# Investigations on Output Parameterizations of Neural Networks for Single Shot 6D Object Pose Estimation

Kilian Kleeberger[1], Markus Völk[1], Richard Bormann[1], and Marco F. Huber[1,2]

*Abstract*—Single shot approaches have demonstrated tremendous success on various computer vision tasks. Finding good parameterizations for 6D object pose estimation remains an open challenge. In this work, we propose different novel parameterizations for the output of the neural network for single shot 6D object pose estimation. Our learning-based approach achieves state-of-the-art performance on two public benchmark datasets. Furthermore, we demonstrate that the pose estimates can be used for real-world robotic grasping tasks without additional ICP refinement.
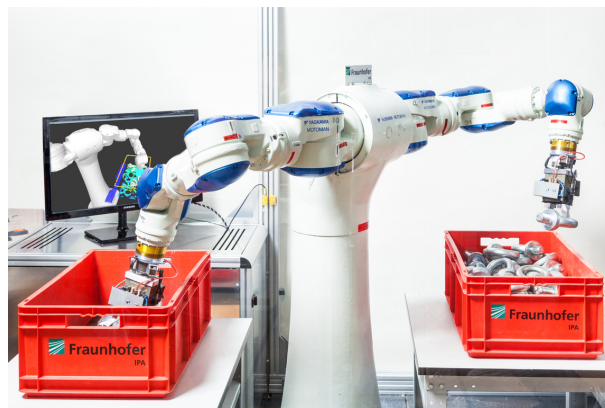
## I. INTRODUCTION

A reliable, fast, and accurate method for 6D object pose estimation is a crucial prerequisite for many robotic grasping and manipulation tasks. Based on a single depth image of known rigid objects, our goal is to estimate the translation vector $\mathbf{t} \in \mathbb{R}^3$ and rotation matrix $\mathbf{R} \in \mathrm{SO}(3)$ of an object relative to the sensor coordinate system. Using this information, a robot can plan a kinematically feasible and collision-free path to pick and place the object as visualized in Fig. 1.
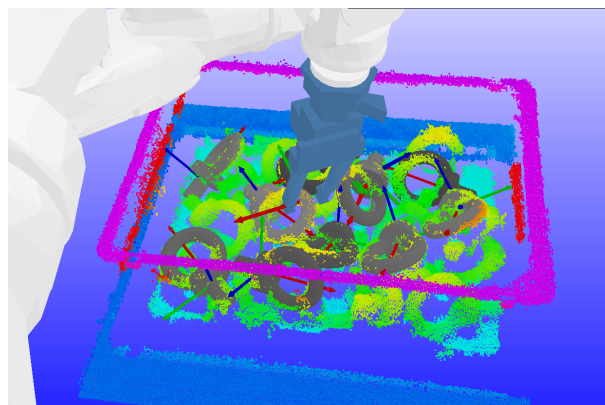
Object pose estimation is challenging because of a potentially very high amount of clutter and occlusion in the scene. Object symmetries propose different annotations for identical observations and typically have to be addressed explicitly for learning-based approaches [1], [2], [3]. Furthermore, occlusions can also lead to pose ambiguities and different lighting conditions affect the appearance of the object in the image. The task is also challenging due to missing, wrong, and noisy depth information.

Learning-based methods have demonstrated tremendous success on various benchmarks [5], [6]. Due to the high amount of data needed for training deep neural networks, using simulations is an attractive choice as they provide an abundant source of data with flawless annotations [7], [8], [9]. Single shot approaches are very fast because they require a single forward pass of the neural network and consider global scene context instead of looking at local image patches only. They have demonstrated tremendous advantages in speed, also for instance segmentation [10], [11].

OP-Net (Object Pose Network) [1] outperformed PPR-Net [2], the winning method of the "Object Pose Estimation Challenge for Bin-Picking" at IROS 2019 [5], on the Siléane

[1]Fraunhofer Institute for Manufacturing Engineering and Automation IPA, Nobelstraße 12, 70569 Stuttgart, Germany `kilian.kleeberger@ipa.fraunhofer.de`
[2]Institute of Industrial Manufacturing and Management IFF, University of Stuttgart, Allmandring 35, 70569 Stuttgart, Germany `marco.huber@ieee.org`

(a)



(b)

Fig. 1. (a) Real-world robot cell for random bin-picking. (b) 3D point cloud (colored) with pose estimates (grey) of OP-Net AP on real-world data without ICP refinement for ring screws. The gripper (blue) is visualized at the chosen grasp pose for execution together with the joint configuration of the robot (white).

dataset [4] based on the commonly used evaluation metric for 6D object pose estimation from Brégier et al. [12], [4], which considers all possible kinds of object symmetries and can handle scenarios in bulk. The metric requires the recovery of the pose of all objects in the image with a visibility of 50% or higher. Furthermore, OP-Net has been extended to model-based robotic pick-and-place [13].

