

ConvPoseCNN2: Prediction and Refinement of Dense 6D Object Poses

Arul Selvam Periyasamy, Catherine Capellen, Max Schwarz, and Sven Behnke

Autonomous Intelligent Systems, University of Bonn, Germany.
periyasamy@ais.uni-bonn.de

Abstract. Object pose estimation is a key perceptual capability in robotics. We propose a fully-convolutional extension of the PoseCNN method, which densely predicts object translations and orientations. This has several advantages such as improving the spatial resolution of the orientation predictions—useful in highly-cluttered arrangements, significant reduction in parameters by avoiding full connectivity, and fast inference. We propose and discuss several aggregation methods for dense orientation predictions that can be applied as a post-processing step, such as averaging and clustering techniques. We demonstrate that our method achieves the same accuracy as PoseCNN on the challenging YCB-Video dataset and provide a detailed ablation study of several variants of our method. Finally, we demonstrate that the model can be further improved by inserting an iterative refinement module into the middle of the network, which enforces consistency of the prediction.

Keywords: Monocular pose estimation · Fully-convolutional architectures · Robotics.

1 Introduction

6D object pose estimation is an important building block for many applications, such as robotic manipulation. While many objects can be grasped without precise pose information, there are many tasks which require 6D pose estimates, for example functional grasping of tools and assembly. Such tasks routinely come up in industrial applications, as evidenced by the Amazon Picking & Robotics Challenges 2015-2017, where pose estimation played a key role for difficult objects, but can also appear in semi-unstructured environments, as in home and assistance robotics.

State-of-the-art pose estimation methods predominantly use CNNs for 6D object pose estimation from RGB(-D) images. One of the notable features of these methods is the joint learning of multiple simultaneous tasks such as object detection, semantic segmentation, and object pose estimation. Although 6D object pose estimation from RGB-D images is an active area of research, for the sake of brevity, we focus on monocular, i.e. RGB only methods. These methods can be broadly classified into two categories: direct regression methods, and 2D-3D correspondence methods. The direct regression methods estimate 6D pose

directly from input images, for example in the form of a 3D vector (translation) and a quaternion (orientation). Examples of these methods include Do et al. [4] and Xiang et al. [22]. In contrast, the correspondence-based methods predict the projection of 3D points in the 2D image and recover the pose of the object by solving the Perspective-n-Point problem. These methods can be further classified into dense correspondence methods and keypoint-based methods. The dense correspondence methods [2, 8] predict the projected 3D coordinates of the objects per pixel while the keypoint-based methods [13, 17, 21, 16] predict projection of 3D keypoints in the 2D image.

Since the CNN architecture we propose is closely related to PoseCNN [22], a direct regression method, we review PoseCNN architecture in detail. PoseCNN learns to predict 6D pose objects jointly with semantic segmentation. The CNN uses a pretrained VGG [20] backbone followed by three branches to predict segmentation class probabilities, direction and distance to center, and orientation (represented as quaternions). The orientation prediction branch uses fully connected layers while the other two branches use fully convolutional layers. The orientation prediction branch takes a fixed size image crop as input. From the segmentation class probabilities, a crop containing a single object is extracted and resized to the fixed orientation prediction branch input size using a RoI pooling layer.

Introduced by Girshick [5], RoI pooling is a powerful mechanism for scale normalization and attention and resulted in significant advancements in object detection and related tasks. However, RoI pooling has drawbacks: Especially in cluttered scenes, its cutting-out operation may disrupt flow of contextual information. Furthermore, RoI pooling requires random access to memory for the cutting-out operation and subsequent interpolation, which may be expensive to implement in hardware circuits and has no equivalent in the visual cortex [7].

Redmon et al. [18] demonstrated that simpler, fully-convolutional architectures can attain the same accuracy, while being tremendously faster and having fewer parameters. In essence, fully-convolutional architectures can be thought of as sliding-window classifiers, which are equivalent to RoI pooling with a fixed window size. While the scale-invariance is lost, fully-convolutional architectures typically outperform RoI-based ones in terms of model size and training/inference speed. When addressing the inherent example imbalances during training with a customized loss function [11], fully-convolutional architectures reach state-of-the-art performance in object detection.

Following this idea, we propose a fully convolutional architecture for pose estimation, which can *densely* predict all required information such as object class and transformation. If required, the dense prediction can be post-processed and aggregated per object to obtain a final prediction.

Given the complex nature of the task, instead of directly predicting pose from the given RGB image of a scene, many approaches formulate pose estimation as an iterative refinement process: Given an initial pose estimate and high quality 3D model of the objects, the objects are rendered as per current pose estimate, a refined pose that minimizes difference between the rendered and the observed

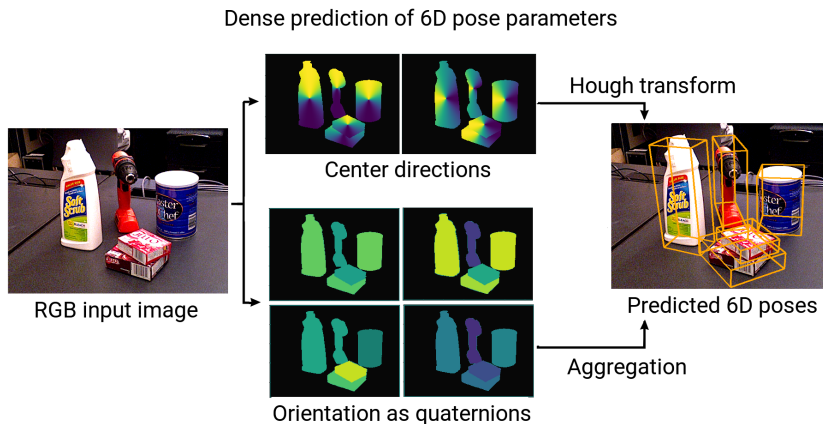


Fig. 1. Dense Prediction of 6D pose parameters inside ConvPoseCNN. The dense predictions are aggregated on the object level to form 6D pose outputs. Source: Capellen., Schwarz., and Behnke. [3].

image is predicted at each step and this step is repeated multiple times. Li et al. [10] trained a CNN that takes RGB image and rendered image of a object as per the current pose estimate as input and predicts the a pose update that refines the current pose update in each step. This step is repeated until the pose update is negligible. Periyasamy, Schwarz, and Behnke [15] used a differentiable renderer to compute pose updates to minimize difference between the rendered and the observed image. Unlike [10] that refines pose of single object at a time, [15] refined poses for all objects in the scene at each iteration. Krull et al. [8] trained a CNN to predict a matching score—how similar are two images— between the rendered and the observed image. The matching score was used to pick one best pose hypothesis among many available pose hypotheses. One prevalent characteristic among the pose refinement approaches is that refinement is done post prediction—refinement model and pose prediction model are decoupled. In contrast, our proposed iterative refinement module is built into the pose estimator. We enhance the ConvPoseCNN architecture from our previous work [3] with an iterative refinement module to learn representations suitable for both translation and orientation predictions instead of refining the predictions from the estimator.

