}

Abstract—In this paper, we introduce a novel single shot approach for 6D object pose estimation of rigid objects based on depth images. For this purpose, a fully convolutional neural network is employed, where the 3D input data is spatially discretized and pose estimation is considered as a regression task that is solved locally on the resulting volume elements. With 65 fps on a GPU, our Object Pose Network (OP-Net) is extremely fast, is optimized end-to-end, and estimates the 6D pose of multiple objects in the image simultaneously. Our approach does not require manually 6D pose-annotated real-world datasets and transfers to the real world, although being entirely trained on synthetic data. The proposed method is evaluated on public benchmark datasets, where we can demonstrate that state-of-the-art methods are significantly outperformed.

I. INTRODUCTION

Knowing the pose of objects is a crucial prerequisite for many robotic grasping and manipulation tasks. 6D object pose estimation (OPE) is a long-standing challenge and an open field of research since the early days of computer vision. The pose of an object is fully described by a translation vector $\mathbf{t} \in \mathbb{R}^3$ and a rotation matrix $\mathbf{R} \in \text{SO}(3)$ of its body-fixed coordinate system relative to a given reference frame. The task of 6D OPE is challenging due to the variety of objects in the real world, sensor noise, clutter and occlusion in the scene and varying lighting conditions, which affect the appearance of the objects in the image. Object symmetries, which are common in man-made and industrial environments, result in pose ambiguities and have to be addressed.

In recent years, research in 6D OPE has been dominated by convolutional neural network (CNN) based approaches. The state-of-the-art approaches solve this as a classification problem, in which the pose space is discretized into bins and a CNN is used to predict a pose bin [1], [2], [3], [4], [5], [6], [7]. The pose space, however, is continuous and thus, we solve the 6D OPE problem by means of a regression task. Recent works regress the 2D image coordinates of the object’s 3D bounding box and use a PnP algorithm to estimate the object’s 6D pose [8], [9], [10]. Because of the more intuitive interpretation, this task seems easier to learn than directly predicting the pose. To the best of our knowledge, we are the first directly predicting the 6D object



Fig. 1. Object pose estimates (black) of our approach on real-world data for ring screws after ICP refinement. Although our models are fully trained on synthetic data, they successfully transfer to the real world.

pose in terms of position and angles for the orientation as a regression task in a single shot framework, without requiring any computational overhead like PnP , clustering, or post-processing steps like the iterative closest point (ICP) algorithm for pose refinement. We propose a solution for dealing with discrete and revolution object symmetries, which need to be considered to avoid inconsistent loss signals during training. Contrary to most other methods using RGB or RGB-D images as input, we use depth data only. Our approach is inspired by [11], [12] and uses a fully convolutional architecture to process the depth image.

Approaches for 6D OPE relying on a prior semantic segmentation step [8], [13], [14], [15], [10] cannot be used for scenes of many parts of the same type in bulk. When the parts overlap, they require instance segmentation, which is hard to obtain for highly occluded and cluttered scenes. Additionally, the quality of the segmentation masks also influences the quality of the pose estimates. As shown in [16], [17] for 2D object detection, a single shot system reasons globally on the image and the whole pipeline is optimized end-to-end, thus being both faster and more accurate compared to multi-stage systems [18], [19].

So far, direct pose regression methods have been less successful than classification [1], [20], [6]. A major reason for this is the lack of precisely annotated training data without data redundancy. Usual datasets show a limited variability of poses relative to the camera due to numerous image acquisitions of the same scene from different views and are therefore less suitable for precise regression tasks [21], [22], [23], [24], [25], [26]. Because creating and

¹Kilian Kleeberger is with the Department Robot and Assistive Systems and Marco F. Huber is with the Center for Cyber Cognitive Intelligence (CCI), Fraunhofer Institute for Manufacturing Engineering and Automation IPA, Nobelstraße 12, 70569 Stuttgart, Germany {kilian.kleeberger, marco.huber}@ipa.fraunhofer.de

²Marco F. Huber is with the Institute of Industrial Manufacturing and Management IFF, University of Stuttgart, Allmandring 35, 70569 Stuttgart, Germany marco.huber@ieee.org

annotating datasets with 6D poses is time-consuming and does not scale, since the process has to be repeated for every new application, we tackle this problem via simulation. This enables generating large-scale datasets with flawless annotations. We use sim-to-real transfer techniques to allow the model generalizing in the real world.

We benchmark our proposed method on the Siléane [27] and Fraunhofer IPA [28] datasets, which are challenging due to multiple instances of the same object type and a high amount of clutter and occlusion. By using our data generator to generate training data, we significantly outperform the state-of-the-art methods without requiring additional post-processing steps like ICP or duplicate removal.

