

PoET: Pose Estimation Transformer for Single-View, Multi-Object 6D Pose Estimation

Thomas Jantos

Control of Networked Systems Group
University of Klagenfurt, Austria
thomas.jantos@aau.at

Mohamed Amin Hammad

Infineon Technologies Austria AG
Villach, Austria
mohamedamin.hammad@infineon.com

Wolfgang Granig

Infineon Technologies Austria AG
Villach, Austria
wolfgang.granig@infineon.com

Stephan Weiss

Control of Networked Systems Group
University of Klagenfurt, Austria
stephan.weiss@aau.at

Jan Steinbrener

Control of Networked Systems Group
University of Klagenfurt, Austria
jan.steinbrener@aau.at

Abstract: Accurate 6D object pose estimation is an important task for a variety of robotic applications such as grasping or localization. It is a challenging task due to object symmetries, clutter and occlusion, but it becomes more challenging when additional information, such as depth and 3D models, is not provided. We present a transformer-based approach that takes an RGB image as input and predicts a 6D pose for each object in the image. Besides the image, our network does not require any additional information such as depth maps or 3D object models. First, the image is passed through an object detector to generate feature maps and to detect objects. Then, the feature maps are fed into a transformer with the detected bounding boxes as additional information. Afterwards, the output object queries are processed by a separate translation and rotation head. We achieve state-of-the-art results for RGB-only approaches on the challenging YCB-V dataset. We illustrate the suitability of the resulting model as pose sensor for a 6-DoF state estimation task. Code is available at <https://github.com/aau-cns/poet>.

Keywords: 6D Pose Estimation, Transformer, Object-Relative Localization

1 Introduction

Accurately estimating the 6D pose of objects from RGB images is essential for robotics tasks such as grasping or localization [1, 2]. Grasping tasks require the robot to know the exact position of the object such that it can place its end-effector effectively. In autonomous driving, it is critical that the vehicle has sufficient knowledge of its surroundings including the relative 6D pose of all objects in its vicinity. For unmanned aerial vehicle (UAV) navigation, especially in close proximity to people or infrastructure, realizing precise control depends on the estimation of 6D object poses. In recent years, vision-based 6D pose estimation with deep learning [3] has been on the rise. Approaches differ in terms of input data, network architecture, post processing and number of viewpoints [4, 5, 6, 7]. Observing objects from multiple viewpoints introduces constraints to the pose of objects and improves estimation [7, 8]. Availability of 3D object models allows for an iterative refinement of an initial pose estimate by either iterative closest point (ICP) [9] matching of pointclouds or by matching keypoints with a perspective-n-point (PnP) [10] algorithm. However, these algorithms are computationally demanding. Besides being used for post processing, 3D models can be utilized as an additional input to the network [5]. This may include model keypoints and corresponding features or information about the object such as symmetry axes and planes. Additionally, prior information

about the object class or a depth map corresponding to the input RGB image, which improves the networks ability to estimate the objects’ distance to the camera [8], can be passed to the network. Although additional input information can greatly benefit the accuracy of the final estimated pose, apart from requiring a more detailed data base containing accurate 3D models and depth maps corresponding to RGB images, it results in higher computational complexity and thus, longer run time in comparison to the same neural network-based architecture taking only RGB images [8].

While the requirements for real-time performance depend on the specific application, more often the availability of high-quality, detailed 3D models or depth maps is limited. Depending on the type of additional input information, additional sensor hardware is also required, which may not be available during inference. In this work, we focus on a pose estimation framework purely based on RGB images and information provided by a backbone object detector. To the best of our knowledge, we are the first to incorporate global image context information into the pose estimation task by passing multi-scale feature maps to a transformer network and relying only on 2D image information. Our approach does not depend on the number of objects present. Our framework, dubbed PoET (Pose Estimation Transformer), can be used on top of any 2D object detector. We evaluate our approach on the YCB-V [3] dataset and compare our results to state-of-the-art approaches. Finally, we illustrate the suitability of the obtained model as a pose sensor for a 6-DoF state estimation task, where ”pose sensor” means the combination of a camera and PoET providing information about the camera pose relative to objects, i.e. its relative position and orientation. Our contributions are the following:

- We present a transformer-based framework that takes a single RGB-image as input, estimates the 6D pose for every object present in the image and can be trained on top of any object detector framework. A detailed ablation study supports our design choices.
- The framework is independent of any additional information which is not contained in the raw RGB image. In particular, it does not depend on depth maps, object symmetries or 3D object models. Hence, our results are achieved without iterative refinement and the whole network can be trained using 3D model independent loss functions.
- We achieve state-of-the-art results on the YCB-V [3] dataset for RGB-only methods and competitive results in comparison to approaches utilizing 3D models.
- We show the feasibility of the resulting model as a pose sensor in a 6-DoF localization task.

The rest of the paper is organized as follows: In Section 2, related work for 6D pose estimation is reviewed. Following the presentation of our method and implementation details in Section 3, the experiments and the corresponding results are discussed in Section 4 including an ablation study investigating our network architecture. Additionally, we illustrate how PoET and its relative 6D pose estimates can be used for localization in Section 4.3. Finally, the limitations are discussed in Section 5 and the paper is concluded in Section 6. We refer the reader to the supplementary material for an extensive ablation study and additional results on the LM-O [11] dataset.

2 Related Work

Classical image-based 6D pose estimation approaches can be split into feature-based methods and template-based methods. For the latter, object pose is determined by matching object templates against the input image [12, 13]. While template-based approaches work well on texture-less objects, they are prone to fail in scenarios where objects are occluded. In feature-based methods, local features are extracted from the image and then matched to the 3D model to determine correspondences [14]. Based on these 2D-3D correspondences, the 6D pose of the object can be derived. While these approaches can handle occlusion of objects, they require textures in order to perform the matching. If RGB-D images are available, the additional information provided by the depth can be used to improve the initial pose estimate by iterative refinement [15, 16] such as e.g., ICP [9]. Even though those methods can achieve state-of-the-art performance, they are computationally expensive and the availability of depth maps is not guaranteed.

In recent years, advancements in deep learning for computer vision tasks have been applied to single-view, image-based 6D pose estimation, either to replace components of classical approaches [17, 18, 19, 20], or as end-to-end learned methods, where the 6D pose is directly estimated from the input using convolutional neural networks (CNNs). Xiang et al. [3] proposed with PoseCNN, a CNN-based method to regress the 3D translation and rotation of each object present in the image using

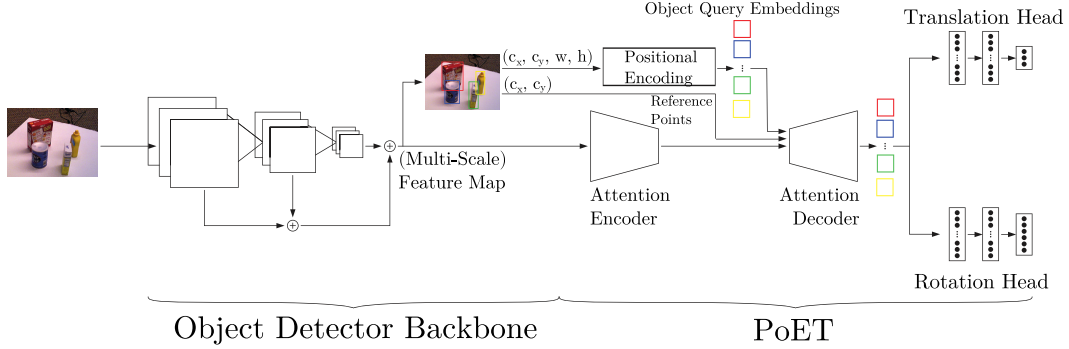


Figure 1: Overview of the PoET network architecture for single-view, multi-object 6D pose estimation. Bounding box information for each detected object is passed to the transformer as an object query. Afterwards, for each object query the egocentric 3D translation and 6D rotation [25] is predicted.

an object symmetry-aware loss function. Li et al. [8] proposed a framework that introduces prior knowledge about the object class into the network and perform pose estimation by discretizing the possible translation and rotation values to unique bins resulting in a classification task.