In summary, our contributions include:

- A network architecture and training regime for dense orientation prediction,
- aggregation & clustering techniques for dense orientation predictions, and
- an iterative refinement module which increases prediction accuracy.

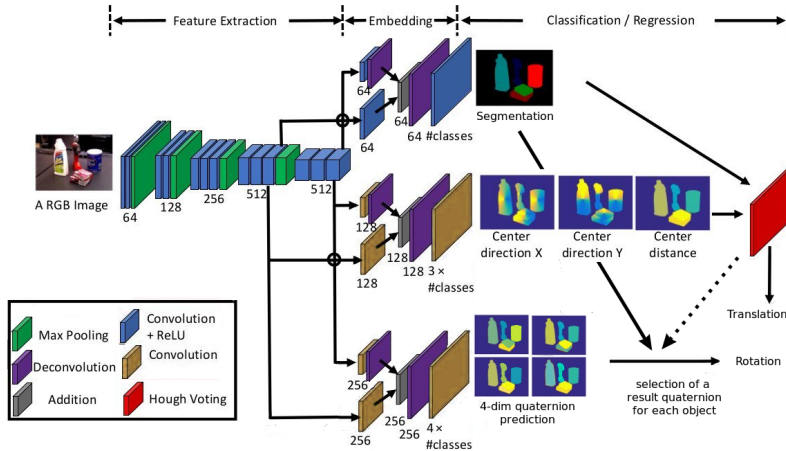


Fig. 2. Our ConvPoseCNN architecture for convolutional pose estimation. During aggregation, candidate quaternions are selected according to the semantic segmentation results or according to Hough inlier information. Source: Capellen., Schwarz., and Behnke. [3].

2 ConvPoseCNN

We propose an extension of the PoseCNN [22] architecture. The PoseCNN network is based on a VGG backbone with three heads: One performing semantic segmentation, one densely predicting object center directions in 2D and object depth, and finally a RoI-Pooling branch with a fully connected head predicting one orientation quaternion for each object. Our proposed network keeps most of this structure, but replaces the orientation prediction branch with a fully convolutional one, which estimates orientation *densely* (see Fig. 2). The architecture of the new branch is modeled after the translation estimation branch.

The dense translation prediction is post-processed during inference as by Xiang et al. [22]: The 2D center predictions are fed into a Hough voting layer, which aggregates them into center hypotheses. The predicted object depth is averaged over all inliers. Finally, the 3D position can be computed through ray projection using the camera intrinsics.

2.1 Aggregation of Dense Orientation Predictions

Estimating the final orientation prediction from pixel-wise quaternion predictions is not as straight-forward, however. We investigate two different approaches for this purpose: averaging and clustering.

Quaternions corresponding to a rotation, by definition, have unit norm. But we do not enforce the quaternion predictions to be of unit norm explicitly during the ConvPoseCNN training. Thus, before aggregating the dense predictions,

we need to scale them to unit norm. Interestingly, we observe that the norm of the quaternion at a pixel prior to scaling corresponds to quality of the prediction. i.e. pixels in the feature-rich regions of the image have higher quaternion norm. Exploiting this observation, we use the norm of the quaternion prediction $w = \|q\|$ as an optional weighting factor in our aggregation step. We extract the quaternions q_1, \dots, q_n corresponding to an object using the segmentation predictions and average them following the the optimization scheme proposed by [12] using the norm w_1, \dots, w_n . The average quaternion \bar{q} is given by

$$\bar{q} = \arg \min_{q \in \mathbb{S}^3} \sum_{i=1}^n w_i \|R(q) - R(q_i)\|_F^2, \quad (1)$$

where $R(q) - R(q_i)$ are the rotation matrices corresponding to the quaternions, \mathbb{S}^3 is the unit 3-sphere, and $\|\cdot\|_F$ is the Frobenius norm. Note that quaternion to rotation matrix conversion eliminates any problems arising from the antipodal symmetry of the quaternion representation. The exact solution to the optimization problem can be found by solving an eigenvalue problem [12]. In case of multiple, overlapping instances of the same object class—here, the predicted segmentation would not be enough to differentiate the instances—we can additionally make use of the Hough voting procedure required for translation estimation to separate the predictions into inlier sets for each object hypothesis.

Averaging based aggregation schemes inherently may from suffer skewed results due to bad outlier predictions. Clustering based aggregation schemes should be less susceptible to outlier predictions.

We follow a weighted RANSAC clustering scheme as an alternative to averaging: For quaternions $Q = \{q_1, \dots, q_n\}$ and their weights w_1, \dots, w_n associated with one object we repeatedly choose a random quaternion $\bar{q} \in Q$ with a probability proportional to its weight and then determines the inlier set $\bar{Q} = \{q \in Q | d(q, \bar{q}) < t\}$, where $d(\cdot, \cdot)$ is the angular distance. Finally, the \bar{q} with largest $\sum_{q_i \in \bar{Q}} w_i$ is selected as the result quaternion.

2.2 Iterative refinement

During the prediction of 6D object poses, translation estimates and orientation estimates influence each other. Predicting translation and orientation components using separate branches as in ConvPoseCNN and PoseCNN does not allow the model to exploit the interdependence between translation and orientation estimates. This motivates in designing network architectures that can refine translation and orientation prediction iteratively to enable the network to model the interdependencies between the predictions. One naive way of doing pose refinement would be to perform refinement after prediction. To this end, we experimented with a simple three layered—three blocks of convolutional layer followed by ReLU activation—fully convolutional model to refine the predictions from ConvPoseCNN model iteratively. At each step, segmentation, translation, and orientation predictions along with the features from the VGG backbone model are provided as input and a refined estimate is computed as depicted

in Figure. 3. The final predictions are obtained after a small fixed number of iterations. We call this approach post-prediction iterative refinement.