In summary, the main contributions of this work are:

- A novel single shot approach for 6D OPE that is highly robust to occlusions between objects and does not rely on any post-processing to get accurate results, even on low resolution input images (128×128 pixel)
- Two novel loss functions that can properly deal with object symmetries and formulate the 6D OPE task as a regression instead of classification problem
- A framework for training the pose estimator entirely in simulation and enabling it to generalize in the real world without requiring 6D pose-annotated real-world training data

The paper is structured as follows. In the next section, related work is reviewed. In Section III the proposed approach is described. Experimental evaluations are provided in Section IV. Pros and cons of our approach are discussed in Section V. The paper closes with a conclusion.

II. RELATED WORK

We first review related work on classical feature and template matching methods before taking a closer look at newer CNN-based methods for 6D OPE using RGB, depth, or RGB-D images as input.

A. Classical Approaches

Traditionally, the problem of 6D OPE is tackled by template or feature matching. In template-based methods [25], [29], rigid templates are constructed by rendering the 3D object model from different viewpoints, which are then matched into the sensor data at different locations. A similarity score is used for the evaluation of the pose hypothesis. These methods are useful for detecting texture-less objects, but occlusions significantly reduce the performance. This is due to the low similarity between the template and the sensor data if the object is occluded. In feature-based methods, local features from the image are matched to the 3D object model to recover the 6D pose based on the spatial relationship [30], [31], [4]. These methods are designed to handle changes in size, viewpoint and illumination. While being robust to occlusion and scene clutter, they can only reliably handle objects with sufficient texture.

B. CNN-based Approaches

In recent years, research in 6D OPE has been dominated by CNN-based approaches. In [2], a binning of the angles and the height for the 6D OPE in depth images is introduced. The accuracy and speed are limited due to the discretization and the multi-stage nature of the system. [20] solves pose estimation as a regression instead of classification task, but focuses on the retrieval of the orientation only based on the object’s bounding box.

Deep-6DPose [32] extends Mask R-CNN [19] with an additional branch for estimating the full 6D pose based on the given region proposal. After a coarse segmentation, BB8 [8] predicts the 2D projections of the corners of the object’s 3D bounding box. The 6D pose is estimated by using a PnP algorithm and a final CNN per object is trained to refine the pose. Due to the usage of multiple separated CNNs, the approach is not optimized end-to-end and is time-consuming for the inference. PoseCNN [13] decouples the translation and rotation estimation based on a prior semantic segmentation step and thus, also requires multiple stages.

SSD-6D [1] extends the SSD [17] detection framework to 6D OPE by predicting the 2D bounding boxes with SSD, hypothesizing 6D poses from the network output (6D pose pool) and running a refinement step (ICP). The rotation estimation is treated as a classification problem by decomposing the 3D rotation space into discrete viewpoints and in-plane rotations.

PPR-Net [33], the winner method of the “*Object Pose Estimation Challenge for Bin-Picking*” at IROS 2019¹, uses PointNet++ [34] and estimates a 6D pose for each point in the point cloud of the object instance to which it belongs. Afterwards, density-based clustering in 6D space is applied and the final pose hypotheses are obtained by averaging the predicted poses for each identified cluster.

A recent result from the Benchmark for 6D Object Pose Estimation (BOP) [35] is that methods based on point-pair features currently perform best on various datasets, outperforming methods based on template matching, 3D local features, and machine learning. This is in contrast to many other computer vision tasks such as image classification [36], [37], object detection [11], [16], [17], semantic segmentation [38], instance segmentation [19], etc. being dominated by deep learning. With this work, we advance the state-of-the-art for learning-based methods by providing a simple pipeline for an end-to-end trainable general purpose 6D pose estimator, which can be trained entirely on synthetic data and generalizes in the real world.

III. OBJECT POSE NETWORK

In this section, we describe the process of data generation for training our neural network, the parameterization of the network’s output, the loss function, the network architecture, the training procedure, and the technique for the transfer of the model from simulation to the real world. Fig. 2 shows an overview of our approach.