Aside from depth images, using 3D object models for pose estimation yielded promising results [21, 4, 5, 6, 7]. The approaches differ in terms of how the 3D model is used. Kehl et al. [21] and Li et al. [4] predict the 6D pose of objects and then refine the estimated pose. With GDR-Net, Wang et al. [22] are able to integrate PnP into an end-to-end trainable network. Similarly, Li et al. [6] use the 3D model of an object for iterative pose refinement. Given an initial pose estimate, a network is trained to iteratively refine the pose by matching a rendered image created from the 3D model to the original image. Labbé et al. [7] apply this approach to multiple viewpoints of the same scene resulting in improved pose estimation. In yet another approach, Billings and Johnson-Roberson [5] provide the 3D model as an additional input to their network, dubbed SilhoNet. Their network predicts the object silhouette and a 3D translation vector derived from the 2D bounding box position. Based on the former, the 3D rotation of the object is estimated and corrected for silhouette symmetric objects - a significant restriction, as the real error in rotation for symmetric objects can be very large.

Similar to our approach, Amini et al. [23] presented a transformer-based architecture to directly regress the 6D pose for multiple objects contained in a single image. By extending the Detection Transformer [24] with translation and rotation heads, they are able to train the whole network in an end-to-end fashion. However, they require the object 3D model as they use a symmetry aware loss [3]. Our approach does not require any additional information such as depth maps, 3D models or known object symmetries, but instead directly estimates the translation and rotation of objects in the camera coordinate frame from a single RGB image. In contrast to other methods that work with regions of interest for pose estimation, we keep the complete image feature map and provide the regions of interest as an additional input to our transformer network.

3 Method

We present a novel transformer-based neural network for the 6D pose estimation task. Taking a single RGB image as its input, the 6D pose of every object detected in the image is predicted simultaneously. After generating (multi-scale) feature maps by passing the image through an object detector, they are processed by a transformer architecture. At the end, translation and rotation are estimated in a decoupled manner. In this section, we first present the general structure of our network. Afterwards, we talk about specific implementation details and data preparation.

3.1 Network Architecture

Fig. 1 shows a detailed overview of our network architecture, which consists of three steps. First, the input image is passed through a backbone object detector network. Both the generated (multi-scale) feature maps and the predicted bounding boxes are used in subsequent processing. PoET can be trained on top of any object detector architecture and thus extend a pre-trained object detection framework to include 6D object pose estimation. Second, the (multi-scale) feature maps used for the object detection step are passed to the encoder of a multi-head attention-based transformer.

Our transformer architecture is a modified version of the *Deformable DETR* transformer module proposed by Zhu et al. [26]. In contrast to the original *DETR* [24], the deformable transformer allows to process multi-scale feature maps similar to state-of-the-art object detectors. By only attending to a limited number of feature map keypoints in the decoder, the transformer achieves a higher pass-through rate. Additionally, a deformable transformer shows faster convergence rates than a regular transformer architecture. The main idea behind using a transformer architecture for feature map refinement is to generate features that capture the global information contained in the image. For example, such additional information might be the image location of other objects present or general information of the overall scene extracted from the image.

While the encoder of the deformable transformer is kept unchanged, we modified the decoder to incorporate more information from the object detection step: First, the learned object query embedding is replaced by bounding box information. For each detected object, the bounding box center coordinates (c_x, c_y) , the width w and the height h are normalized and then position-encoded [27] to generate the object query embeddings. The embedding dimension L is chosen such that $2 \cdot n_p \cdot L$ equals the hidden dimension d_h of the transformer. In our case, the number of parameters n_p equals to 4. Moreover, the inter-query attention heads ensure that information is properly propagated between the different object queries. Second, the *Deformable DETR* originally only attends to a limited number of keypoints which are randomly sampled around reference points. Instead of predicting reference points from query embeddings by a trainable fully connected layer, we directly feed the normalized center coordinates (c_x, c_y) as the reference points to the decoder. By feeding this additional information to the decoder along with the encoder-refined image feature maps and the inter-query attention heads, the decoder generates new object query embeddings which not only contain local information regarding the object but also global information extracted from the image. Third, the object queries outputted by the transformer are passed through a translation and rotation head. This allows us to simultaneously estimate the pose for multiple objects independent of how many objects are present and which class they belong to. As we approach the pose estimation problem from a global image context perspective by extracting features from the whole image, the network directly estimates the translation and rotation with respect to the camera.

Our translation head is a simple multi-layer perceptron (MLP) with input dimension d_h , one hidden layer and output dimension 3. We directly predict the translation $\tilde{t} = (\tilde{t}_x, \tilde{t}_y, \tilde{t}_z)$ with respect to the camera frame. Given the ground-truth translation t , our translation head is trained with a simple L2-loss defined as

$$L_t = \|t - \tilde{t}\|_2. \quad (1)$$

The rotation prediction head is identical to the translation head besides the output dimension. For estimating the rotation, we use the 6D rotation representation proposed by Zhou et al. [25] as this representation does not suffer from discontinuities with respect to learning as e.g., quaternions do. Hence, the network predicts a 6-dimensional output vector. Afterwards, the 6D representation is used to determine the estimated rotation matrix $\tilde{R} \in SO(3)$ as described in [25]. The rotation head is trained using a geodesic loss [28] given by

$$L_{rot} = \arccos \frac{1}{2} \left(\text{Tr} \left(R \tilde{R}^T \right) - 1 \right), \quad (2)$$

where R is the ground-truth rotation and $\text{Tr}(\cdot)$ is the matrix trace operator. To ensure numerical stability of the loss during training, the argument of the arccos is clamped between $-1 + \epsilon$ and $1 - \epsilon$, where $\epsilon = 1e - 6$. Our whole network is then trained with a weighted multi-task loss expressed as

$$L = \lambda_t L_t + \lambda_{rot} L_{rot}, \quad (3)$$

where the loss is calculated for each object and then averaged across all objects present across all images in the batch. λ_t and λ_{rot} are the weighting parameters for the translation and rotation loss respectively. PoET can be trained either class-specific or class-agnostic. In the class-specific case with n_{cls} different classes, the translation and rotation output dimension are changed to $3 \cdot n_{cls}$ and $6 \cdot n_{cls}$, respectively. For each object query, one hypothesis for each class is regressed and the final output is chosen depending on the class predicted for the bounding box. However, no class information is fed into the transformer.