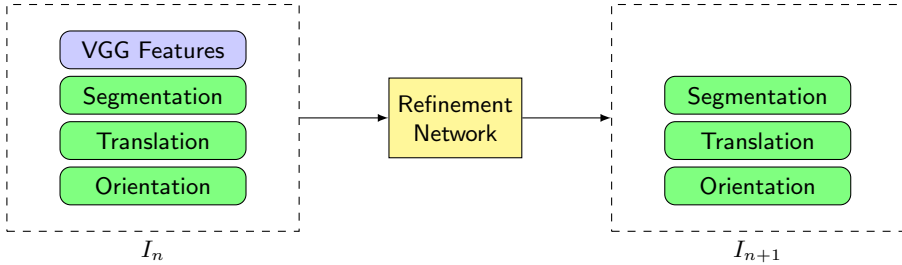


Fig. 3. Naive post-prediction iterative refinement of segmentation probabilities, translation predictions, and rotation predictions. The dense predictions (green) are refined using a small network, which can be applied repeatedly for further refinement. High-level features (blue) can be fed into the network to provide additional context information.

However, naive post-prediction refinement might be a challenging task because the predictions might be in a suitable form for a simple three layer model. To address this concern, we experimented with a pre-prediction iterative refinement of intermediate representation shown in Figure 4. The features from the pretrained VGG backbone model are refined before providing them as input to ConvPoseCNN network enabling ConvPoseCNN model to learn joint intermediate representations suitable for both translation and orientation predictions. The refinement module is akin to residual blocks in ResNet architecture (Ren et al. [19]). Each iteration refinement module computes $\Delta(x)$ that is added to the input with the use of skip connections.

$$f^{i+1}(x) = f^i(x) + \Delta(x)$$

In detail, refinement blocks takes two set of features maps F_A , and F_B each of dimension $512 \times 60 \times 80$, and $512 \times 30 \times 40$ respectively as input. F_B is upsampled with transposed convolution and concatenated with F_A . The resulting $1024 \times 60 \times 80$ is passed through a sequence of convolutional, ReLU, and convolutional layers. All the convolutions have a window size of 3 and stride of 1. Zero padding of one pixels is applied to maintain the spatial resolution of the features.

Then the features are split to two equal parts. One of them is downsampled. Thus we arrive at ΔF_A and ΔF_B having same spatial dimensions as of F_A , and F_B respectively. Using skip connections ΔF_A and ΔF_B and added to F_A , and F_B respectively.

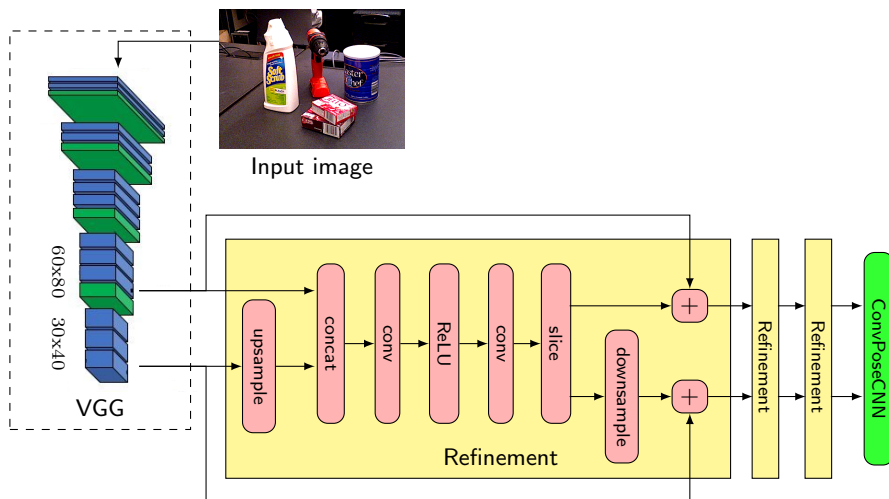


Fig. 4. Pre-prediction iterative refinement of the features extracted from the VGG network. The refined features are then fed into the ConvPoseCNN head networks. Note that there is only one set of weights for the refinement module, i.e. it is applied iteratively.

3 Evaluation

3.1 Dataset

We evaluate our method on the YCB-Video dataset [22]. The dataset contains 133,936 images of VGA resolution (640×480) extracted from 92 video sequences. Each image contains a varying number of objects selected from a larger set of 21 objects, some of which exhibit symmetry or are texture-poor. The first frame of each of the 92 video sequences was manually annotated with ground truth pose information, while the rest of the frames were automatically annotated by tracking the camera pose via SLAM. In each sequence, the objects are arranged in various spatial configurations resulting in varying degrees of occlusions making it a challenging dataset. High quality 3D models and downsampled point models containing 2620 points each are made available. The real images are supplemented with simple synthetic renderings of the object models.

3.2 Training procedure

We implemented ConvPoseCNN using PyTorch [14] framework. Except for the novel dense orientation estimation branch, we based our implementation on the openly available PoseCNN implementation¹. The openly available official implementation has minor changes compared to the model described by [22]. We

¹ <https://github.com/yuxng/PoseCNN>

noted that these minor design choices were helpful and incorporated them in our implementation as well. Additionally, we implemented Hough voting layer—non differentiable layer that computes inlier pixels—using Numba [9]. Although Numba is CPU only, the Hough Voting layer implementation with Numba is faster than other GPU implementations.

The segmentation and translation branches of ConvPoseCNN are trained with the standard pixelwise negative log-likelihood (NLL) and L2 loss respectively. The depth component of the translation branch is of a smaller scale compared to the other two components. To balance this discrepancy we scale the depth component loss by 100.

We use the ShapeMatch loss (SMLoss) proposed by [22] to train the orientation branch of ConvPoseCNN. SMLoss handles objects with and without symmetry using two different loss definitions as follows.

$$\text{SMLoss}(\tilde{q}, q) = \begin{cases} \text{SLoss}(\tilde{q}, q), & \text{if object is symmetric,} \\ \text{PLoss}(\tilde{q}, q), & \text{otherwise.} \end{cases}, \quad (2)$$

Given a set of 3D points \mathbb{M} , where $m = |\mathbb{M}|$ and $R(q)$ and $R(\tilde{q})$ are the rotation matrices corresponding to ground truth and estimated quaternion, respectively, and PLoss and SLoss are defined as follows:

$$\text{PLoss}(\tilde{q}, q) = \frac{1}{2m} \sum_{x \in \mathbb{M}} \|R(\tilde{q})x - R(q)x\|^2, \quad (3)$$

$$\text{SLoss}(\tilde{q}, q) = \frac{1}{2m} \sum_{x_1 \in \mathbb{M}} \min_{x_2 \in \mathbb{M}} \|R(\tilde{q})x_1 - R(q)x_2\|^2. \quad (4)$$

Similar to the ICP objective, SLoss does not penalize rotations of symmetric objects that lead to equivalent shapes.