¹<http://www.bin-picking.ai/en/competition.html>

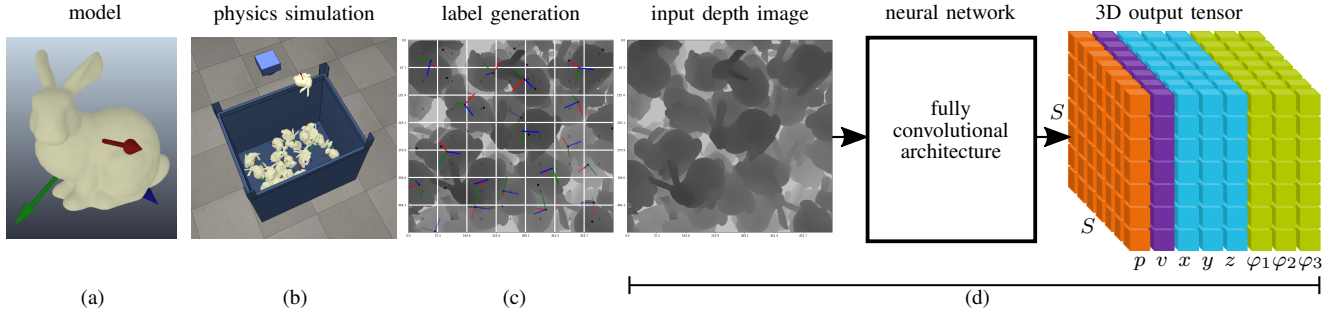


Fig. 2. Overview: (a) 3D object model with body-fixed coordinate system (can be a CAD or previously scanned model). (b) Physics simulation for the generation of training data. (c) Rendered depth image from the simulation with overlaid $S \times S$ grid. Based on the 3D position of the object coordinate system, the objects in the scene are assigned to the volume elements. (d, left) Rendered depth image used as input for the neural network. (d, middle) The proposed CNN architecture based on a DenseNet-BC [37]. (d, right) 3D output tensor of the neural network, where each spatial location consists of a vector comprising the probability p that the cell contains an origin, visibility v , positions x, y, z , and angles $\varphi_1, \varphi_2, \varphi_3$.

A. Data Generation

In order to generate the training data for OP-Net, we use a physics simulation that produces scenes of relevant scenarios as depicted in Fig. 2 (b). The depth images and segmentation masks are generated in perspective projection. The ground truth annotation consists of a class label, the 6D pose, i.e., translation vector \mathbf{t} and rotation matrix \mathbf{R} relative to the coordinate system of the 3D sensor, a segmentation label, and a visibility score of each object instance in the scene. The visibility $v \in [0, 1]$ is the ratio between the number of visible pixels and the number of pixels without any occlusion.

B. Parameterization of the Output

The output of our model is inspired by the parameterization proposed in [12], [11]. Instead of predicting (oriented) rectangles in 2D, we estimate the 6D pose of the objects in the image. By discretizing the 3D scene in $S \times S$ volume elements, the 6D pose regression task is solved locally, i.e., individually for each volume element. Each volume element comprises an 8-dimensional vector containing the probability p , visibility v , positions x, y, z , and angles $\varphi_1, \varphi_2, \varphi_3$ in a given convention. For the ground truth generation, the pose of an object is assigned to a volume element if the origin of the body-fixed coordinate system of the 3D object model is located within the volume element as visualized in Fig. 2 (c). If multiple origins fall into the same volume element, the object with higher visibility v is assigned as ground truth. All volume elements not containing an object are filled with a zero vector. The output of the network is a $S \times S \times 8$ tensor as depicted in Fig. 2 (d, right).

The probability channel reflects how confident the model is that the volume element contains an origin of a 3D object model. Both the probability and visibility prediction at each spatial location are used to filter the results at test time.

The parameterization can naturally be extended to class probabilities of the objects. For instance, a variant of YOLO [16] can differentiate between 9,000 object classes alongside with confidence scores and bounding box coordinates. Also an estimate of how well an object can be grasped,

can naturally be added. Thus, the parameterization is general and can be used for any rigid 3D object.

C. Loss Functions

To train the network, the multi-task loss function

$$\mathcal{L} = \sum_{i=1}^{S^2} \left(\lambda_1 (p_i - \hat{p}_i)^2 + \left[\lambda_2 (v_i - \hat{v}_i)^2 + \lambda_3 (\mathcal{L}_{\text{pos}} + \lambda_4 \mathcal{L}_{\text{ori}}) \right] p_i \right) \quad (1)$$

is optimized, where a hat indicates estimates of the network. $\lambda_1, \lambda_2, \lambda_3$, and λ_4 are manually tuned weights for the different loss terms. While $\lambda_1 = 0.1$, $\lambda_2 = 0.25$, and λ_4 are constant, $\lambda_3 = 8v^3$ is a function of the ground truth visibility v . The less an object is visible, the lower the loss of the pose error is weighted. The idea is to let the network focus on the well visible objects and alleviate convergence issues due to pose ambiguities because of occlusions.

To stabilize the training, the position $\mathbf{x} = [x, y, z]^T$ of the object is estimated relative to the volume element, i.e. $x, y, z \in (0, 1)$, while z is the position between the near and far clipping plane of the 3D sensor. We use

$$\mathcal{L}_{\text{pos}} = \|\mathbf{x} - \hat{\mathbf{x}}\|^2, \quad (2)$$

with $\|\cdot\|$ being the L^2 norm.