3.2 Implementation Details & Data Preparation

While PoET can be trained on top of any object detector, we used Scaled-YOLOv4 [29] as the backbone object detector as it offers a good trade-off between speed and accuracy. An MS-

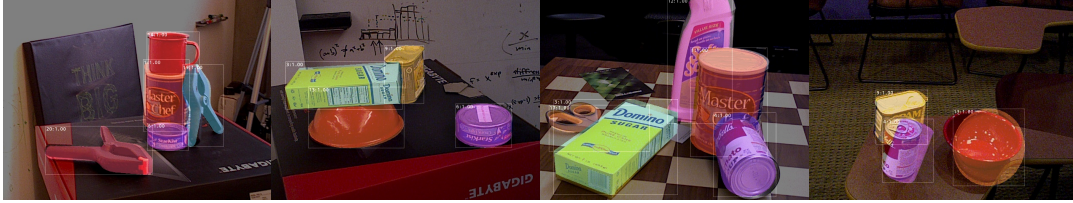


Figure 2: Qualitative results of the relative 6-DoF object poses predicted by PoET for the YCB-V dataset.

COCO [30] pre-trained Scaled-YOLOv4 is fine tuned for 10 epochs on the YCB-V dataset for the object detection task. During the training of PoET, the weights of the object detector backbone are frozen.

We implement PoET using PyTorch [31] and train it for 50 epochs using AdamW [32] with a learning rate of $2e - 5$ and a batch size of 16. Our best performing network has five encoder and decoder layers, $d_h = 256$, 16 attention heads and a positional embedding dimension of $L = 32$. The network is simultaneously trained for 3D translation and 3D rotation estimation and the weighting parameters are set to $\lambda_t = 2$ and $\lambda_{rot} = 1$, such that both losses are in the same value range.

Given the commonly used data split for YCB-V [3, 5, 6, 7, 8, 33], we train our network on 80 of the available 92 video sequences and reserve the remaining 12 sequences and their challenging keyframes for testing and evaluation. Moreover, we include 80,000 synthetic images generated from 2D projections of the 3D object models provided by the original authors [3]. We refer the reader to the supplementary material, for more details regarding our implementation.

4 Results & Experiments

In this section, we present PoET’s performance on the YCB-V benchmark dataset for 6D pose estimation. For evaluation, we use the AUC of ADD-S metric [3] and additionally report the average translation and rotation error in *cm* and degree, respectively, as done by [5]. We compare our results to state-of-the-art, single-view, RGB-based approaches and list the results of other approaches as reported in the corresponding work. We conduct an ablation study on the network architecture and our modifications to the transformer part. The ablation study investigates the influence of the object detector backbone, transformer modifications, data augmentation, network architecture and rotation representation on PoET’s performance. While we present here our most important findings, we refer the reader to our supplementary material for additional ablation results. Finally, we illustrate how PoET’s relative pose estimates can be used for camera localization.

4.1 6D Pose Estimation

Our best performing network is class-specific, consists of 5 encoder and decoder layers with 16 attention heads, and was trained according to Section 3.2. Methods that rely on an object detector to predict regions of interest (ROIs) usually do not elaborate whether and, if so, how multiple predictions, bad predictions and missing predictions are treated within their approach. Therefore, to allow for a fair comparison, we also report the results of PoET and other approaches given ground-truth bounding boxes in addition to the results on predicted bounding boxes. We present representative qualitative results in Fig. 2.

In Table 1 we report our results for the AUC of ADD-S metric. Both for predicted and ground-truth bounding boxes, PoET outperforms (overall and also for most individual classes) other state-of-the-art RGB-based methods that either work on the whole input image [3, 23] or that predict ROIs for pose regression [8]. We also outperform SilhoNet [5] which feeds additional information from the 3D model to its network and reduces predicted rotations by known object symmetries. This shows that PoET and its feature maps containing global image context reduce the need for any additional inputs. This is especially highlighted when comparing the performance of PoET to the ROI-based approaches SilhoNet and MCN [8] in the case where all networks are provided with ground-truth bounding boxes. Only models that explicitly utilize 3D object models during inference, either through PnP during the estimation [22] or by performing iterative refinement after estimating an initial pose [7, 6], achieve slightly better results but this comes at the cost of significantly increased computational complexity and the need for accurate 3D models to be known a priori.

Table 1: Comparison with state-of-the-art on YCB-V. We report the AUC of ADD-S. *gt* denotes results achieved with providing ground-truth ROIs to network. The 3D model row indicates how the object model is used: either to calculate the loss function, as an additional input, for symmetry reduction or for PnP or iterative refinement (IR) based pose estimation. (*) denotes symmetric objects. Bold and italic values indicate state-of-the-art results for methods not based on PnP/IR, using ground-truth or predicted ROIs, respectively.

Method	PoseCNN [3]	SilhoNet [5]	SilhoNet _{gt}	MCN[8]	MCN _{gt}	T6D [23]	PoET _{gt}	PoET	GDR-Net [22]	CosyPose [7]	DeepIM[6]
3D Model	Loss	Input + Sym	Input + Sym	2D	2D	Loss	2D	2D	PnP	IR	IR
master chef can	84.0	84.0	83.6	87.8	91.2	<i>91.9</i>	92.9	88.4	96.6	-	93.1
cracker box	76.9	73.5	88.4	64.3	78.5	<i>86.6</i>	90.4	80.5	84.9	-	91.0
sugar box	84.3	86.6	88.8	82.4	85.1	90.3	94.5	<i>92.4</i>	98.3	-	96.2
tomato soup can	80.9	88.7	89.4	87.9	93.3	88.9	94.0	<i>91.4</i>	96.1	-	92.4
mustard bottle	90.2	89.8	91.0	92.5	91.9	<i>94.7</i>	94.8	91.7	99.5	-	95.1
tuna fish can	87.9	89.5	89.9	84.7	95.2	92.2	94.0	90.4	95.1	-	96.1
pudding box	79.0	60.1	89.1	51.0	84.9	<i>85.1</i>	93.8	<i>89.0</i>	94.8	-	90.7
gelatin box	87.1	92.7	94.6	86.4	92.1	86.9	92.7	91.7	95.3	-	94.3
potted meat can	78.5	78.8	84.8	83.1	90.8	83.5	94.1	<i>91.2</i>	82.9	-	86.4
banana	85.9	80.7	88.7	79.1	70.0	<i>93.8</i>	94.3	89.5	96.0	-	72.3
pitcher base	76.8	91.7	91.8	84.8	91.1	92.3	94.3	91.7	98.8	-	94.6
bleach cleanser	71.9	73.6	72.0	76.0	86.8	83.0	92.6	<i>85.4</i>	94.4	-	90.3
bowl*	69.7	79.6	72.5	76.1	85.0	<i>91.6</i>	92.1	90.5	84.0	-	81.4
mug	78.0	86.8	92.1	<i>91.4</i>	91.9	89.8	94.1	<i>91.4</i>	96.9	-	91.3
power drill	72.8	56.5	82.9	76.0	87.2	<i>88.8</i>	94.3	88.8	91.9	-	92.3
wood block*	65.8	66.2	79.2	54.0	87.2	<i>90.7</i>	92.0	75.7	77.3	-	81.9
scissors	56.2	49.1	78.3	71.6	80.2	<i>83.0</i>	92.5	75.2	68.4	-	75.4
large marker	71.4	75.0	83.1	60.1	66.4	74.9	81.6	<i>81.2</i>	87.4	-	86.2
large clamp*	49.9	69.2	84.5	66.8	86.5	78.3	95.7	88.6	69.3	-	74.3
extra large clamp*	47.0	72.3	88.4	61.1	79.5	<i>54.7</i>	96.0	83.5	73.6	-	73.2
foam brick*	87.8	77.9	88.4	60.9	79.2	<i>89.9</i>	89.7	81.3	90.4	-	81.9
All	75.9	79.6	85.8	75.1	86.9	86.2	92.8	<i>87.1</i>	89.1	89.8	88.1