Despite providing robust pose estimates at a very high speed, the method comes with shortcomings regarding the parameterization of the output which limit the performance. Several limitations are visualized in Fig. 2 (see Section III for more details). While the design of the network output is trivial for tasks such as semantic segmentation (one-hot encoded class vector for each pixel),
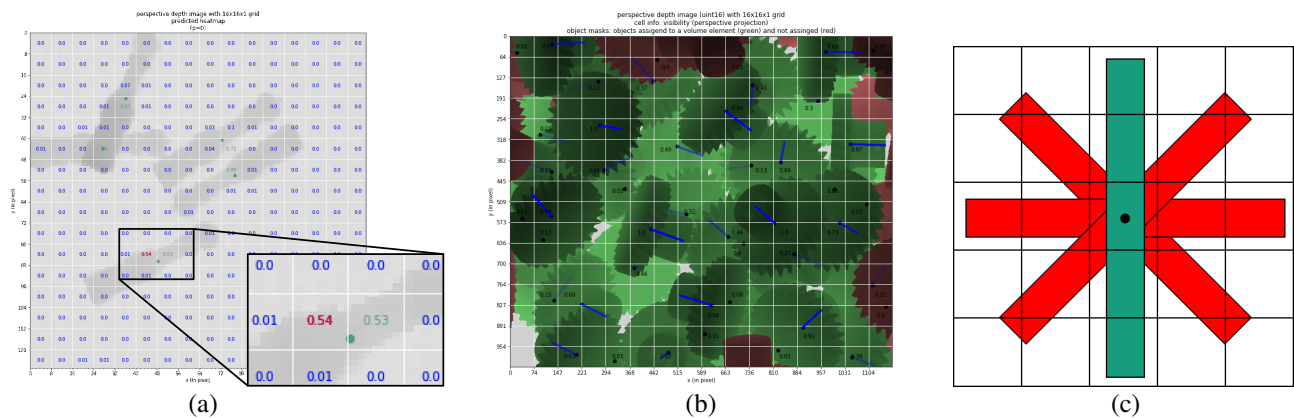
Fig. 2. Major limitations of single shot approaches to 6D object pose estimation using an assignment based on the origin of the object: (a) Exemplary prediction for the presence of an object origin for the pepper object from the Siléane dataset [4]: If an object origin is close to the border of a volume element, it might also be incorrectly detected by the neighboring spatial location (red), potentially with a higher confidence. Both cells detect the object with low confidence ($\hat{p} \in [0, 1]$). (b) Exemplary ground truth sample for the gear object from the Siléane dataset [4]: Objects, which are assigned to the ground truth tensor, are masked in green. Some objects at the border of the image are well visible but cannot be predicted because their object origin is outside of the image (masked in red). (c) The pose of well visible objects (red) cannot be predicted because each spatial location can only predict one pose. The object with the highest visibility (green) is used in the ground truth tensor.

finding a suitable output is challenging for tasks such as object detection [14], [15], [16], object pose estimation [1], [17], [18], or instance segmentation [19], [10], [11]. In this work, we address these challenges and give solutions by proposing different novel parameterizations for the output of the neural network. Fig. 1 (b) shows pose estimates of our approach on real-world data, which are accurate enough for reliable grasping with a robot even without ICP refinement. Videos of our experiments are available at https://owncloud.fraunhofer.de/index.php/s/6QsRj5sSty7fI9c.

In summary, the main contributions of this work are:

- Propose different novel parameterizations of the output of the neural network which solve problems of current methods using an origin-based assignment
- Extensive evaluation on public benchmark datasets
- Real-world robot demonstration of the approach for bin-picking

## II. RELATED WORK

Object pose estimation is a fundamental task in computer vision and an active field of research [3]. Challenges such as the "Object Pose Estimation Challenge for Bin-Picking" [5], SIXD Challenge [20], or BOP (Benchmark for 6D Object Pose Estimation) Challenge [21], [22], [6], [23] aim to capture and advance the state of the art in the field. Contrary to the SIXD Challenge [20] in 2017 and the BOP Challenge in 2019 [22], in the BOP Challenge 2020 [6] learning-based approaches caught up with classical approaches, including methods based on point pair features [24], which dominated previous editions of the challenge. CNN-based methods gained popularity in recent years, are a dominant research direction, and outperformed classical methods [6] based on feature [24] or template [25] matching.

PPR-Net [2], the winning method of the "Object Pose Estimation Challenge for Bin-Picking" at IROS 2019 [5] uses PointNet++ [26], estimates a 6D pose for each point in the point cloud and gets the final pose hypothesis by averaging the resulting pose clusters in 6D space. OP-Net [1] introduces a spatial discretization of the measurement volume of the sensor and solves a regression task for each resulting volume element. Other single shot approaches to object pose estimation output the 2D projections of the 3D bounding box and use a P$n$P algorithm [27] to compute the 6D pose [17], [18], [28], [29]. This results in a suboptimal two-stage process [30]. Our single shot variants directly regress the full 6D pose. The methods provide the pose of multiple objects simultaneously (single shot). Approaches with pixel- or point-wise predictions [2] are slower because of using a less compact parameterization and using post processing for averaging the predicted poses.