We consider two loss functions for the estimation of the orientation. First,

$$\mathcal{L}_{\text{ori1}} = \|\boldsymbol{\varphi} - \hat{\boldsymbol{\varphi}}\|^2, \quad \lambda_4 = 1, \quad (3)$$

with $\boldsymbol{\varphi} = [\varphi_1, \varphi_2, \varphi_3]^T$. Second,

$$\mathcal{L}_{\text{ori2}} = \min_{\mathbf{pR} \in \mathcal{R}_{\mathbf{R}}(\mathcal{P})} \|\mathbf{pR} - \hat{\mathbf{pR}}\|, \quad \lambda_4 = 0.5, \quad (4)$$

with \mathbf{pR} being the orientation part of the pose representative $\mathbf{p} \in \mathcal{R}$ of pose \mathcal{P} and comprising the relevant axis vectors of the rotation matrix \mathbf{R} depending on the proper symmetry group of the object, based on [39], [27]. This loss function computes the Euclidean distance between the prediction and the closest ground truth. Since it gives better results, we bound the angles for both loss functions, i.e., $\varphi_1, \varphi_2 \in [0, 2\pi)$ and $\varphi_3 \in [0, 2\pi/k)$ are mapped to $[0, 1)$ where $k \in \mathbb{N}$

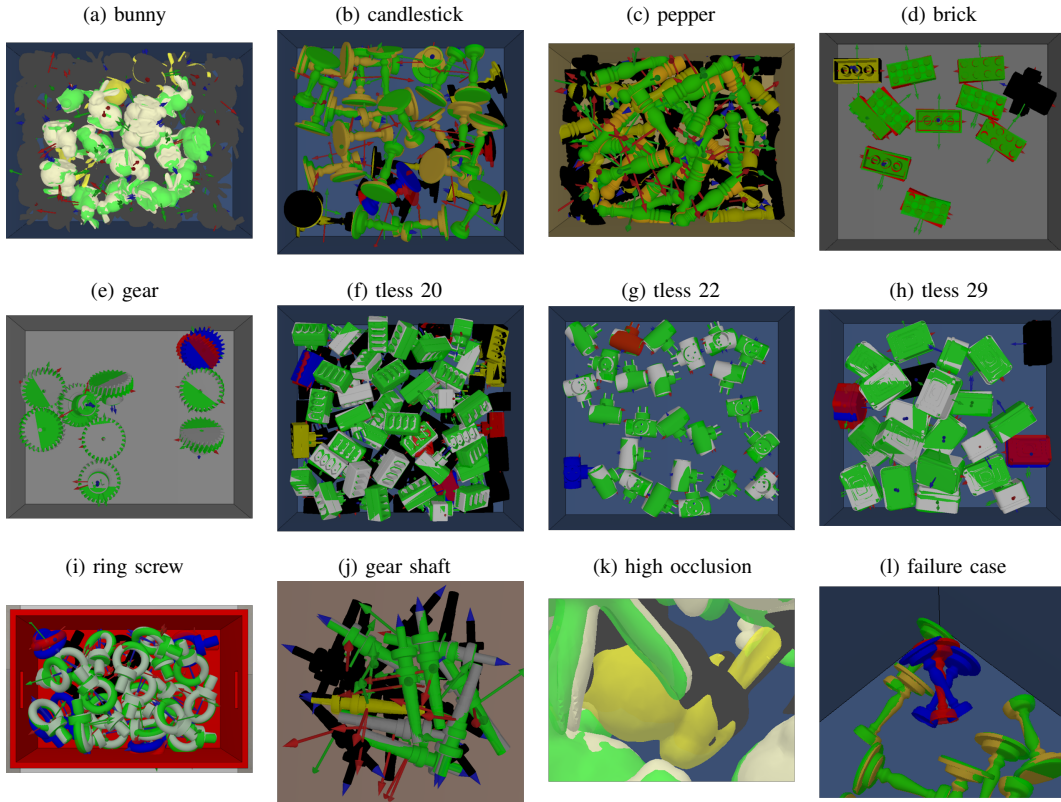


Fig. 3. Some qualitative results on scenes from the Siléane [27] (a–h) and Fraunhofer IPA [28] (i, j) datasets without any post-processing, a successful retrieval of an object with $v = 0.2$ (k), and a failure case due to a limitation of the parameterization (l). All images show a 3D visualization of the ground truth and the result of our method in a simulation. Object instances that do not need to be found ($v \leq 0.5$) are visualized in black, while especially some objects at the border of the 3D visualization are black because they are truncated in the depth image. Objects that need to be found and whose pose was correctly estimated are visualized in their original color. Otherwise they are visualized in blue. Correctly found objects are visualized in green (true positives). Red indicates that the object is either a duplicate or too far off from the ground truth pose (false positive). Objects that were correctly found, but their pose was not of interest for retrieval (too much occluded) are displayed in yellow and are neither considered as true nor false positive.

represents the order of the cyclic symmetry. For objects with a revolution symmetry, we omit the output feature-map for φ_3 . In Section IV-B we compare the performance of the loss functions for the regression of the angles.

D. Network Architecture

The input of our model is a single normalized depth image in perspective projection. The data is processed with a fully convolutional architecture and mapped to a 3D output tensor as shown in Fig. 2 (d). In the experiments, we use an input resolution of 128×128 pixel, a DenseNet-BC [37] with 40 layers and a growth rate of 50, which represents the number of feature-maps being added per layer, and $S = 16$. We choose a DenseNet-BC because it promotes gradient propagation by introducing direct connections between any two layers with the same feature-map size and has a high parameter efficiency. The network architecture consists of four dense blocks and downsampling is performed three times via 2×2 average pooling to reduce the size of the feature-maps from 128×128 to 16×16 and preserve the spatial information. ReLU activation functions are employed in the dense blocks and sigmoid functions for the 3D output

tensor. With this architecture, forward passes are performed with an average frame rate of 65 fps on a Nvidia Tesla V100.