During the training phase, computing SMLoss per pixels is computationally infeasible. Thus, we resort to aggregate dense predictions for each object before computing loss functions. We experimented with the aggregation mechanisms discussed in Section 2.1 and observed poor convergence. We hypothesize that this could be because of weighting quaternions with their norm before aggregation results in pixels with smaller quaternion prediction norm receiving smaller gradients. Empirically, we found that using simple numerical averaging to arrive at \tilde{q} alleviates the issue of uneven gradient distribution and contributes to convergence of the training process. Additionally, numerical averaging is computationally less expensive.

Alternatively, we also experimented with training the orientation branch with pixel-wise L2 loss and QLoss [1]

For two quaternions \tilde{q} and q it is defined as:

$$\text{QLoss}(\tilde{q}, q) = \log(\epsilon + 1 - |\tilde{q} \cdot q|), \quad (5)$$

where ϵ is introduced for stability. QLoss is designed to handle the quaternion symmetry.

The final loss function used during training is, similarly to PoseCNN, a linear combination of segmentation, translation, and orientation loss:

$$L = \alpha_{\text{seg}}L_{\text{seg}} + \alpha_{\text{trans}}L_{\text{trans}} + \alpha_{\text{rot}}L_{\text{rot}}. \quad (6)$$

where α_{seg} , α_{trans} , are set to 1. α_{rot} is set to 1 and 100 in the case of L2 loss and QLoss, and SMLoss respectively. We train ConvPoseCNN model using SGD with learning rate 0.001 and momentum 0.9.

3.3 Evaluation Metrics

We report area under the accuracy curve (AUC) metrics AUC-P and AUC-S for varying area threshold between 0 and 0.1m on ADD and ADD-S metrics as introduced along with YCB-Video Dataset [22]. ADD is average distance between corresponding points of the 3D object model in predicted and ground truth pose. ADD-S is the average of distance between each 3D point in predicted pose to the closest point in ground truth pose. ADD-S penalizes objects with symmetry less than ADD metric.

3.4 Results

Prediction Averaging To aggregate the dense pixel-wise predictions into a single orientation estimate, we use weighted quaternion averaging [12]. In the case of ConvPoseCNN, there are two possible sources of the pixel-wise weighting: segmentation score, and predicted quaternion norm. In Table 1, we show the comparison between the two weighting schemes. The norm weighting showed better results than both no averaging and using segmentation score as weighting. This suggests the predictions with smaller norms are less precise. Encouraged by this observation, we experimented further with pruning the predictions before aggregation. We sorted the predictions based on the norm and pruned varying percentile number (λ) of them.

Table 2 shows results of pruning with percentile ranging from 0 (no pruning) to 1 (extreme case of discarding all but one prediction). Pruning improves the results by a small factor overall but considerably for the objects with symmetry. This can be explained by the fact that averaging shape-equivalent orientations might result in a non-equivalent orientation and thus averaging schemes are not suitable for handling objects with symmetry.

Table 1. Weighting strategies for ConvPoseCNN L2

| Method | 6D pose ¹ | | Rotation only | |
|----------------------|----------------------|--------------|---------------|--------------|
| | AUC P | AUC S | AUC P | AUC S |
| PoseCNN ² | 53.71 | 76.12 | 78.87 | 93.16 |
| unit weights | 56.59 | 78.86 | 72.87 | 90.68 |
| norm weights | 57.13 | 79.01 | 73.84 | 91.02 |
| segm. weights | 56.63 | 78.87 | 72.95 | 90.71 |

¹ Following Xiang et al. [22].² Calculated from the published PoseCNN model.

Source: Capellen., Schwarz., and Behnke. [3].

Table 2. Quaternion pruning for ConvPoseCNN L2

| Method | 6D pose ¹ | | Rotation only | |
|--------------|----------------------|--------------|---------------|--------------|
| | AUC P | AUC S | AUC P | AUC S |
| PoseCNN | 53.71 | 76.12 | 78.87 | 93.16 |
| pruned(0) | 57.13 | 79.01 | 73.84 | 91.02 |
| pruned(0.5) | 57.43 | 79.14 | 74.43 | 91.33 |
| pruned(0.75) | 57.43 | 79.19 | 74.48 | 91.45 |
| pruned(0.9) | 57.37 | 79.23 | 74.41 | 91.50 |
| pruned(0.95) | 57.39 | 79.21 | 74.45 | 91.50 |
| single | 57.11 | 79.22 | 74.00 | 91.46 |

¹ Following Xiang et al. [22].

Source: Capellen., Schwarz., and Behnke. [3].

Prediction Clustering As an alternative to averaging schemes we experimented with RANSAC-based clustering schemes where we chose a quaternion at random and cluster the other quaternions into inliers and outliers based on the angular distance between corresponding rotations as the threshold. We repeat the process 50 times and select the quaternion prediction with the maximum inlier count. As opposed to the L2 distance in quaternion space, angular distance function invariant to the antipodal symmetry of the quaternion orientation representation. The results are shown in Table 3. Similar to averaging schemes, weighted variant of RANSAC performs better than non-weighted variants. Overall, clustering schemes outperform averaging schemes slightly on AUC S metric but perform slightly worse on AUC P. This is expected as the clustering schemes can handle object symmetries well.

Table 3. Clustering strategies for ConvPoseCNN L2

| Method | 6D pose | | Rotation only | |
|----------------|--------------|--------------|---------------|--------------|
| | AUC P | AUC S | AUC P | AUC S |
| PoseCNN | 53.71 | 76.12 | 78.87 | 93.16 |
| RANSAC(0.1) | 57.18 | 79.16 | 74.12 | 91.37 |
| RANSAC(0.2) | 57.36 | 79.20 | 74.40 | 91.45 |
| RANSAC(0.3) | 57.27 | 79.20 | 74.13 | 91.35 |
| RANSAC(0.4) | 57.00 | 79.13 | 73.55 | 91.14 |
| W-RANSAC(0.1) | 57.27 | 79.20 | 74.29 | 91.45 |
| W-RANSAC(0.2) | 57.42 | 79.26 | 74.53 | 91.56 |
| W-RANSAC(0.3) | 57.38 | 79.24 | 74.36 | 91.46 |
| pruned(0.75) | 57.43 | 79.19 | 74.48 | 91.45 |
| most confident | 57.11 | 79.22 | 74.00 | 91.46 |

RANSAC uses unit weights, while W-RANSAC is weighted by quaternion norm. PoseCNN and the best performing averaging methods are shown for comparison. Numbers in parentheses describe the clustering threshold in radians. Source: Capellen., Schwarz., and Behnke. [3].