## III. PROBLEM STATEMENT

Based on a single depth image, our goal is to estimate the 6D pose of known rigid objects from a single category relative to the sensor coordinate system. Current single shot systems for 6D object pose estimation [1], [18], [13] discretize the measurement volume of the sensor in $S_x \times S_y$ volume elements. Due to the perspective projection, the shape of these volume elements corresponds to truncated oblique rectangular pyramids. At each spatial location, a feature vector of the object is assigned, which results in a 3D output tensor. The feature vector $\mathbf{x} = [p, v, x, y, z, \varphi_1, \varphi_2, \varphi_3]$ comprises the presence $p = 1$ of an object, the visibility $v \in [0, 1]$, the positions $x$, $y$, $z$, and the Euler angles $\varphi_1$, $\varphi_2$, $\varphi_3$. For objects with a revolution symmetry, it is sufficient to predict the angles $\varphi_1$ and $\varphi_2$, as the rotation around the $z$-axis with $\varphi_3$ results in identical observations. In case multiple objects fall into the same volume element, the object with highest visibility is chosen for assignment. For all spatial locations without an object, a zero vector is assigned.

The parameterization of these single shot systems for 6D object pose estimation [1], [18], [13] comes with several limitations, which are presented in Fig. 2. In case the origin of an object is close to the border of a volume element, the object might also be detected by the neighboring volume element (also with higher confidence). This results in duplicates and the neighboring cell has not been trained to output useful pose information. In Fig. 2 (a), both cells detect the object with low confidence. Furthermore, the single shot approaches to object pose estimation cannot predict objects with their object origin outside of the image (see objects masked in red in Fig. 2 (b)). Moreover, the parameterization is not unique when multiple objects fall into the same spatial location with their object origin (see Fig. 2 (c)). The objects are neglected in the ground truth and therefore cannot be predicted at test time. In this work, we address these shortcomings mentioned above by providing different novel parameterizations of the output to overcome the limitations of previous works and improve the performance of single shot approaches.

## IV. APPROACH

This section describes various novel output parameterizations of the neural network, the loss function, the averaging of poses for variants with multiple predictions per object, and the technique for a robust sim-to-real transfer.

### A. Output Parameterizations

In this section, we describe different parameterizations of the neural network output.

*1) Extended Volume Elements (OP-Net EVE):* For this variant, we also assign the feature vector **x** to the neighboring volume elements if the origin of the object coordinate system is close to the border of a volume element because the object might also get incorrectly localized by the neighboring spatial location as visualized in Fig. 2 (a). With this modification, the model does not get a jumping loss for challenging border cases during training and the accepted duplicate pose estimates can be averaged (see Section IV-C).

With these extensions to neighboring volume elements, we do not overwrite origins of other objects. In case multiple objects extend to the same volume element, we assign the feature vector of the object with highest visibility. In our experiments, we use a distance threshold from the border of 0.2 of the object origin relative to the spatial location. In the original version, the neighboring volume elements are not trained to propose correct pose information. Fig. 3 visualizes an exemplary ground truth sample with extended volume elements.

*2) Additional Points (OP-Net AP):* An object may be well visible in the image (i.e., high visibility), but another more dominant object (higher visibility) is already assigned to the volume element as visualized in Fig. 2 (c). As the spatial locations in the surrounding of the object do not have the task to predict a pose, they can be used for estimating additional ground truth information based on additional reference points on the object which are defined relative to the object coordinate system. Objects which are visible at the border
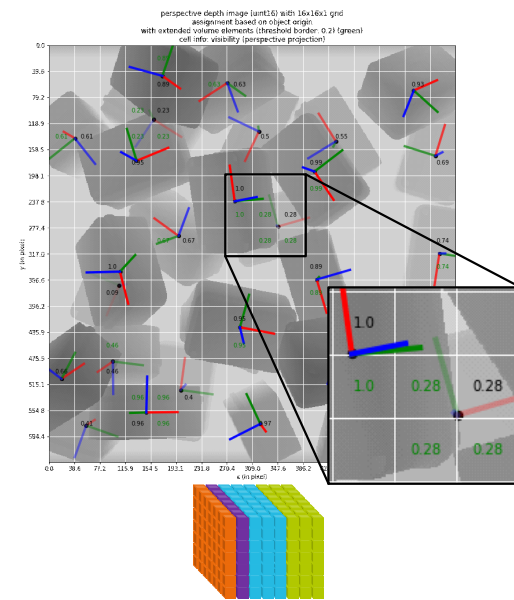


Fig. 3. (top) Exemplary ground truth sample for extended volume elements (EVE) for the T-Less 29 object from the Siléane dataset [4]. The feature vector of an object is also assigned to the neighboring elements if the origin of the object is close to the border of the cell. The extended elements are highlighted in green. The cell information indicates the visibility of the object. (bottom) Resulting output tensor (shape unchanged).

of an image but with their origin outside can be predicted by this variant, because the additional points are inside the image (contrary to the object origin itself). Therefore, this variant addresses the limitation in Fig. 2 (b) and (c) plus implicitly also (a) because it is very unlikely that the object origin and all additional points are very close to a border. Fig. 4 visualizes the concept and demonstrates that much fewer objects will be missed by the used metric for objects pose estimation.