The improved performance of PoET compared to other RGB-based models in terms of the ADD-S score is due to a better estimate of the 3D translation as can be seen in Table 2. Again, ground-truth-based results are also presented. In contrast to directly estimating the 3D translation like PoET, PoseCNN as well as SilhoNet determine the 3D translation by estimating the depth and center pixel coordinates of an object and then reprojecting them using the known camera intrinsic parameters, which is considered the easier task to learn [3]. MCN treats the translation estimation as a classification problem by binning the translation space. In addition to the global image information provided by the multi-scale feature maps, our approach also learns to model the camera intrinsics, which results in a more accurate estimation of translation.

For 3D rotation estimation, we achieve the same average error as regular PoseCNN. Not surprisingly, networks that reduce the possible rotation space based on known object symmetries achieve a better result but with limited applicability to real-world scenarios. The mean average 3D rotation error is shown in Table 2. The influence of reducing the rotation by geometrical symmetries (†) is highlighted for PoseCNN. Since PoET makes use of the full multi-scale RGB feature maps, it outperforms SilhoNet for objects that have no symmetries in the silhouette space and only performs significantly worse for objects with rotational symmetries around an axis in 3D space, but without requiring a priori knowledge about the 3D shape of objects or restricting the rotation space based on symmetries. The main source for rotational errors are objects with rotational symmetries around one axis (master chef can, tomato soup can, tuna fish can). We have investigated the axes of our rotation errors and compared them to the symmetry axes for symmetric objects. The average tilt of rotation error axes with respect to symmetry axes across all test images and symmetric objects is only 15 degrees. If we ignore rotational errors about symmetry axes, our average rotation error reduces to 11.24 degrees, outperforming all RGB-only competitors, see Table 2. We refer the reader to the supplementary material for additional ablation experiments as well as PoET’s performance for the stricter BOP[33] and AUC of ADD [3] metrics and for the LM-O benchmark dataset [11].

4.2 Ablation Study

The ablation study investigates the influence of different components on the performance of our PoET framework. By assuming a perfect object detector that provides ground truth bounding boxes to PoET, we ensure that the evaluation includes every object present in the image even for those which might be not detected by the object detector. During the ablation study, we focus on the AUC of ADD/ADD-S metric, the average translation error and rotation error. All results reported in this section are for the same hyperparameter configuration as described in Section 3.2, the same fixed seed and trained for the same number of epochs. The final results are summarized in Table 3. We compare our best performing network from Section 4.1 (Baseline) to a class-agnostic version (Agnostic) and a network with less layers (Small). Finally, we investigate the influence of the integration of bounding box information into the transformer on the performance of PoET. We train PoET with learnable reference points (RP), trainable query embeddings (Q) or the combination of

Table 2: Comparison of average translation in cm and rotation error in degrees on YCB-V. Bold and italic values indicate the state-of-the-art for results, when working on ground-truth or predicted ROIs respectively. † indicates that the rotation predictions are reduced by geometric symmetries as described in [5].

Method	PoseCNN[3]	SilhoNet[5]	SilhoNet _{gt}	PoET _{gt}	PoET	PoseCNN	PoseCNN [†]	SilhoNet [†]	SilhoNet _{gt}	PoET _{gt}	PoET
master chef can	3.29	3.02	3.14	1.37	<i>2.26</i>	50.7	7.57	<i>1.21</i>	1.11	89.25	80.12
cracker box	4.02	5.24	2.38	1.48	<i>3.14</i>	19.69	19.69	19.86	9.53	9.68	21.87
sugar box	3.06	2.10	1.67	0.94	<i>1.42</i>	9.29	9.29	12.28	11.50	3.95	<i>4.40</i>
tomato soup can	3.02	2.40	2.24	1.09	<i>1.62</i>	18.23	8.40	<i>1.91</i>	1.82	50.97	49.29
mustard bottle	1.72	1.65	1.41	0.94	<i>1.42</i>	9.94	9.59	5.78	5.07	23.71	27.73
tuna fish can	2.41	<i>1.57</i>	1.49	0.95	1.79	32.80	12.74	<i>1.46</i>	1.50	60.30	63.72
pudding box	3.69	7.15	1.91	1.01	<i>1.94</i>	10.20	10.20	20.95	18.39	6.36	<i>6.87</i>
gelatin box	2.49	<i>1.09</i>	0.79	1.20	1.41	5.25	5.25	12.52	8.48	6.69	7.19
potted meat can	3.65	4.30	2.74	1.13	<i>1.75</i>	28.67	19.74	7.27	10.93	5.06	6.75
banana	2.43	4.12	2.59	1.06	<i>1.95</i>	<i>15.48</i>	<i>15.48</i>	16.29	5.70	7.90	20.40
pitcher base	4.43	<i>1.31</i>	1.29	0.95	1.55	11.98	11.98	<i>6.64</i>	6.61	7.51	8.04
bleach cleanser	4.86	3.60	3.99	1.09	2.47	20.85	20.85	51.28	48.42	16.32	21.93
bowl*	5.23	3.30	4.08	1.51	<i>1.76</i>	75.53	75.53	49.95	53.95	16.06	25.71
mug	4.00	2.61	1.43	1.28	<i>1.85</i>	19.44	19.44	18.14	7.02	3.86	5.59
power drill	4.59	6.77	3.19	0.98	2.29	9.91	9.91	30.54	10.66	5.92	6.45
wood block*	6.34	5.59	3.23	1.41	<i>4.75</i>	23.63	23.63	25.52	23.23	5.88	<i>14.32</i>
scissors	6.40	9.91	2.59	1.38	<i>3.72</i>	43.98	43.98	155.53	154.82	3.19	6.27
large marker	3.89	3.24	2.31	2.68	2.75	92.44	13.59	<i>10.44</i>	10.72	24.95	25.91
large clamp*	9.79	6.27	3.51	0.98	2.33	38.12	38.12	<i>3.54</i>	6.03	2.61	4.88
extra large clamp*	8.36	4.86	2.12	0.91	<i>3.10</i>	34.18	34.18	29.18	7.30	2.38	26.01
foam brick*	2.48	3.98	2.31	1.90	3.42	22.67	22.67	<i>13.84</i>	17.36	37.20	36.34
All	4.16	3.49	2.45	1.20	2.12	27.79	17.82	<i>16.04</i>	13.48	23.65	27.26

Table 3: Ablation study results of PoET on YCB-V. We report the AUC of ADD/S and the average translation and rotation error. Ablation of class mode (Agnostic), network size (Small) and transformer modifications (RP, Q, RP + Q). The exact meaning of the tags are described in the text.