Table 4. Results for ConvPoseCNN Shape

| | 6D Pose | | Rotation only | |
|---------------|--------------|--------------|---------------|--------------|
| | AUC P | AUC S | AUC P | AUC S |
| PoseCNN | 53.71 | 76.12 | 78.87 | 93.16 |
| average | 54.27 | 78.94 | 70.02 | 90.91 |
| norm weighted | 55.54 | 79.27 | 72.15 | 91.55 |
| pruned(0.5) | 55.33 | 79.29 | 71.82 | 91.45 |
| pruned(0.75) | 54.62 | 79.09 | 70.56 | 91.00 |
| pruned(0.85) | 53.86 | 78.85 | 69.34 | 90.57 |
| pruned(0.9) | 53.23 | 78.66 | 68.37 | 90.25 |
| RANSAC(0.2) | 49.44 | 77.65 | 63.09 | 88.73 |
| RANSAC(0.3) | 50.47 | 77.92 | 64.53 | 89.18 |
| RANSAC(0.4) | 51.19 | 78.09 | 65.61 | 89.50 |
| W-RANSAC(0.2) | 49.56 | 77.73 | 63.33 | 88.85 |
| W-RANSAC(0.3) | 50.54 | 77.91 | 64.78 | 89.21 |
| W-RANSAC(0.4) | 51.33 | 78.13 | 65.94 | 89.56 |

Table from Capellen., Schwarz., and Behnke. [3].

Loss Variants The choice of aggregation method did not have a big impact on the models trained with QLoss and thus we show only the results for Shape

variant in Table 4. Among the averaging methods, norm weight improves the result, whereas pruning does not. This suggests that there are less-confident but important predictions with higher distance from the mean and removing them significantly affects the average. This could be an effect of training with the average quaternion, where such behavior is not discouraged. Both RANSAC variants—with and without weighting—resulted in comparatively worse results. We conclude that the pixel-wise losses obtain superior performance, and average-before-loss scheme is not advantageous. Also, a fast dense version of SMLoss would need to be found in order to apply it in our architecture.

ConvPoseCNN Final Results We start the discussion about ConvPoseCNN with the qualitative result shown in Fig. 5. We visualize the 3D ground truth and predictions for all the objects in the input scene as well as orientation error and predicted orientation norm per pixel. Dense pixel-wise orientation prediction makes it easier to visualize error at each pixel and to analyze them closely. A major observation from the visualizations is that the pixels in the feature-rich regions—close to object boundaries or distinctive textures—have lower orientation error while the pixels in the feature-poor regions exhibit higher angular error. A similar phenomenon is also observed in the prediction norm visualization. The pixels in feature-rich regions have higher norm orientation predictions while the pixels in feature-poor regions have lower norm orientation predictions. We hypothesize that in feature-rich regions, the network is confident of the predictions and thus the predictions are encouraged on one specific direction, whereas in the feature-poor regions, the predictions are pulled towards various possible directions resulting in predictions with a smaller norm.

In Table 5, we report the quantitative results of ConvPoseCNN models trained with three different loss functions—L2 and QLoss, and Shape— and compare it with the PoseCNN baseline model provided in the YCB-Video Toolbox².

We provide AUC P and AUC S metric for all models including results from our own implementation of PoseCNN model in Table 5. All the three variants of ConvPoseCNN perform slightly better than PoseCNN on both AUC P and AUC S metrics. Moreover, the ConvPoseCNN variant trained with L2 yields the best results among the ConvPoseCNN variants. QLoss variant performed comparative to L2 variant on AUC P metric, whereas Shape variant performed comparative to L2 variant on AUC S loss.

Moreover, to understand the influence of translation and orientation components on the overall AUC P and AUC S metric we report AUC P and AUC S metric computed for rotation only and translation error (computed in Meters) separately. Although all the ConvPoseCNN variants perform slightly better than PoseCNN on the AUC metrics, only QLoss variant performs better than PoseCNN on the rotation only AUC metrics. Analyzing the translation error suggests that the translation estimate influences the AUC losses more than

² https://github.com/yuxng/YCB_Video_toolbox

Table 5. 6D pose, translation, rotation, and segmentation results

| | | 6D pose | | Rotation only | | NonSymC | SymC | Transl. | Segm. |
|--------------|----------------------|---------|-------|---------------|-------|---------|-------|-----------|--------|
| | | AUC P | AUC S | AUC P | AUC S | AUC P | AUC S | Error [m] | IoU |
| full network | PoseCNN | 53.71 | 76.12 | 78.87 | 93.16 | 60.49 | 63.28 | 0.0520 | 0.8369 |
| | PoseCNN ¹ | 53.29 | 78.31 | 69.00 | 90.49 | 60.91 | 57.91 | 0.0465 | 0.8071 |
| | ours, QLoss | 57.16 | 77.08 | 80.51 | 93.35 | 64.75 | 53.95 | 0.0565 | 0.7725 |
| | ours, Shape | 55.54 | 79.27 | 72.15 | 91.55 | 62.77 | 56.42 | 0.0455 | 0.8038 |
| | ours, L2 | 57.42 | 79.26 | 74.53 | 91.56 | 63.48 | 58.85 | 0.0411 | 0.8044 |
| GT segm. | PoseCNN ¹ | 52.90 | 80.11 | 69.60 | 91.63 | 76.63 | 84.15 | 0.0345 | 1 |
| | ours, QLoss | 57.73 | 79.04 | 81.20 | 94.52 | 88.27 | 90.14 | 0.0386 | 1 |
| | ours, Shape | 56.27 | 81.27 | 72.53 | 92.27 | 77.32 | 89.06 | 0.0316 | 1 |
| | ours, L2 | 59.50 | 81.54 | 76.37 | 92.32 | 80.67 | 85.52 | 0.0314 | 1 |

The average translation error, the segmentation IoU and the AUC metrics for different models. The AUC results were achieved using weighted RANSAC(0.1) for ConvPoseCNN QLoss, Markley’s norm weighted average for ConvPoseCNN Shape and weighted RANSAC(0.2) for ConvPoseCNN L2. *GT segm.* refers to ground truth segmentation (i.e. only pose estimation). Source: Capellen., Schwarz., and Behnke. [3].

¹ Our own reimplemention.