With additional points, we do not overwrite origins of other objects and assign the feature vector of the object with highest visibility for conflicting spatial locations. For objects with a revolution symmetry (e.g., candlestick, pepper, gear, gear shaft), the additional points have to lie on the axis of symmetry. For cyclic symmetries (e.g., brick, tless 20, tless 29, ring screw), the additional points have to permute the location when applying the steps of rotation.

*3) Discretization in z-Direction (OP-Net Z):* In addition to the already used spatial discretization in $x$- and $y$-direction, we introduce a variant with additional discretization in $z$-direction from the near to the far clipping plane of the sensor. Fig. 5 visualizes this concept for the perspective projection. With this, we address the limitation in Fig. 2 (c) because the objects lie at different heights. This parameterization misses much fewer objects, which are relevant for the evaluation metric, i.e., objects with visibility $v \geq 0.5$. With the 8-dimensional feature vector, the output tensor is of shape $S_x \times S_y \times 8 \cdot S_z$ for this variant. In our experiments, we choose $S_x = S_y = S_z = 16$.

*4) Multiple Poses (OP-Net MP):* As the vanilla OP-Net can predict one 6D pose per volume element only, we
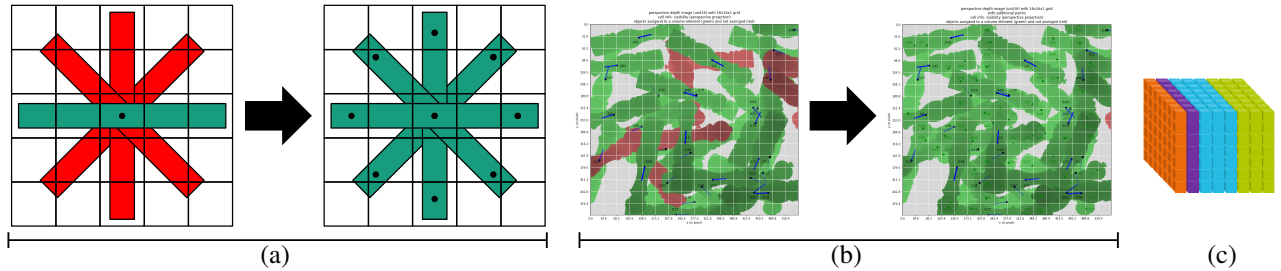
Fig. 4. (a) Concept for additional points (AP): Objects masked in green are assigned to a volume element and objects masked in red are missed. With additional points much fewer objects are missed by the parameterization and, therefore, also at test time. (b) An exemplary ground truth sample for the pepper object from the Siléane dataset [4] with two additional points. (c) Resulting output tensor (shape unchanged).
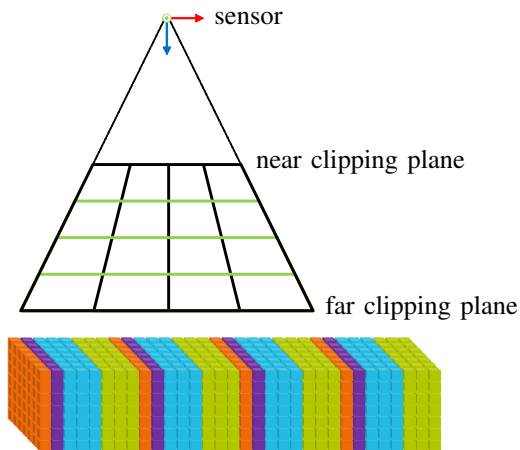


Fig. 5. (top) Discretization of the measurement volume of the 3D sensor in $z$-direction in addition to the $x$- and $y$-direction. The green lines indicate the introduced discretization in $z$-direction. (bottom) Resulting output tensor ($S_z$ times the vanilla OP-Net tensor).
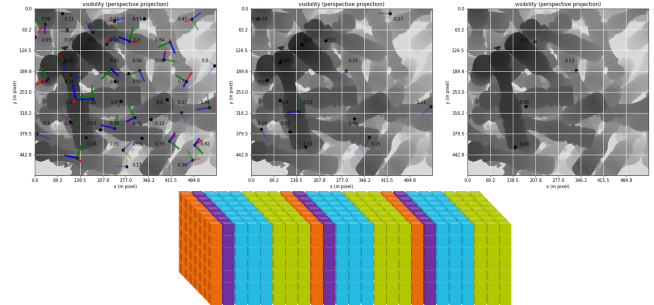


Fig. 6. (top) Exemplary ground truth sample for predicting multiple poses (MP) per spatial location for the pepper object from the Siléane dataset [4]. (bottom) Resulting output tensor ($P$ times the vanilla OP-Net tensor).

propose a variant that outputs multiple poses per spatial location to be able to predict the pose of less dominant objects in terms of visibility. This addresses the limitation visualized in Fig. 2 (c). Each spatial location predicts $P$ feature vectors. For our experiments, we use $S_x = S_y = 8$, $S_z = 1$ (i.e., no discretization in $z$-direction), and $P = 3$. For practical scenarios, e.g., random bin-picking, it is very unlikely that more than three objects with $v \geq 0.5$ fall into the same volume element with the chosen configuration (this does not happen in the datasets).