Metric	Baseline	Agnostic	Small	RP	Q	RP + Q
AUC of ADD-S	92.8	88.9	91.8	87.6	82.2	42.0
AUC of ADD	80.8	73.2	78.1	66.8	59.5	12.2
Avg. T. Error [cm]	1.20	1.95	1.48	1.92	2.99	9.06
Avg. Rot. Error [°]	23.65	24.92	25.64	37.31	35.39	74.26

both (RP + Q). In all three cases the performance is reduced in comparison to a version of PoET that uses bounding box information to generate query embeddings and reference points. This shows that providing a transformer with bounding box information can greatly benefit its training process leading to improved performance for the same training duration. We kindly refer the reader to the supplementary material for an in-depth analysis and discussion of the ablation study.

4.3 Localization

PoET is well suited for vision-based object-relative localization. A transformation between coordinate frames A and B expressed in coordinate frame C is fully defined by the translation ${}^C t_{ab}$ and the rotation R_{ab} . Given a set of landmarks with known ground-truth pose (R_{wo}^i, t_{wo}^i), a single frame that captures at least one of those landmarks can be used in combination with PoET to localize the camera by estimating the relative pose ($\tilde{R}_{co}^i, \tilde{t}_{co}^i$) to all landmarks present in the frame. For each landmark, the estimated camera pose ($\tilde{R}_{wc}^i, \tilde{t}_{wc}^i$) can be determined by

$$\tilde{R}_{wc}^i = R_{wo}^i \tilde{R}_{co}^{iT} \quad \text{and} \quad {}^W \tilde{t}_{wc}^i = {}^W t_{wo}^i - R_{wo}^i \tilde{R}_{co}^{iT} {}^C \tilde{t}_{co}^i. \quad (4)$$

The final camera pose ($\tilde{R}_{wc}, \tilde{t}_{wc}$) can then be determined by taking the average over all landmarks present in the image. YCB-V’s test sequences offer 12 different camera trajectories and each with a different constellation of multiple objects serving as landmarks. For each individual frame we estimate the camera pose and compare it to the ground-truth. In Fig. 3 we show an example trajectory for a single sequence. For further examples we refer the reader to the supplementary material.

We compare three different approaches: using all detected objects (all), perform simple outlier rejection by taking and choosing the hypothesis the majority of objects agree on (out) or by incorporating the camera pose estimate from the previous frame into the outlier rejection in cases with multiple hypotheses having the same number of votes (prev). For the sake of comparison, we also calculate the camera pose by choosing the estimated camera pose being closest to the current ground-truth camera pose (best). Moreover, we differentiate between PoET being provided either object detections from the backbone (bb) or ground-truth bounding box information (gt). We calculate the sample error mean and standard deviation across all frames for the position and attitude and summarize them in Table 4. Simple outlier rejection greatly benefits the localization error. This is due to wrong relative object pose estimates having a large impact on the final estimated camera pose when taking a simple average. Incorporating the camera pose estimate from the previous frame only slightly improves the localization as it only helps in multi-hypothesis

cases. In comparison to the estimated trajectories using the best estimated camera pose, our outlier rejection based trajectories perform worse with around factor two. Nonetheless, the localization results when considering only individual frames are remarkable. The difference in performance between backbone and ground-truth detections is due to the backbone potentially missing an object or assigning a wrong class and thus throwing off the estimate.

To further motivate the use of PoET as a pose sensor in mobile robotics, we have integrated it in a state-of-the-art sensor fusion framework [34] and performed state estimation experiments with YCB-V objects in our motion capture room using inertial data for propagation and PoET for the pose update. For a representative trajectory, the average error in rotation was (roll, pitch, yaw) = (4.0, 6.9, 14.8) degrees and in translation (x, y, z) = (0.084, 0.179, 0.032) meters. The performance for our real data is only slightly worse compared to the performance on the benchmark dataset, even though we were using a completely different camera, than the one used to record the original dataset, and the objects have slightly changed in appearance in comparison to the original YCB objects. Nonetheless, PoET is able to provide sufficiently accurate pose data. Summarizing, PoET achieves sufficient localization accuracy such that it can be used as a pose sensor for robotics.

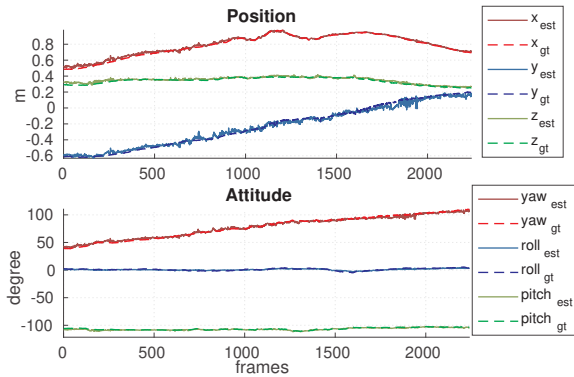


Figure 3: Example trajectory of a single sequence for gt deviation of camera pose across all 12 test and out. We plot the ground-truth and estimated position sequence frames. Position and attitude are reported in *m* and degree.

		(x, y, z) ± (σ _x , σ _y , σ _z) [mm]
bb	all	(119, 157, 52) ± (121, 131, 59)
	out	(56, 54, 28) ± (100, 106, 36)
	prev	(42, 43, 26) ± (56, 59, 27)
	best	(28, 27, 23) ± (31, 35, 22)
gt	all	(102, 138, 41) ± (107, 117, 53)
	out	(38, 34, 19) ± (80, 80, 28)
	prev	(32, 28, 19) ± (45, 39, 21)
	best	(21, 19, 17) ± (28, 32, 18)
		(yaw, roll, pitch) ± (σ _y , σ _r , σ _p) [deg]
bb	all	(19.6, 7.1, 12.4) ± (18.6, 11.8, 14.9)
	out	(5.4, 1.6, 2.2) ± (10.4, 2.8, 4.3)
	prev	(4.1, 1.4, 1.9) ± (5.0, 2.2, 2.1)
	best	(2.6, 1.5, 1.7) ± (3.0, 1.8, 1.7)
gt	all	(18.0, 6.7, 11.0) ± (16.5, 10.8, 12.9)
	out	(3.7, 1.2, 3.0) ± (8.9, 2.0, 3.0)
	prev	(2.8, 1.2, 1.4) ± (3.7, 2.2, 1.5)
	best	(2.0, 1.3, 1.2) ± (3.0, 1.4, 1.3)

Table 4: Sample error mean and standard deviation of camera pose across all 12 test and out. We plot the ground-truth and estimated position sequence frames. Position and attitude are reported in *m* and degree.

5 Limitations

In Section 4.1 it was discussed that PoET is outperformed by methods utilizing 3D object information for rotation estimation in particular for objects that have rotation-symmetric silhouettes. Wrongly estimated rotations lead to the assumption that the camera views the objects from a different angle resulting in wrong hypotheses for the localization task as discussed in Section 4.3. Moreover, the low resolution of the images in the YCB-V dataset means that RGB textures are not as dominant, especially when augmentation is used. This is the main reason why PoET has difficulties to estimate silhouette rotation-symmetric objects that are not symmetric in RGB space, see e.g. Table 2.