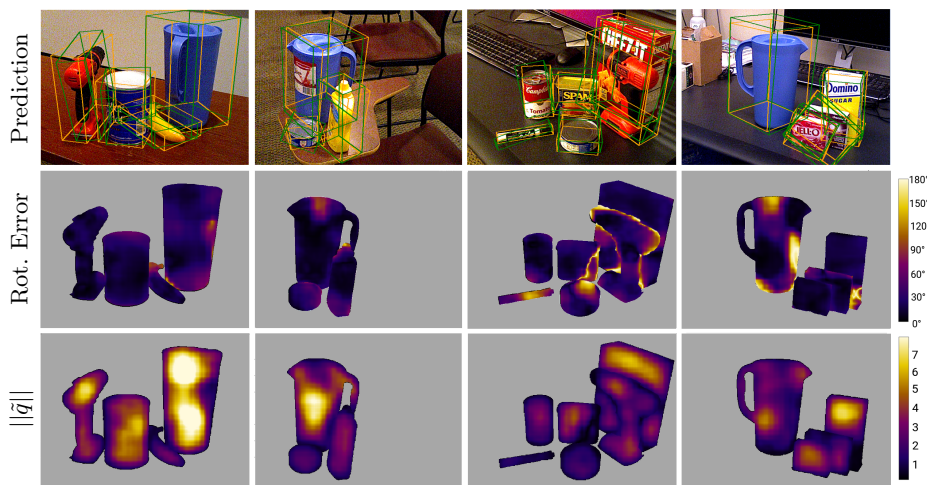


Fig. 5. Qualitative results from ConvPoseCNN L2 on the YCB-Video test set. Top: The orange boxes show the ground truth bounding boxes, the green boxes the 6D pose prediction. Middle: Angular error of the dense quaternion prediction \tilde{q} w.r.t. ground truth, masked by ground truth segmentation. Bottom: Quaternion prediction norm $\|\tilde{q}\|$ before normalization. This measure is used for weighted aggregation. Note that the prediction norm is low in high-error regions and high in regions that are far from occlusions and feature-rich. Source: Capellen., Schwarz., and Behnke. [3].

the orientation estimate. However, the models that achieve better translation estimation, performs worse with the orientation estimate.

Furthermore, to analyze the performance of the models on objects with and without symmetry we report the average per-class AUC P metric for objects without symmetry and average per-class AUC S for objects with symmetry. ConvPoseCNN performed a bit better than PoseCNN for the objects without symmetry but worse for the ones with symmetry. This can be explained by the use loss functions—QLoss and L2 loss—that are not designed to handle symmetry. But, surprisingly, the model trained with SMLoss also performs worse for the symmetric objects compared to PoseCNN.

This might be due to different reasons: First, we utilize an average before calculating the loss; therefore during training the average might penalize predicting different shape-equivalent quaternions, in case their average is not shape-equivalent. Secondly, there are only five symmetric objects in the dataset and we noticed that two of those, the two clamp objects, are very similar and thus challenging, not only for the orientation but as well for the segmentation and vertex prediction. This is further complicated by a difference in object coordinate systems for these two objects.

While aggregating the dense pixel-wise orientation predictions to a single orientation prediction per-class, we use segmentation results. Thus, the segmentation results also influence the final metrics. To quantify the influence of segmentation results we report metrics for the all the five models—three ConvPoseCNN, and two PoseCNN variants—using the ground truth segmentation as well. Using ground truth segmentation improves translation and orientation for all the models. Hu et al. [6] also report a similar observation.

Table 6. Comparison to Related Work

| | Total | | Average |
|---------------------|-------|-------|------------------|
| | AUC P | AUC S | AUC ¹ |
| PoseCNN | 53.7 | 75.9 | 61.30 |
| ConvPoseCNN L2 | 57.4 | 79.2 | 62.40 |
| HeatMaps without FM | | | 61.41 |
| ConvPoseCNN+FM | 58.22 | 79.55 | 61.59 |
| HeatMaps with FM | | | 72.79 |

Comparison between PoseCNN (as reported by Xiang et al. [22]), ConvPoseCNN L2 with pruned(0.75), and HeatMaps [13] without and with Feature Mapping (FM). Source: Capellen., Schwarz., and Behnke. [3].

¹ As defined by Oberweger, Rad, and Lepetit [13].

Table 7. Detailed Class-wise Results

| Class | Ours | | PoseCNN | |
|-------------------|-------|-------|---------|-------|
| | AUC P | AUC S | AUC P | AUC S |
| master_chef_can | 62.32 | 89.55 | 50.08 | 83.72 |
| cracker_box | 66.69 | 83.78 | 52.94 | 76.56 |
| sugar_box | 67.19 | 82.51 | 68.33 | 83.95 |
| tomato_soup_can | 75.52 | 88.05 | 66.11 | 80.90 |
| mustard_bottle | 83.79 | 92.59 | 80.84 | 90.64 |
| tuna_fish_can | 60.98 | 83.67 | 70.56 | 88.05 |
| pudding_box | 62.17 | 76.31 | 62.22 | 78.72 |
| gelatin_box | 83.84 | 92.92 | 74.86 | 85.73 |
| potted_meat_can | 65.86 | 85.92 | 59.40 | 79.51 |
| banana | 37.74 | 76.30 | 72.16 | 86.24 |
| pitcher_base | 62.19 | 84.63 | 53.11 | 78.08 |
| bleach_cleanser | 55.14 | 76.92 | 50.22 | 72.81 |
| bowl | 3.55 | 66.41 | 3.09 | 70.31 |
| mug | 45.83 | 72.05 | 58.39 | 78.22 |
| power_drill | 76.47 | 88.26 | 55.21 | 72.91 |
| wood_block | 0.12 | 25.90 | 26.19 | 62.43 |
| scissors | 56.42 | 79.01 | 35.27 | 57.48 |
| large_marker | 55.26 | 70.19 | 58.11 | 70.98 |
| large_clamp | 29.73 | 58.21 | 24.47 | 51.05 |
| extra_large_clamp | 21.99 | 54.43 | 15.97 | 46.15 |
| foam_brick | 51.80 | 88.02 | 39.90 | 86.46 |

Source: Capellen., Schwarz., and Behnke. [3].

Comparison to Related Work In Table 6, we show the comparisons between ConvPoseCNN, PoseCNN, and HeatMaps [13] approaches. Oberweger, Rad, and Lepetit [13] report class-wise area under the accuracy curve metric (AUC) instead of AUC P and AUC S metrics. To make the methods comparable, we provide AUC for both ConvPoseCNN and PoseCNN. [13] proposed Feature Mapping (FM) technique that significantly improves their results. Without feature mapping, we perform slightly better than both PoseCNN and HeatMaps. However, the difference is negligible considering the variations due to the choice of hyperparameters and minor implementations details. Detailed class-wise AUC metrics for both the best performing ConvPoseCNN and PoseCNN models are shown in 7.