We condition the additional output feature maps for the presence of an object on the channel, i.e., each channel predicts whether the next channel comprises further objects. The objects are ranked and assigned to the individual channels based on their visibility. Fig. 6 shows an exemplary ground truth sample for predicting multiple poses per volume element.

*5) Resized Segmentation Image (OP-Net SI):* Instead of assigning the objects to the ground truth tensor based on the origin, we propose to use the segmentation masks. We resize the segmentation image (comprises the IDs of the objects) to $32 \times 32$ pixel and assign the features of the objects to the tensor based on the ID of the pixels. This variant addresses

the limitations in Fig. 2 (a), (b), and (c). Because of the dense parameterization, we use a linear weighting for the pose error, i.e., $\lambda_3 = v$.

### B. Loss Function

To train the network, the multi-task loss function

$$\mathcal{L} = \sum_{i=1}^{S_x}\sum_{j=1}^{S_y}\sum_{k=1}^{S_z}\sum_{l=1}^{P}\left(\lambda_1\mathcal{L}_\mathrm{p} + \left[\lambda_2\mathcal{L}_\mathrm{v} + \lambda_3\left(\mathcal{L}_\mathrm{pos} + \lambda_4\mathcal{L}_\mathrm{ori}\right)\right]p_{ijkl}\right) \tag{1}$$

is optimized. The scalars $\lambda_1 = 0.1$, $\lambda_2 = 0.25$, $\lambda_3 = 8v^3$, and $\lambda_4 = 1$ are used for weighting the different loss terms.

We use a squared L2-loss for the presence of an object and the visibility to compute $\mathcal{L}_\mathrm{p}$ and $\mathcal{L}_\mathrm{v}$, respectively. Using the squared L2-norm, the positions $x, y, z \in [0, 1]$ of the object are estimated relative to the volume element, except for the variants, which make multiple predictions for one object (OP-Net EVE, AP, SI) because the object origin can be located outside of the spatial location. For these variants, the position is estimated relative to an image reference which is enlarged by the object diameter (diameter of the smallest bounding sphere of the object). This also allows to predict origins outside of the image. We also use the squared L2-norm for the regression of the Euler angles with bounded values, i.e., $\varphi_1, \varphi_2 \in [0, 2\pi)$ and $\varphi_3 \in [0, 2\pi/k)$ are mapped to $[0, 1]$, where $k \in \mathbb{N}$ represents the order of the cyclic symmetry (similar to vanilla OP-Net [1]).
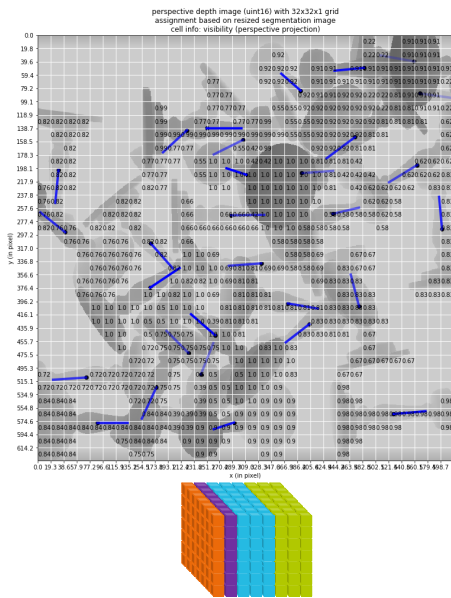
Fig. 7. (top) Exemplary ground truth sample for the assignment based on resized segmentation images (SI) for the candlestick object from the Siléane dataset [4]. (bottom) Resulting output tensor.

### C. Pose Averaging

The variants OP-Net EVE, AP, and SI potentially output multiple predictions per object. Since these should be close to each other in the 6D pose space and in order not to have many duplicates, we use unsupervised learning and make use of density-based clustering [31], [32]. To form the final pose predictions, we simply average the resulting pose clusters.

### D. Sim-to-Real Transfer

For a robust transfer of the models from simulation to the real world, we use domain randomization [33]. The technique to resize the image from the original size to the input size of the neural network ($128 \times 128$ pixel) is randomized during training of the model as a kind of noise by randomly selecting a resizing option available in OpenCV. Furthermore, we apply augmentations in a random order and with varying intensity, such as elastic transformations, blurring, and adding noise to the pixels of the image.