6 Conclusion

In this work we presented a novel, transformer-based framework for multi-object 6D pose estimation. It can be used on top of any object detector and the only input required is a single RGB image. The image is passed through an object detector backbone to create (multi-scale) image feature maps and to detect objects. Bounding box information of detected objects is fed into the transformer decoder which improves the learning. By taking the whole image into consideration during the estimation process, our framework does not rely on any additional information. We outperform other RGB-based methods by a wide margin on the YCB-V dataset. This is especially important for scenarios where no detailed 3D models or prior object information is available, or where computational efficiency is required and thus any input besides the RGB image has to be dropped. Moreover, we highlighted how PoET’s relative 6D object pose estimation can be used as a pose sensor for robot localization tasks.

Acknowledgments

This work was supported by the Federal Ministry for Climate Action, Environment, Energy, Mobility, Innovation and Technology (BMK) under the grant agreement 881082 (MUKISANO).

References

- [1] H. Song, W. Choi, and H. Kim. Robust vision-based relative-localization approach using an rgb-depth camera and lidar sensor fusion. *IEEE Transactions on Industrial Electronics*, 63(6): 3725–3736, 2016.
- [2] V. Loing, R. Marlet, and M. Aubry. Virtual training for a real application: Accurate object-robot relative localization without calibration. *International Journal of Computer Vision*, 126(9):1045–1060, 2018.
- [3] 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)*, 2018.
- [4] Z. Li, G. Wang, and X. Ji. Cdpn: Coordinates-based disentangled pose network for real-time rgb-based 6-dof object pose estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7678–7687, 2019.
- [5] G. Billings and M. Johnson-Roberson. Silhonet: An rgb method for 6d object pose estimation. *IEEE Robotics and Automation Letters*, 4(4):3727–3734, 2019.
- [6] Y. Li, G. Wang, X. Ji, Y. Xiang, and D. Fox. Deepim: Deep iterative matching for 6d pose estimation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 683–698, 2018.
- [7] Y. Labbé, J. Carpentier, M. Aubry, and J. Sivic. Cosypose: Consistent multi-view multi-object 6d pose estimation. In *European Conference on Computer Vision*, pages 574–591. Springer, 2020.
- [8] C. Li, J. Bai, and G. D. Hager. A unified framework for multi-view multi-class object pose estimation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 254–269, 2018.
- [9] P. J. Besl and N. D. McKay. Method for registration of 3-d shapes. In *Sensor Fusion IV: Control Paradigms and Data Structures*, volume 1611, pages 586–606. International Society for Optics and Photonics, 1992.
- [10] M. A. Fischler and R. C. Bolles. Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 24(6):381–395, 1981.
- [11] A. Krull, E. Brachmann, F. Michel, M. Y. Yang, S. Gumhold, and C. Rother. Learning analysis-by-synthesis for 6d pose estimation in rgb-d images. In *Proceedings of the IEEE international conference on computer vision*, pages 954–962, 2015.
- [12] S. Hinterstoisser, C. Cagniart, S. Ilic, P. Sturm, N. Navab, P. Fua, and V. Lepetit. Gradient response maps for real-time detection of textureless objects. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(5):876–888, 2011.
- [13] Z. Cao, Y. Sheikh, and N. K. Banerjee. Real-time scalable 6dof pose estimation for textureless objects. In *2016 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2441–2448. IEEE, 2016.
- [14] G. Pavlakos, X. Zhou, A. Chan, K. G. Derpanis, and K. Daniilidis. 6-dof object pose from semantic keypoints. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2011–2018. IEEE, 2017.
- [15] 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*, pages 548–562. Springer, 2012.

- [16] F. Michel, A. Kirillov, E. Brachmann, A. Krull, S. Gumhold, B. Savchynskyy, and C. Rother. Global hypothesis generation for 6d object pose estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 462–471, 2017.
- [17] K. Park, T. Patten, and M. Vincze. Pix2pose: Pixel-wise coordinate regression of objects for 6d pose estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7668–7677, 2019.
- [18] T. Hodan, D. Barath, and J. Matas. Epos: Estimating 6d pose of objects with symmetries. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11703–11712, 2020.
- [19] 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 *Proceedings of the IEEE International Conference on Computer Vision*, pages 3828–3836, 2017.
- [20] S. Peng, Y. Liu, Q. Huang, X. Zhou, and H. Bao. Pvnnet: Pixel-wise voting network for 6dof pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4561–4570, 2019.
- [21] 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 *Proceedings of the IEEE International Conference on Computer Vision*, pages 1521–1529, 2017.
- [22] G. Wang, F. Manhardt, F. Tombari, and X. Ji. Gdr-net: Geometry-guided direct regression network for monocular 6d object pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16611–16621, 2021.
- [23] A. Amini, A. S. Periyasamy, and S. Behnke. T6d-direct: Transformers for multi-object 6d pose direct regression. In *DAGM German Conference on Pattern Recognition*, pages 530–544. Springer, 2021.
- [24] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko. End-to-end object detection with transformers. In *European Conference on Computer Vision*, pages 213–229. Springer, 2020.
- [25] Y. Zhou, C. Barnes, J. Lu, J. Yang, and H. Li. On the continuity of rotation representations in neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5745–5753, 2019.
- [26] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, and J. Dai. Deformable detr: Deformable transformers for end-to-end object detection. In *International Conference on Learning Representations (ICLR)*, 2021.
- [27] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European Conference on Computer Vision*, pages 405–421. Springer, 2020.
- [28] S. Mahendran, H. Ali, and R. Vidal. 3d pose regression using convolutional neural networks. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 2174–2182, 2017.
- [29] C.-Y. Wang, A. Bochkovskiy, and H.-Y. M. Liao. Scaled-yolov4: Scaling cross stage partial network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13029–13038, 2021.
- [30] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, pages 740–755. Springer, 2014.
- [31] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.

- [32] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. *International Conference on Learning Representations (ICLR)*, 2019.
- [33] 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*, pages 577–594. Springer, 2020.
- [34] C. Brommer, R. Jung, J. Steinbrener, and S. Weiss. MaRS: A modular and robust sensor-fusion framework. *IEEE Robotics and Automation Letters*, 6(2):359–366, 2020.

PoET: Pose Estimation Transformer for Single-View, Multi-Object 6D Pose Estimation – Supplementary Material –

Thomas Jantos

Control of Networked Systems Group
University of Klagenfurt, Austria
thomas.jantos@aau.at

Mohamed Amin Hamdad

Infineon Technologies Austria AG
Villach, Austria
mohamedamin.hammad@infineon.com

Wolfgang Granig

Infineon Technologies Austria AG
Villach, Austria
wolfgang.granig@infineon.com

Stephan Weiss

Control of Networked Systems Group
University of Klagenfurt, Austria
stephan.weiss@aau.at

Jan Steinbrener

Control of Networked Systems Group
University of Klagenfurt, Austria
jan.steinbrener@aau.at

1 Introduction

In this supplementary material, we present additional results for our PoET framework that allow other works to perform detailed comparisons in the future. While the paper discusses the core results and main findings for our approach, we conduct and present the results of an extensive ablation study, which investigates the influence of the network architecture, the rotation representation, data augmentation, transformer modifications and the backbone on PoET’s performance. Moreover, we provide results for additional metrics and visualize qualitative results for a better understanding. Besides that, we provide additional details regarding our implementation and the dataset preparation. Finally, we evaluate PoET on the benchmark dataset Linemod-Occluded (LM-O) [1] and investigate different quaternion losses.