We also investigated applying the Feature Mapping technique [13] to our model. Following the process, we render synthetic images with poses corresponding to the real training data. We selected the features from backbone VGG-16 for the mapping process and thus have to transfer two feature maps with 512 features each. We replaced the fullyconnected network architecture for feature mapping as done by [13], with a convolutional set-up and mapped the feature from the different stages to each other with residual blocks based on (1×1) convolutions. The results are presented in 6. However, we did not observe the large gains reported by [13] for our architecture. We hypothesize that the feature mapping technique is highly dependent on the quality and distribution of the rendered synthetic images, which are maybe not of sufficient quality in our case.

Table 8. Training performance & model sizes

| | Iterations/s ¹ | Model size |
|--------------------|---------------------------|------------|
| PoseCNN | 1.18 | 1.1 GiB |
| ConvPoseCNN L2 | 2.09 | 308.9 MiB |
| ConvPoseCNN QLoss | 2.09 | 308.9 MiB |
| ConvPoseCNN SMLoss | 1.99 | 308.9 MiB |

¹ Using a batch size of 2. Averaged over 400 iterations.
Source: Capellen., Schwarz., and Behnke. [3].

Time Comparisons We used NVIDIA GTX 1080 Ti GPU with 11 GB of memory to benchmark the training and inference time for ConvPoseCNN and PoseCNN models. In table 8 we report number of iterations per second. All the variants of ConvPoseCNN are significantly faster. Additionally, size of the saved ConvPoseCNN models are significantly smaller compared to the PoseCNN models.

Unfortunately, this advantage in speed during the training process is not observed during the inference as shown in 9. Averaging methods, on average, consume time comparable to the PoseCNN. But the RANSAC based clustering

methods more time consuming; the forward pass of ConvPoseCNN takes about 65.5 ms, the Hough transform around 68.6 ms. We attribute the comparable inference time consumption to the highly optimized ROI pooling layers in the modern deep learning frameworks.

Table 9. Inference timings

| Method | Time [ms] ¹ | Aggregation [ms] |
|----------------------|------------------------|------------------|
| PoseCNN ² | 141.71 | |
| ConvPoseCNN | | |
| - naive average | 136.96 | 2.34 |
| - average | 146.70 | 12.61 |
| - weighted average | 146.92 | 13.00 |
| - pruned w. average | 148.61 | 14.64 |
| - RANSAC | 158.66 | 24.97 |
| - w. RANSAC | 563.16 | 65.82 |

¹ Single frame, includes aggregation.

² Xiang et al. [22].

Source: Capellen., Schwarz., and Behnke. [3].

Iterative Refinement Post-prediction iterative refinement module is trained with segmentation, translation, and orientation estimates from ConvPoseCNN as well as VGG16 features as input. At each iteration, the model refines segmentation, translation, and orientation estimates. VGG16 features provide contextual information about the input scene. We experimented with varying number of refinement steps. Similar to ConvPoseCNN, we used same combined loss function as discussed in Section 3.2. But, we observed both training and validation loss plateauing very early on the training process and the resulting model also performed worse quantitatively compared to ConvPoseCNN on the test set.

This could be because the estimates are in a form that is not a suitable for a simple three layer network. Exploring complex architectures is not an option for us since we focus on keeping the overhead of iterative refinement minimal.

In contrast to the post-prediction refinement, pre-prediction refinement not only performed well during but also improved the AUC metrics on the test set. This suggests that in the case of ConvPoseCNN, refining the features at an early stage helps the network in learning representations better suitable for pose estimation. We trained the refinement module with a various number of iterations and in Fig. 6, we present the AUC metrics achieved by various number of refinement iterations. Overall, the iterative refinement improves the prediction and different number of iterations results in slightly different AUC metrics. Interestingly, the performance peaks at three iterations. If there are any gains with more iterations, they are not significant. We attribute this fact to the small depth of

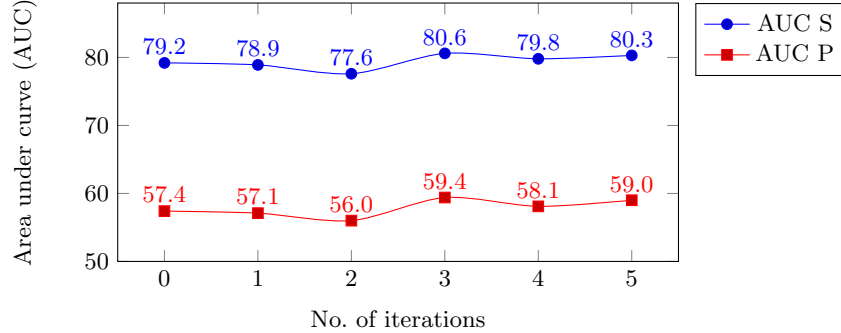


Fig. 6. Results of pre-prediction feature refinement process for various number of iterations. The variant with zero iterations corresponds to ConvPoseCNN without any refinement (Table 6).

Table 10. Class-wise Results ConvPoseCNN without refinement, with three iterations of refinement and with five iterations of refinement.

| Class | PoseCNN | | IR 3 | | IR 5 | |
|-------------------|---------|-------|-------|-------|-------|-------|
| | AUC P | AUC S | AUC P | AUC S | AUC P | AUC S |
| master_chef_can | 62.32 | 89.55 | 62.69 | 90.93 | 61.58 | 91.09 |
| cracker_box | 66.69 | 83.78 | 51.64 | 79.02 | 62.48 | 82.53 |
| sugar_box | 67.19 | 82.51 | 63.16 | 80.81 | 68.95 | 84.16 |
| tomato_soup_can | 75.52 | 88.05 | 78.70 | 90.70 | 75.12 | 88.65 |
| mustard_bottle | 83.79 | 92.59 | 83.66 | 92.09 | 83.99 | 91.65 |
| tuna_fish_can | 60.98 | 83.67 | 71.10 | 88.15 | 72.68 | 90.37 |
| pudding_box | 62.17 | 76.31 | 67.72 | 84.73 | 66.11 | 83.25 |
| gelatin_box | 83.84 | 92.92 | 83.38 | 91.45 | 86.98 | 93.18 |
| potted_meat_can | 65.86 | 85.92 | 69.52 | 87.56 | 68.21 | 86.22 |
| banana | 37.74 | 76.30 | 42.96 | 70.24 | 42.75 | 70.34 |
| pitcher_base | 62.19 | 84.63 | 68.31 | 86.79 | 66.51 | 86.55 |
| bleach_cleanser | 55.14 | 76.92 | 50.86 | 71.48 | 52.28 | 75.61 |
| bowl | 3.55 | 66.41 | 7.21 | 73.04 | 8.24 | 69.29 |
| mug | 45.83 | 72.05 | 58.31 | 81.68 | 62.11 | 82.64 |
| power_drill | 76.47 | 88.26 | 73.12 | 86.57 | 71.60 | 85.68 |
| wood_block | 0.12 | 25.90 | 0.785 | 27.70 | 1.07 | 31.86 |
| scissors | 56.42 | 79.01 | 62.41 | 80.96 | 51.22 | 75.56 |
| large_marker | 55.26 | 70.19 | 64.16 | 76.35 | 60.15 | 71.96 |
| large_clamp | 29.73 | 58.21 | 35.66 | 62.34 | 33.14 | 61.93 |
| extra_large_clamp | 21.99 | 54.43 | 23.16 | 55.74 | 23.91 | 55.91 |
| foam_brick | 51.80 | 88.02 | 51.31 | 88.62 | 47.69 | 87.08 |