## V. EXPERIMENTAL EVALUATION

### A. Datasets and Evaluation Metric

We evaluate the performance of our variants of OP-Net on the Siléane [4] and Fraunhofer IPA [8] datasets. The Fraunhofer IPA dataset is a large-scale dataset that comprises synthetic data for training deep neural networks and fully pose annotated real-world data for performance evaluation. We use the synthetic datasets for training our models. These datasets show scenarios that are typical for bin-picking (see Fig. 1 (a)) and are very challenging because of a high amount of clutter and heavy occlusions in combination with sensor noise, wrong, or missing depth information.

For evaluation, we use the metric from Brégier et al. [12], [4], which considers all possible kinds of object symmetries,

can handle scenarios in bulk, requires the recovery of all objects with a visibility of 50% or higher, and breaks down the performance of a method to a single scalar value (average precision).

### B. Setup

We use a fully convolutional network architecture to learn the mapping from normalized input depth images in perspective projection to the output tensors, where $S_x$ and $S_y$ correspond to the size of the feature map of the last layer of the neural network. This method is a single shot approach because only one forward pass of the neural network is needed instead of multiple evaluations at different locations.

Analogously to OP-Net [1], we use a DenseNet-BC [35], [36] with 40 layers and a growth rate of 50, which indicates the number of feature maps being added per layer. For best comparability, we also use the same loss function for the regression of the orientation (see Sec. IV-B) and an input resolution of $128 \times 128$ pixel. For the variants OP-Net EVE and AP, we use $S_x = S_y = 16$, similar to vanilla OP-Net for best comparability.

### C. Results

Our variants provide robust pose estimates although they are fully trained on synthetic images and annotations. The test datasets were recorded with different 3D sensors, i.e., providing robust pose estimates is independent of the actual 3D sensor technology being used. Pixels with missing depth information are bilinearly interpolated before being fed into the neural network. Table I reports the performance of the different variants in terms of average precision, where all values of OP-Net and its variants are without ICP refinement. Our variants outperform other classical and learning-based approaches on multiple datasets. Fig. 8 visualizes qualitative results of OP-Net AP on the different datasets.

The variant OP-Net EVE can give slightly better results because it also outputs useful pose information for challenging border cases. OP-Net AP performs significantly better than other state-of-the-art approaches because it misses fewer objects in the parameterization and addresses all major limitations in Fig. 2. OP-Net Z and MP also provide good results but not consistently along all datasets. A key limitation is the very high-dimensional output tensor. For OP-Net SI, we encountered the transfer from synthetic to real-world data as a key limitation due to the dense parameterization of the output.

The OP-Net variants are approaches which require a single forward pass of the network only (single shot) with 13 or 17 ms per depth image for the forward pass on a Nvidia Tesla V100 or GTX 1080 Ti, respectively. With 200 ms for the forward pass (GTX 1060) plus the time for the additional clustering and averaging of many datapoints, PPR-Net [2] is much slower than OP-Net and its variants.

## VI. CONCLUSIONS

In this work, we propose different novel parameterizations of the output for single shot approaches to 6D object pose

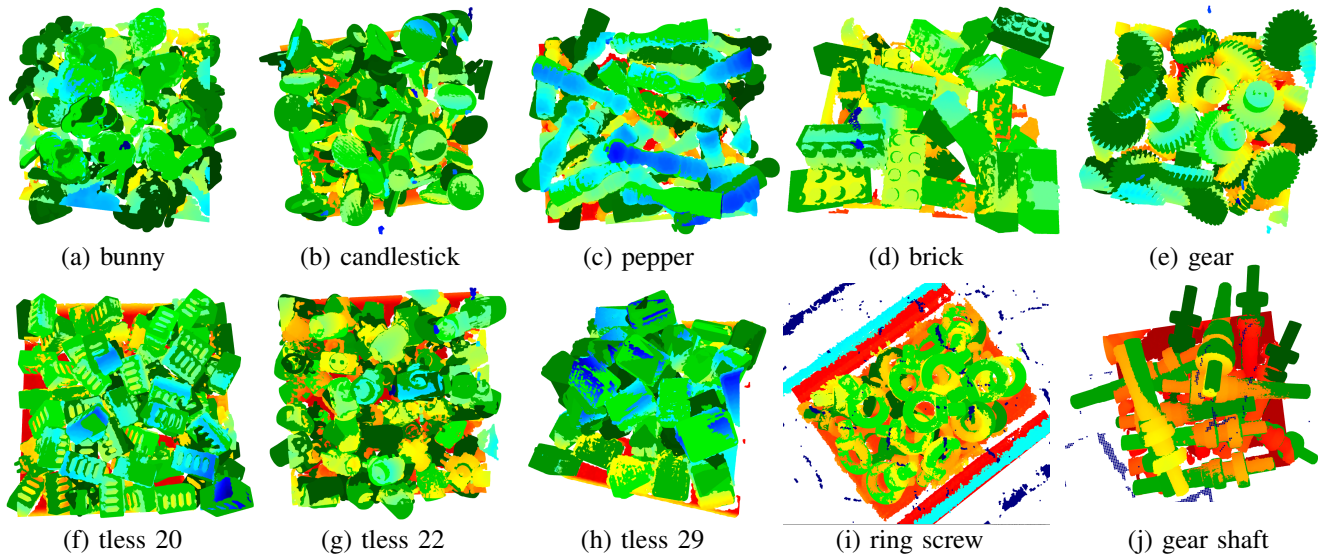| (a) bunny | (b) candlestick | (c) pepper | (d) brick | (e) gear |
| (f) tless 20 | (g) tless 22 | (h) tless 29 | (i) ring screw | (j) gear shaft |