2 Implementation Details & Data Preparation

As discussed in the main paper, PoET can be trained on top of any object detector. For the YCB-V dataset [2] a Scaled-YOLOv4 [3] is used as the backbone object detector as it offers a good trade-off between speed and accuracy. An MS-COCO [4] pre-trained Scaled-YOLOv4 is fine tuned for 10 epochs on the YCB-V dataset for the object detection task. For the LM-O dataset we train PoET on top of a publicly available, pre-trained Mask R-CNN [5]¹ network. During the training of PoET, the weights of the object detector backbone are frozen.

In order to utilize the benefits of batch processing, the number of input object queries has to be constant across all images. This is usually not the case, as the number of objects present in an image can vary significantly. Therefore, the number of object queries n_q is fixed for a specific dataset to the maximum number of objects present in any of its images. If fewer objects are present in an image, the remainder of the object queries are filled up with dummies. Such a dummy object is assigned a bounding box $(c_x, c_y, w, h) = (-1, -1, -1, -1)$, a dummy class and a dummy query embedding.

¹<https://github.com/y1labbe/cosypose>

In case that the object detector predicts more objects than the number of allowed object queries, the top- n_q predictions are chosen based on the classification score. Due to the transformer’s attention mechanism, it should learn not to focus on dummy embedding feature vectors. For the loss calculation and the evaluation, dummy object queries are disregarded, as object queries are matched to ground-truth objects based on bounding box center distance, predicted class and generalized intersection over union (GIoU) using an adjusted Hungarian matcher similar to [6]. Dummy embeddings are not needed in inference, as only single images are processed. During training, the ground-truth bounding box and class are used instead of the object detector predictions.

We implement PoET using PyTorch [7] and train it on a single NVIDIA GeForce RTX 3090 for 50 epochs using AdamW [8] with a learning rate of $2e - 5$ and a batch size of 16. Our best performing network has five encoder and decoder layers, $d_h = 256$, 16 attention heads and a positional embedding dimension of $L = 32$. The network is simultaneously trained for 3D translation and 3D rotation estimation and the weighting parameters are set to $\lambda_t = 2$ and $\lambda_{rot} = 1$, such that both losses are in the same value range after scaling. During evaluation, our whole pipeline consumes around 5.3GB of VRAM and runs with 71 FPS.

YCB-V [2] contains 21 objects [9] and consists of 92 video sequences totalling 133,827 frames (real). Each video sequence contains multiple objects and the level of object occlusion and scene clutter varies between the sequences. Throughout the whole dataset at most 9 objects are contained in one scene and thus we fix n_q to 10. Following the proposed and commonly used data split [2, 10, 11, 12, 13, 14], we train our network on 80 of those 92 video sequences and reserve the remaining 12 sequences for testing. Out of these 12 sequences, the challenging keyframes are used for evaluation. The original authors of the dataset provide 80,000 synthetic images generated from 2D projections of the 3D object models (synt). We include these synthetic images during training for better network generalization. The background of the synthetic images is blank and thus, to improve network generalization even further, we randomly choose an image from the MS COCO dataset [4] as a background. Even though the BOP challenge also provides photorealistic images (pbr) for YCB-V, we do not utilize this data at all.

LM-O [1] is a single video sequence of the LM dataset [15] containing 1214 frames and the annotated ground-truth pose for 8 objects. In this specific scene, there is significant occlusion between the objects. Following the BOP challenge [10], we only use the 50,000 publicly available synthetic images, which were generated using physically-based rendering (pbr), for training PoET. Given that per image only a single instance of each object is present, the number of object queries n_q is set to 10. For evaluation we use the ADD(-S) metric as described by Hinterstoisser et al. [15].

We also perform random RGB augmentation during training by not only modifying the image color, sharpness, brightness and contrast, but also blurring the image and converting it to grayscale. However, images are not scaled, flipped or cropped.

3 Main Ablation Study on YCB-V

In this ablation study, we want to investigate the influence of different components on the performance of our PoET framework. By assuming a perfect object detector that provides ground-truth bounding boxes to PoET, we ensure that the evaluation includes every object present in the image even for those which might be not detected by the object detector. During the ablation study, we focus on the AUC of ADD/ADD-S metric, the average translation error and rotation error. All results reported in this section are for the same hyperparameter configuration as described in the main paper and the same fixed seed. The final results are summarized in Table 1. Detailed results for each class along with an in-depth analysis are provided in Section 4.

Network Architecture. We utilize our best performing model from the main paper as our baseline model (Baseline). First of all, we compare it to its class-agnostic counterpart (Agnostic). We observe that having dedicated outputs in the rotation and translation head for each class improves the results across all four metrics. Especially the average translation error is improved. We also

Table 1: Ablation study results of PoET on YCB-V. We report the AUC of ADD/-S and the average translation and rotation error. Ablation of class mode (`Agnostic`), network size (`Small`), rotation representation (`(Silho)Quat`), data augmentation (`w/o Aug.`) and transformer modifications (`RP`, `Q`, `RP + Q`). The exact meaning of the tags are described in the text.

Metric	Baseline	Agnostic	Small	Jitter	Quat	SilhoQuat	w/o Aug.	RP	Q	RP + Q
AUC of ADD-S	92.8	88.9	91.8	90.7	91.8	88.8	85.6	87.6	82.2	42.0
AUC of ADD	80.8	73.2	78.1	73.5	77.3	71.6	63.2	66.8	59.5	12.2
Avg. T. Error [cm]	1.20	1.95	1.48	1.58	1.25	2.11	2.36	1.92	2.99	9.06
Avg. Rot. Error [°]	23.65	24.92	25.64	30.35	28.63	27.54	37.95	37.31	35.39	74.26

investigate the influence of a smaller transformer network (`Small`) by training a PoET that only has 3 encoder and decoder layers; everything else is kept the same. While the smaller network results in an expected loss of performance across all metrics, it is still very competitive and leads to faster processing times (88 FPS). Therefore, by accepting minimal drops in performance a more efficient PoET can be trained in terms of memory consumption and computational time, which is beneficial for systems with hardware limitations.

Rotation representation. One of the design choices made was to use a 6D rotation representation. However, quaternions are still a popular rotation representation due to their ability to express a rotation with just 4 values and also they do not suffer from gimbal lock such as the Euler angle representation. Therefore, we want to show that PoET can be trained using quaternions (`Quat`). The output dimension of the rotation head is adjusted to $4 \cdot n_{cls}$, the output is L_2 normalized to ensure unit vector requirement for quaternions and the loss function is replaced by

$$L_{rot} = -\log(\langle q, \tilde{q} \rangle^2 + \epsilon), \quad (1)$$

where $\langle \cdot, \cdot \rangle$ represents the regular vector dot product and ϵ is a small number for numerical stability. The 6D representation achieves only slightly better results than the quaternion representation. The main reason is that `Quat` has on average a bigger rotation error. Important to note, the choice of loss function influences the performance of PoET in case of the quaternions. We train PoET also using the same quaternion loss function as Billings and Johnson-Roberson [14] (`SilhoQuat`). We can observe that while the average rotation error is slightly better, the average translation error is almost twice as big. The main reason for a decreased performance with respect to the translation estimation might be due to the different loss landscapes of the two different quaternion rotation loss functions. A more detailed analysis is provided in Section 9.

Data Augmentation. In previous work [16, 12] it was mentioned that data augmentation in the RGB space is beneficial for the pose estimation task. We investigated the influence of our data augmentations as described in Section 2 on PoET’s performance (`w/o Aug`). The results clearly indicate that including simple RGB data augmentation has a significant impact on PoET’s ability to learn the task of 6D pose estimation.