E. Training

During training, the error for the entire probability channel and the error of the visibility and 6D pose for the specific entries in the 3D output tensor that contain ground truth poses are backpropagated. Thanks to the utilized simulation, there is an abundant source of data that allows training the network from scratch. We use the Adam optimizer with an initial learning rate of 0.01, monitor the validation loss, reduce the learning rate by a factor of 10 if the loss did not improve for three epochs, and train the network for about 50 epochs on the data generated by the physics simulation.

F. Sim-to-Real Transfer

To transfer the model from simulation to the real world, we use a technique called domain randomization [40]. We randomize various aspects of the simulation, e.g., pose and size of distractor objects that must be ignored by the network, and apply different augmentations with varying intensity to the rendered training images, e.g., adding noise, blurring, elastic transformations, dropout, etc. This allows the 6D pose estimator generalizing in the real world, although being

trained entirely on synthetic data. The main advantage of this technique is that it requires no samples from the real world. This is in contrast to domain adaptation [41], [42] techniques, which are usually based on Generative Adversarial Networks (GANs) [43]. Apart from being hard to train and often yielding fragile training results, GANs require the acquisition of (unlabeled) real-world samples for new applications, which negatively impacts scalability.

IV. EXPERIMENTAL EVALUATION

We evaluate our single shot approach for 6D OPE on bin-picking scenarios. These consist of multiple rigid objects of the same type, which are stored chaotically in a bin. 6D OPE is challenging due to the cluttered scene with multiple and heavy occlusions. We benchmark the performance of our method on the noisy data from the Siléane dataset [27] and the real-world data from the Fraunhofer IPA dataset [28]. These show highly chaotic scenes typical for bin-picking. Fig. 1 and Fig. 3 depict qualitative results of OP-Net.

A. Evaluation Metric

A common evaluation metric for 6D OPE is ADD by Hinterstoisser et al. [25], which accepts a pose hypothesis if the average distance of model points between the ground truth and estimated pose is less than 0.1 times the diameter of the smallest bounding sphere of the object. As this metric cannot handle symmetric objects, ADI [25] was introduced for handling those. The ADI metric is widely used, but can fail to reject false positives as demonstrated in [27]. For evaluation, we use the metric provided by Brégier et al. [39], [27], which is suitable for rigid objects, for scenes of many parts in bulk, and properly considers cyclic and revolution object symmetries. A pose representative $\mathbf{p} \in \mathcal{R}(\mathcal{P})$ comprises a translation vector \mathbf{t} and the relevant axis vectors of the rotation matrix \mathbf{R} depending on the object’s proper symmetry group. The distance between a pair of poses \mathcal{P}_1 and \mathcal{P}_2 is defined as the minimum of the Euclidean distance between their respective pose representatives:

$$d(\mathcal{P}_1, \mathcal{P}_2) = \min_{\mathbf{p}_1 \in \mathcal{R}(\mathcal{P}_1), \mathbf{p}_2 \in \mathcal{R}(\mathcal{P}_2)} \|\mathbf{p}_1 - \mathbf{p}_2\| \quad (5)$$

A pose hypothesis is accepted (considered as true positive) if the minimum distance to the ground truth is less than 0.1 times the object’s diameter. Following [27], only the pose of objects that are less than 50% occluded are relevant for the retrieval. The metric breaks down the performance of a method to a single scalar value named average precision (AP) by taking the area under the precision-recall curve.

B. Comparison

Depending on the object, the Siléane dataset [27] comprises 46 to 325 images and is too small for training deep neural networks. Therefore, we precisely rebuild the setup from the dataset in simulation.

To allow OP-Net generalizing on real-world data without being trained on real-world samples, we randomize the bin pose in the simulation during the data generation process and

augment the synthetic images during training. As the real-world 3D sensor operates in perspective mode, the obtained point clouds are projected to perspective depth images because it shows fewer pixels with missing depth information compared to the orthogonal projection. The missing values are interpolated before the image is fed into the neural network.

Table I shows the performance of our method against the state-of-the-art in terms of the metric by Brégier et al. [39], [27]. We provide performance results for different loss functions for the regression of the angles. OP-Net outperforms other classical [44], [25] and learning-based [2], [33] approaches without requiring post-processing steps (PP) like duplicate removal or ICP, even on low resolution input images. Applying randomizations in the simulation and augmentations to the training data allows the model generalizing on real-world data without ever being trained on samples from the target domain. Additional PP steps like duplicate removal and ICP further improve the results.

V. DISCUSSION

In the following, we summarize strengths and discuss limitations of our approach.

A. Strengths

Due to learning plausible object pose configurations, OP-Net is highly robust to occlusions. Fig. 3 (k) exemplarily shows a successful estimate of an object that is barely visible in the input depth image.