Fig. 8. Qualitative results on the Siléane [4] (a–h) and Fraunhofer IPA [8] (i, j) datasets of OP-Net AP without ICP refinement: 3D point cloud (colored) with pose estimates in green (the brighter, the higher the confidence). Note that the models were trained on synthetic images and annotations only and provide robust results on noisy real-world data.

TABLE I

AVERAGE PRECISION VALUES OF DIFFERENT METHODS (BEST RESULTS IN BOLD) ON THE OBJECTS FROM THE SILÉANE [4] AND FRAUNHOFER IPA [8] DATASETS. FOR OUR VARIANTS, WE DID NOT USE ICP REFINEMENT. RESULTS MARKED WITH * ARE TAKEN FROM THE "OBJECT POSE ESTIMATION CHALLENGE FOR BIN-PICKING" AT IROS 2019 [5].

| object<br>object symmetry based on [12], [4] | bunny [4]<br>(no proper symmetry) | candlestick [4]<br>(revolution) | pepper [4]<br>(revolution) | brick [4]<br>(cyclic, $k=2$) | gear [4]<br>(revolution) | tless 20 [4]<br>(cyclic, $k=2$) | tless 22 [4]<br>(no proper symmetry) | tless 29 [4]<br>(cyclic, $k=2$) | ring screw [8]<br>(cyclic, $k=2$) | gear shaft [8]<br>(revolution) | average |
|---|---|---|---|---|---|---|---|---|---|---|---|
| PPF [24], [4] | 0.29 | 0.16 | 0.06 | 0.08 | 0.62 | 0.20 | 0.08 | 0.19 | - | - | - |
| PPF PP [24], [4] | 0.37 | 0.22 | 0.12 | 0.13 | 0.63 | 0.23 | 0.12 | 0.23 | - | - | - |
| LINEMOD+ [25], [4] | 0.39 | 0.38 | 0.04 | 0.31 | 0.44 | 0.25 | 0.19 | 0.20 | - | - | - |
| LINEMOD+ PP [25], [4] | 0.45 | 0.49 | 0.03 | 0.39 | 0.50 | 0.31 | 0.21 | 0.26 | - | - | - |
| Sock et al. [34] | 0.74 | 0.64 | 0.43 | - | - | - | - | - | - | - | - |
| PPR-Net [2] | 0.82 | 0.91 | 0.80 | - | - | 0.81 | - | - | 0.95* | 0.99* | - |
| PPR-Net with ICP [2] | 0.89 | 0.95 | 0.84 | - | - | 0.85 | - | - | - | - | - |
| OP-Net with $\mathcal{L}_{ori1}$ [1] | **0.92** | 0.94 | 0.98 | 0.41 | **0.82** | 0.85 | 0.77 | 0.51 | 0.88 | 0.99 | 0.81 |
| OP-Net with $\mathcal{L}_{ori2}$ [1] | 0.74 | 0.95 | 0.92 | **0.79** | 0.58 | 0.56 | 0.53 | 0.36 | 0.73 | **1.0** | 0.72 |
| OP-Net EVE (ours) | **0.92** | 0.94 | 0.97 | 0.45 | 0.77 | **0.88** | 0.81 | **0.56** | 0.94 | 0.98 | 0.82 |
| OP-Net AP (ours) | **0.92** | **0.98** | **0.99** | 0.45 | **0.82** | 0.87 | **0.84** | **0.56** | **0.96** | 0.99 | **0.84** |
| OP-Net Z (ours) | 0.85 | 0.92 | 0.95 | 0.41 | 0.72 | 0.84 | **0.84** | 0.43 | 0.92 | 0.95 | 0.78 |
| OP-Net MP (ours) | 0.92 | 0.82 | 0.91 | 0.32 | 0.80 | 0.70 | 0.66 | **0.56** | 0.84 | 0.95 | 0.75 |
| OP-Net SI (ours) | 0.30 | 0.87 | 0.97 | 0.36 | 0.37 | 0.46 | 0.51 | 0.36 | 0.56 | 0.88 | 0.56 |

estimation. Our experiments demonstrate that extending the spatial locations or adding additional points to the object significantly improves the performance in terms of average precision, allows overcoming the limitations of previous works, and improves the performance of single shot approaches. Our models demonstrate state-of-the-art performance on various public datasets and can be used for real-world robotic grasping tasks without ICP refinement.