our refinement network which limits the operations it can perform. In Table 10 we compare the class-wise AUC metrics for ConvPoseCNN (without refinement), three and five iterations of refinement. For most objects, the AUC metrics are improved but for some objects, there is a drop in accuracy. The maximum gain of 12.48 AUC-P and 9.63 is observed for the *mug* object while a severe drop of 15.05 AUC-P and 5.54 is observed for *cracker_box* and *bleach_cleanser* respectively—both relatively big objects, where information needs to be communicated and fused over larger regions.

4 Conclusion

We presented ConvPoseCNN, a fully convolutional architecture for object pose estimation and demonstrated that, similar to translation estimation, direct regression of the orientation estimation can be done in a dense pixel-wise manner. This helps in not only simplifying neural networks architectures for 6D object pose estimation but also reducing the size of the models and faster training. To further the performance of fully convolutional models for pose estimation, scalable dense pixel-wise loss function needs to be explored. As a next step, we plan to evaluate ConvPoseCNN on highly cluttered scenes where we expect the dense predictions to be especially beneficial, since disambiguation of close objects should be more direct than with RoI-based architectures.

Moreover, we demonstrated that the pose predictions can be refined even with a small network to boost the performance, provided the refinement is done at the right level of abstraction. In case of ConvPoseCNN, refining intermediate representations yielded better performance than post-prediction refinement. Thus, network architecture designs that imbue refinement modules should be favoured for object pose estimation. In the future, we plan to combine iterative refinement with other state-of-the-art architectures and further investigate the design of refinement modules. The key challenge is to balance the need for a larger powerful module capable of iterative refinement with keeping the processing time and memory overhead introduced by the refinement module low.

References

- [1] Gideon Billings and Matthew Johnson-Roberson. “SilhoNet: An RGB Method for 3D Object Pose Estimation and Grasp Planning”. In: *arXiv preprint arXiv:1809.06893* (2018).
- [2] Eric Brachmann et al. “Uncertainty-driven 6D pose estimation of objects and scenes from a single RGB image”. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016, pp. 3364–3372.

- [3] Catherine Capellen., Max Schwarz., and Sven Behnke. “ConvPoseCNN: Dense Convolutional 6D Object Pose Estimation”. In: *Proceedings of the 15th International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications - Volume 5: VISAPP, INSTICC*. SciTePress, 2020, pp. 162–172. ISBN: 978-989-758-402-2. DOI: 10.5220/0008990901620172.
- [4] Thanh-Toan Do et al. “Deep-6DPose: Recovering 6D Object Pose from a Single RGB Image”. In: *European Conference on Computer Vision (ECCV)*. 2018.
- [5] Ross Girshick. “Fast R-CNN”. In: *IEEE International Conference on Computer Vision (ICCV)*. 2015, pp. 1440–1448.
- [6] Yinlin Hu et al. “Segmentation-driven 6d object pose estimation”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2019, pp. 3385–3394.
- [7] Eric R Kandel et al. *Principles of neural science*. Vol. 4. McGraw-hill New York, 2000.
- [8] Alexander Krull et al. “Learning analysis-by-synthesis for 6D pose estimation in RGB-D images”. In: *International Conference on Computer Vision (ICCV)*. 2015, pp. 954–962.
- [9] Siu Kwan Lam, Antoine Pitrou, and Stanley Seibert. “Numba: A LLVM-based Python JIT Compiler”. In: *Second Workshop on the LLVM Compiler Infrastructure in HPC*. Austin, Texas: ACM, 2015.
- [10] Yi Li et al. “DeepIM: Deep Iterative Matching for 6D Pose Estimation”. In: *European Conference on Computer Vision (ECCV)*. 2018.
- [11] Tsung-Yi Lin et al. “Focal loss for dense object detection”. In: *IEEE International Conference on Computer Vision (ICCV)*. 2017, pp. 2980–2988.
- [12] F Landis Markley et al. “Averaging quaternions”. In: *Journal of Guidance, Control, and Dynamics* 30.4 (2007), pp. 1193–1197.
- [13] Markus Oberweger, Mahdi Rad, and Vincent Lepetit. “Making Deep Heatmaps Robust to Partial Occlusions for 3D Object Pose Estimation”. In: *European Conference on Computer Vision (ECCV)*. 2018, pp. 125–141.
- [14] Adam Paszke et al. “Automatic differentiation in PyTorch”. In: *NIPS*. 2017.
- [15] Arul Selvam Periyasamy, Max Schwarz, and Sven Behnke. “Refining 6D Object Pose Predictions using Abstract Render-and-Compare”. In: *2019 IEEE-RAS 19th International Conference on Humanoid Robots (Humanoids)*. IEEE. 2019, pp. 739–746.
- [16] Mahdi Rad and Vincent Lepetit. “BB8: A Scalable, Accurate, Robust to Partial Occlusion Method for Predicting the 3D Poses of Challenging Objects without Using Depth”. In: *International Conference on Computer Vision (ICCV)*. 2017.
- [17] Mahdi Rad, Markus Oberweger, and Vincent Lepetit. “Feature Mapping for Learning Fast and Accurate 3D Pose Inference from Synthetic Images”. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2018.

- [18] Joseph Redmon et al. “You only look once: Unified, real-time object detection”. In: *Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016, pp. 779–788.
- [19] Shaoqing Ren et al. “Faster R-CNN: Towards real-time object detection with region proposal networks”. In: *Advances in neural information processing systems*. 2015, pp. 91–99.
- [20] Karen Simonyan and Andrew Zisserman. “Very deep convolutional networks for large-scale image recognition”. In: *arXiv preprint arXiv:1409.1556* (2014).
- [21] Bugra Tekin, Sudipta N. Sinha, and Pascal Fua. “Real-Time Seamless Single Shot 6D Object Pose Prediction”. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2018.
- [22] Yu Xiang et al. “PoseCNN: A Convolutional Neural Network for 6D Object Pose Estimation in Cluttered Scenes”. In: *Robotics: Science and Systems (RSS)*. 2018.