For classical approaches and CNN-based methods relying on region proposals [32], [1], [6], [7], [20] or a prior segmentation step [8], [13], [14], [15] for pose estimation, the runtime increases linearly with the number of objects present in the scene, which is not optimal for very cluttered scenes. Our approach only requires a single feed-forward to estimate the 6D pose of multiple objects in the image simultaneously at 65 fps on a GPU. PPR-Net [33] requires 200 ms (GTX 1060) for the forward pass and additional clustering in 6D space plus pose averaging afterwards. OP-Net is faster because of using a more compact parameterization and not requiring any post-processing steps. Contrary to multi-stage methods, single shot approaches allow a global reasoning and still make accurate predictions locally. We frame 6D OPE as a regression problem and are therefore not limited in the accuracy given by the discretization of the pose space.

Our proposed loss functions are designed such that they can properly handle object symmetries. This avoids inconsistent loss signals during training causing convergence issues, because different orientations may generate identical observations. PoseCNN [13] uses the ADD and ADI metric [25] for symmetric objects as loss function, while ADI does not properly handle object symmetries as shown in [27]. Using the metric provided by Brégier et al. [39] as loss function, which properly handles object symmetries and compares the retrieved object pose against every possible ground truth, allows an efficient pose distance computation contrary to

TABLE I

AVERAGE PRECISION (AP) VALUES OF DIFFERENT METHODS AND COMPARISON OF THE PERFORMANCE OF DIFFERENT LOSS FUNCTIONS FOR OUR METHOD (BEST RESULTS IN BOLD). THE MODELS ARE TRAINED ON THE DATA GENERATED WITH OUR PHYSICS SIMULATION AND ALL IMAGES FROM THE SILÉANE [27] AND FRAUNHOFER IPA [28] DATASETS ARE USED FOR TESTING. RESULTS MARKED WITH * ARE TAKEN FROM THE “*Object Pose Estimation Challenge for Bin-Picking*” AT IROS 2019.

object	bunny [27] (no proper symmetry)	candlestick [27] (revolution)	pepper [27] (revolution)	brick [27] (cyclic, $k = 2$)	gear [27] (revolution)	tless 20 [27] (cyclic, $k = 2$)	tless 22 [27] (no proper symmetry)	tless 29 [27] (cyclic, $k = 2$)	ring screw [28] (cyclic, $k = 2$)	gear shaft [28] (revolution)
OP-Net (ours) with \mathcal{L}_{ori1}	0.92	0.94	0.98	0.41	0.82	0.85	0.77	0.51	0.88	0.99
OP-Net (ours) with \mathcal{L}_{ori1} and PP	0.94	0.97	0.98	0.42	0.84	0.88	0.86	0.58	0.93	0.99
OP-Net (ours) with \mathcal{L}_{ori2}	0.74	0.95	0.92	0.79	0.58	0.56	0.53	0.36	0.73	1.0
OP-Net (ours) with \mathcal{L}_{ori2} and PP	0.76	0.96	0.93	0.80	0.60	0.58	0.55	0.39	0.75	1.0
PPF [44], [27]	0.29	0.16	0.06	0.08	0.62	0.20	0.08	0.19	-	-
PPF PP [44], [27]	0.37	0.22	0.12	0.13	0.63	0.23	0.12	0.23	-	-
LINEMOD+ [25], [27]	0.39	0.38	0.04	0.31	0.44	0.25	0.19	0.20	-	-
LINEMOD+ PP [25], [27]	0.45	0.49	0.03	0.39	0.50	0.31	0.21	0.26	-	-
Sock et al. [2]	0.74	0.64	0.43	-	-	-	-	-	-	-
PPR-Net [33]	0.82	0.91	0.80	-	-	0.81	-	-	0.95*	0.99*
PPR-Net with ICP [33]	0.89	0.95	0.84	-	-	0.85	-	-	-	-

ADD and ADI, which depend on the sampling of the model points and require an iteration over all model points.

Manually creating and annotating a dataset with 6D poses is a very time-consuming process that does not scale because it has to be repeated for every new application. Therefore, we tackle the problem via simulation which allows to easily generate large-scale 6D pose-annotated datasets. We apply various augmentations to the synthetic images and show that our approach transfers to the real world, although being trained entirely on synthetic data. As we train our network on a bunch of synthetic variations that generalize in the real world, we are also independent of the 3D sensor technology being used. The Siléane [27] and Fraunhofer IPA [28] datasets were recorded with different sensors. In our experiments, we used the same network architecture, parameter configuration, and augmentation techniques which emphasises the scalability of our approach.