REFERENCES

[1] K. Kleeberger and M. F. Huber, "Single shot 6d object pose estimation," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2020.

[2] 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)*, 2019.

[3] Y. Labbé, J. Carpentier, M. Aubry, and J. Sivic, "CosyPose: Consistent multi-view multi-object 6d pose estimation," in *European Conference on Computer Vision (ECCV)*, 2020.

[4] 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)*, 2017.

[5] K. Kleeberger and M. F. Huber, "Object pose estimation challenge for bin-picking," 2019. [Online]. Available: http://www.bin-picking.ai/en/competition.html

[6] T. Hodaň, M. Sundermeyer, B. Drost, Y. Labbé, E. Brachmann, F. Michel, C. Rother, and J. Matas, "BOP challenge 2020 on 6d object localization," in *European Conference on Computer Vision (ECCV) Workshop*, 2020.

[7] K. Kleeberger, R. Bormann, W. Kraus, and M. F. Huber, "A survey on learning-based robotic grasping," *Current Robotics Reports*, vol. 1, no. 4, pp. 239–249, 2020.

[8] 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)*, 2019.

[9] M. Denninger, M. Sundermeyer, D. Winkelbauer, D. Olefir, T. Hodany, Y. Zidan, M. Elbadrawy, M. Knauer, H. Katam, and A. Lodhi, "BlenderProc: Reducing the reality gap with photorealistic rendering," in *Robotics: Science and Systems (RSS) Workshop*, 2020.

[10] D. Bolya, C. Zhou, F. Xiao, and Y. J. Lee, "YOLACT: Real-time instance segmentation," in *IEEE International Conference on Computer Vision (ICCV)*, 2019.

[11] ——, "YOLACT++: Better real-time instance segmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2020.

[12] 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, 2018.

[13] K. Kleeberger, M. Völk, M. Moosmann, E. Thiessenhusen, F. Roth, R. Bormann, and M. F. Huber, "Transferring experience from simulation to the real world for precise pick-and-place tasks in highly cluttered scenes," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2020.

[14] 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)*, 2016.

[15] J. Redmon and A. Farhadi, "YOLO9000: Better, faster, stronger," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.

[16] 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)*, 2016.

[17] 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)*, 2017.

[18] 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)*, 2018.

[19] K. He, G. Gkioxari, P. Dollar, and R. Girshick, "Mask R-CNN," in *IEEE International Conference on Computer Vision (ICCV)*, 2017.

[20] T. Hodaň, F. Michel, C. Sahin, T.-K. Kim, J. Matas, and C. Rother, "SIXD challenge 2017," 2017. [Online]. Available: http://cmp.felk.cvut.cz/sixd/challenge_2017/

[21] T. Hodaň, M. Sundermeyer, E. Brachmann, B. Drost, F. Michel, J. Matas, and C. Rother, "BOP challenge 2019/2020." [Online]. Available: https://bop.felk.cvut.cz/challenges/bop-challenge-2020/

[22] T. Hodaň, E. Brachmann, B. Drost, F. Michel, M. Sundermeyer, J. Matas, and C. Rother, "BOP challenge 2019." [Online]. Available: https://bop.felk.cvut.cz/media/bop_challenge_2019_results.pdf

[23] T. Hodaň, F. Michel, E. Brachmann, W. Kehl, A. Glent Buch, 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)*, 2018.

[24] 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)*, 2010.

[25] S. Hinterstoisser, V. Lepetit, S. Ilic, S. Holzer, K. 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)*, 2012.

[26] 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)*, 2017.

[27] V. Lepetit, F. Moreno-Noguer, and P. Fua, "Epnp: An accurate o(n) solution to the pnp problem," *International Journal of Computer Vision (IJCV)*, vol. 81, no. 2, pp. 155–166, 2009.

[28] 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)*, 2019.

[29] J. Tremblay, T. To, B. Sundaralingam, Y. Xiang, D. Fox, and S. Birchfield, "Deep object pose estimation for semantic robotic grasping of household objects," in *Conference on Robot Learning (CoRL)*, 2018.

[30] Y. Hu, P. Fua, W. Wang, and M. Salzmann, "Single-stage 6d object pose estimation," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

[31] M. Ester, H. P. Kriegel, J. Sander, and X. Xiaowei, "A density-based algorithm for discovering clusters in large spatial databases with noise," in *International Conference on Knowledge Discovery and Data Mining*, 1996.

[32] E. Schubert, J. Sander, M. Ester, H. P. Kriegel, and X. Xu, "DBSCAN revisited, revisited: Why and how you should (still) use DBSCAN," *ACM Transactions on Database Systems*, vol. 42, no. 3, pp. 1–21, 2017.

[33] 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)*, 2017.

[34] 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)*, 2018.

[35] 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)*, 2017.

[36] G. Huang, Z. Liu, G. Pleiss, L. van der Maaten, and K. Q. Weinberger, "Convolutional networks with dense connectivity," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2019.