B. Limitations

The parameterization of the network output faces several limitations. If multiple objects fall into the same volume element, the object with higher visibility is assigned as ground truth. A failure case is exemplary visualized in Fig. 3 (l), where two candlesticks, which have their origin in the center of mass, fall into the same volume element with their origin. Since this scenario can hinder better convergence behaviour during training, a possible solution would be to predict multiple 6D poses per volume element or to introduce a discretization in z -direction resulting in a 3D probability map and therefore a 4D output tensor. A further limitation is that the origin of the 3D object model has to be inside the image in order to be detected by the proposed pose estimator. In case the estimation of the 6D pose of those objects is of relevance, one could add pixels at the border of the image to include the origin in the image. When the origin of the 3D object model is very close to the border of a volume element, the object might also get detected by neighboring volume elements. These can be filtered out via setting a higher threshold for \hat{p} or an additional duplicate

removal step based on the distance between the origins or pose representatives of the 6D object poses in the result.

VI. CONCLUSION AND FUTURE WORK

In this paper, we proposed a novel single shot approach for 6D object pose estimation of rigid objects based on depth images that is fast by design, end-to-end trainable, and facilitates direct 6D pose regression of the visible objects in the image simultaneously. Our experiments demonstrate that our method is highly robust to occlusions, can handle symmetric objects, and provides accurate pose estimates, even in highly cluttered scenes. OP-Net outperforms state-of-the-art methods on the Siléane dataset [27] by a large margin without any post-processing steps, even on low resolution input images. The proposed 6D pose estimator can be transferred to the real world, although being entirely trained on synthetic images and annotations by randomizing the simulation and augmenting the synthetic images. In doing so, pose estimation becomes independent of the 3D sensor technology being used in the real world. In future work, we will focus on giving a quality estimate for each predefined grasp pose on the 3D object model based on grasping trials in a physics simulation.

ACKNOWLEDGMENT

This work was partially supported by the Baden-Württemberg Stiftung gGmbH (Deep Grasping – Grant No. NEU016/1) and the Ministry of Economic Affairs of the state Baden-Württemberg (Center for Cyber Cognitive Intelligence (CCI) – Grant No. 017-192996). We would like to thank our colleagues for helpful discussions and comments.

REFERENCES

- [1] W. Kehl, F. Manhardt, F. Tombari, S. Ilic, and N. Navab, “SSD-6D: Making rgb-based 3d detection and 6d pose estimation great again,” in *IEEE International Conference on Computer Vision (ICCV)*, October 2017.
- [2] J. Sock, K. I. Kim, C. Sahin, and T.-K. Kim, “Multi-task deep networks for depth-based 6d object pose and joint registration in crowd scenarios,” in *British Machine Vision Conference (BMVC)*, September 2018.

- [3] S. Gupta, P. Arbelaez, R. Girshick, and J. Malik, "Aligning 3d models to rgb-d images of cluttered scenes," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.
- [4] S. Tulsiani and J. Malik, "Viewpoints and keypoints," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.
- [5] H. Su, C. R. Qi, Y. Li, and L. J. Guibas, "Render for CNN: Viewpoint estimation in images using cnns trained with rendered 3d model views," in *IEEE International Conference on Computer Vision (ICCV)*, December 2015.
- [6] M. Sundermeyer, Z.-C. Marton, M. Durner, M. Brucker, and R. Triebel, "Implicit 3d orientation learning for 6d object detection from rgb images," in *European Conference on Computer Vision (ECCV)*, September 2018.
- [7] X. Deng, A. Mousavian, Y. Xiang, F. Xia, T. Bretl, and D. Fox, "PoseRBPF: A rao-blackwellized particle filter for 6d object pose estimation," in *Robotics: Science and Systems (RSS)*, June 2019.
- [8] M. Rad and V. Lepetit, "BB8: A scalable, accurate, robust to partial occlusion method for predicting the 3d poses of challenging objects without using depth," in *IEEE International Conference on Computer Vision (ICCV)*, October 2017.
- [9] B. Tekin, S. N. Sinha, and P. Fua, "Real-time seamless single shot 6d object pose prediction," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [10] S. Peng, Y. Liu, Q. Huang, H. Bao, and X. Zhou, "PVNet: Pixel-wise voting network for 6dof pose estimation," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [11] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [12] J. Redmon and A. Angelova, "Real-time grasp detection using convolutional neural networks," in *IEEE International Conference on Robotics and Automation (ICRA)*, May 2015.
- [13] Y. Xiang, T. Schmidt, V. Narayanan, and D. Fox, "PoseCNN: A convolutional neural network for 6d object pose estimation in cluttered scenes," in *Robotics: Science and Systems (RSS)*, June 2018.
- [14] A. S. Periyasamy, M. Schwarz, and S. Behnke, "Robust 6d object pose estimation in cluttered scenes using semantic segmentation and pose regression networks," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, October 2018.
- [15] J. Wu, B. Zhou, R. Russell, V. Kee, S. Wagner, M. Hebert, A. Torralba, and D. M. S. Johnson, "Real-time object pose estimation with pose interpreter networks," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, October 2018.
- [16] J. Redmon and A. Farhadi, "YOLO9000: Better, faster, stronger," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [17] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "SSD: Single shot multibox detector," in *European Conference on Computer Vision (ECCV)*, October 2016.
- [18] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," in *Neural Information Processing Systems (NIPS)*, December 2015.
- [19] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask R-CNN," in *IEEE International Conference on Computer Vision (ICCV)*, October 2017.
- [20] S. Mahendran, H. Ali, and R. Vidal, "3d pose regression using convolutional neural networks," in *IEEE International Conference on Computer Vision (ICCV)*, October 2017.
- [21] T. Hodaň, P. Haluza, Š. Obdržálek, J. Matas, M. Lourakis, and X. Zabulis, "T-LESS: An rgb-d dataset for 6d pose estimation of texture-less objects," in *IEEE Winter Conference on Applications of Computer Vision (WACV)*, March 2017.
- [22] A. Tejani, D. Tang, R. Kouskouridas, and T.-K. Kim, "Latent-class hough forests for 3d object detection and pose estimation," in *European Conference on Computer Vision (ECCV)*, September 2014.
- [23] A. Doumanoglou, R. Kouskouridas, S. Malassiotis, and T.-K. Kim, "Recovering 6d object pose and predicting next-best-view in the crowd," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [24] U. Bonde, V. Badrinarayanan, and R. Cipolla, "Robust instance recognition in presence of occlusion and clutter," in *European Conference on Computer Vision (ECCV)*, September 2014.
- [25] S. Hinterstoisser, V. Lepetit, S. Ilic, S. Holzer, G. Bradski, K. Konolige, and N. Navab, "Model based training, detection and pose estimation of texture-less 3d objects in heavily cluttered scenes," in *Asian Conference on Computer Vision (ACCV)*, November 2012.
- [26] E. Brachmann, A. Krull, F. Michel, S. Gumhold, J. Shotton, and C. Rother, "Learning 6d object pose estimation using 3d object coordinates," in *European Conference on Computer Vision (ECCV)*, September 2014.
- [27] R. Brégier, F. Devernay, L. Leyrit, and J. L. Crowley, "Symmetry aware evaluation of 3d object detection and pose estimation in scenes of many parts in bulk," in *IEEE International Conference on Computer Vision (ICCV)*, October 2017.
- [28] K. Kleeberger, C. Landgraf, and M. F. Huber, "Large-scale 6d object pose estimation dataset for industrial bin-picking," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, November 2019.
- [29] 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 *International Conference on Computer Vision (ICCV)*, November 2011.
- [30] D. G. Lowe, "Object recognition from local scale-invariant features," in *International Conference on Computer Vision (ICCV)*, September 1999.
- [31] F. Rothganger, S. Lazebnik, C. Schmid, and J. Ponce, "3d object modeling and recognition using local affine-invariant image descriptors and multi-view spatial constraints," *International Journal of Computer Vision (IJCV)*, vol. 66, no. 3, pp. 231–259, March 2006.
- [32] T.-T. Do, M. Cai, T. Pham, and I. D. Reid, "Deep-6DPose: Recovering 6d object pose from a single rgb image," *arXiv e-prints*, vol. abs/1802.10367, February 2018.
- [33] Z. Dong, S. Liu, T. Zhou, H. Cheng, L. Zeng, X. Yu, and H. Liu, "PPR-Net: point-wise pose regression network for instance segmentation and 6d pose estimation in bin-picking scenarios," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, November 2019.
- [34] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "PointNet++: Deep hierarchical feature learning on point sets in a metric space," in *Neural Information Processing Systems (NIPS)*, December 2017.
- [35] T. Hodaň, F. Michel, E. Brachmann, W. Kehl, A. GlentBuch, D. Kraft, B. Drost, J. Vidal, S. Ihrke, X. Zabulis, C. Sahin, F. Manhardt, F. Tombari, T.-K. Kim, J. Matas, and C. Rother, "BOP: Benchmark for 6d object pose estimation," in *European Conference on Computer Vision (ECCV)*, September 2018.
- [36] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Neural Information Processing Systems (NIPS)*, December 2012.
- [37] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [38] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.
- [39] R. Brégier, F. Devernay, L. Leyrit, and J. L. Crowley, "Defining the pose of any 3d rigid object and an associated distance," *International Journal of Computer Vision (IJCV)*, vol. 126, no. 6, pp. 571–596, June 2018.
- [40] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel, "Domain randomization for transferring deep neural networks from simulation to the real world," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, September 2017.
- [41] A. Shrivastava, T. Pfister, O. Tuzel, J. Susskind, W. Wang, and R. Webb, "Learning from simulated and unsupervised images through adversarial training," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [42] K. Bousmalis, N. Silberman, D. Dohan, D. Erhan, and D. Krishnan, "Unsupervised pixel-level domain adaptation with generative adversarial networks," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [43] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Neural Information Processing Systems (NIPS)*, December 2014.
- [44] B. Drost, M. Ulrich, N. Navab, and S. Ilic, "Model globally, match locally: Efficient and robust 3d object recognition," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2010.