Transformer Modifications. Finally, we investigated the influence of replacing learnable reference points and query embeddings by directly passing bounding box information to the transformer decoder as described in our main paper’s methods section. We either make the reference points (`RP`), the query embeddings (`Q`) or both (`RP + Q`) learnable. In the latter case, the transformer architecture does not perform at all in comparison to the baseline network given that they were trained for the same number of epochs. In general, directly feeding bounding box information as reference points and object query embeddings does not only yield an improvement in terms of performance, but it also saves computational complexity as it means less trainable components.

Backbone. Besides the ablation study regarding PoET’s network architecture, we also want to investigate how the quality of the predictions of the object detector influences the performance. We measure quality of bounding boxes in terms of IoU with respect to their corresponding ground-truth. In order to model the decreasing object detector quality, we add jitter to the ground-truth bounding boxes before evaluating on them by sampling from a truncated normal distribution. Given a ground-truth bounding box, we sample a new center such that it is contained within the original bounding box. Furthermore, we sample new widths and heights from another truncated normal distribution

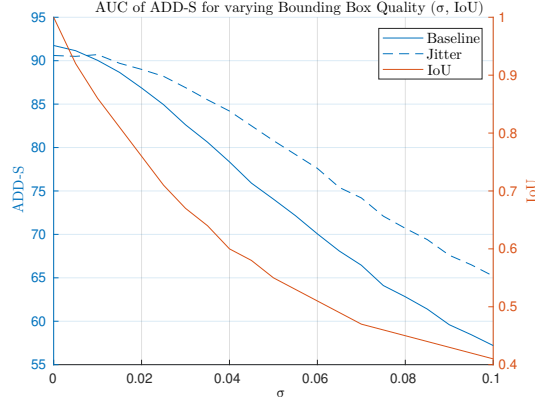


Figure 1: Comparison between Baseline and PoET trained with jitter bounding boxes for different bounding box qualities. With varying σ the IoU of the ground-truth boxes decreases and so the performance of PoET.

limiting them to be between 0.7 and 1.3 of the original size. The degree of noise is determined by the standard deviation σ . For both truncated normal distributions we use the same σ and calculate the AUC of ADD-S and the average IoU across all bounding boxes by varying σ between 0 and 0.1. While $\sigma = 0$ corresponds to no additional noise, a $\sigma = 0.1$ results in 99% of the sampled values being within a range of 150 pixels with respect to their original value. The results of this analysis can be found in Fig. 1. With decreasing object detector quality, the performance of PoET decreases. While the peak performance of Baseline cannot be matched by Jitter, we can see that the latter is not as influenced by the decreasing IoU of the object detector. Even though PoET’s performance depends on the quality of the object detector, we would still achieve state-of-the-art AUC of ADD-S when using non-perfect object detectors. For reference, FCOS [17], Faster R-CNN [18] and Scaled-YOLOv4 [3] evaluated on the keyframes achieve an average IoU across all detections of 0.86, 0.86 and 0.92 respectively. The main reasons for the discrepancy between the expected AUC of ADD-S score of PoET and the true one achieved by PoET working on Scaled-YOLOv4 is due to the fact that an object detector does not always detect all the objects present. However, the achieved score is close to the expected score. Moreover, we also compare our baseline to a PoET trained with jittered bounding boxes and $\sigma = 0.02$ (Jitter), but evaluated with ground-truth bounding boxes to make the comparison of the ablation networks independent of the quality of the object detector.

4 In-depth Analysis of Ablation Study

This section serves as an extension to our ablation study on the YCB-V dataset [2]. We present detailed results for each class for the AUC of ADD-S, AUC of ADD, average translation and average rotation error in Table 2, Table 3, Table 4 and Table 5 respectively.

The difference in performance in terms of the AUC of ADD-S metric for our different ablation networks mostly stems from their performance with respect to the translation estimation. When comparing the achieved ADD-S scores, reported in Table 2 of the ablation networks to their average translation and rotation errors, as reported in Table 4 and Table 5 respectively, one can see that the better the translation estimation is, the better the ADD-S score is, while the rotation seems to have a slightly smaller influence, e.g. for Agnostic and Small. On the other hand, a better estimation of the object rotation improves the performance on the AUC of ADD metric. Nevertheless, for the AUC of ADD metric a good translation estimation is required as can be seen by comparing Quat to SilhoQuat. We refer the reader to Section 9 for an in-depth comparison between Quat and SilhoQuat.

The class-agnostic version of PoET (Agnostic) mostly achieves two to three points less than the Baseline model for the AUC of ADD-S metric. However, the performance drops especially for

the wood block and the foam brick, two objects that show symmetries in the RGB as well as silhouette space. While `Agnostic` can better estimate the rotation of these two objects, the translation estimation suffers heavily due to not learning class-specific outputs.

As already mentioned in Section 3, `Small` achieves slightly lower scores across all metrics than our `Baseline`. To no surprise, this is due to the smaller network architecture offering less parameters to better capture objects’ characteristics. The performance drop is especially noticeable for the banana. There are many objects in the dataset that share similar shapes and thus can be easier learned. However, the banana has a unique shape and the smaller network seems to be not able to capture important characteristics of the banana.

Applying simple RGB data augmentation during the training improves the performance of PoET drastically. For the translation we can observe an improvement across all classes. While most classes also benefit from RGB augmentation in terms of rotation error (especially objects symmetric in RGB and silhouette space), we have observed the opposite effect for the master chef and the tomato soup can. These two objects are symmetric in silhouette space, but contain enough features in the RGB space to precisely determine the orientation of the object. Therefore, it is apparent that RGB augmentation benefits the network to focus more on geometric features of objects rather than RGB features.

In the main ablation study it was already discussed that replacing learnable reference points and object query embeddings by bounding box information yields better results for PoET. Given the results from Table 2, it is clear that using positionally-encoded bounding box information as the query embedding contains the most information. Namely, making this part of the transformer learnable (Q) performs significantly worse than our `Baseline`. A version of PoET trained with learnable reference points (RP) does also not achieve the same scores as the `Baseline` but it still performs better than Q. This shows that the positional bounding box encoding contains a lot of information, as RP can still predict sufficient reference points from the query embeddings and thus achieve higher scores. However, there are differences between RP and Q when comparing on the class level. While RP performs better for objects that have either (borderline) rotational symmetries in the silhouette space or multiple symmetry axes and planes in the RGB as well as silhouette space, having learnable query embeddings (Q) benefits PoET for objects that have either unique shapes (scissors, pitcher base and bananas) or that have features in RGB space that break the symmetries present in the silhouette space (gelatin and pudding box). Nonetheless, there are still outliers to this observation, e.g., the cracker box. Learning both reference points and query embeddings (RP + Q) significantly degrades performance across all classes. Further investigation is needed as to why this is the case.

These results illustrate that the design decisions for the baseline model indeed lead to a significant improvement in performance.

Table 2: Ablation study results of PoET on YCB-V. We report the AUC of ADD-S scores for each class and average score over all classes. Ablation of class mode, network size, rotation representation, data augmentation and transformer modifications.

Object	Baseline	Agnostic	Small	Jitter	Quat	SilhoQuat	w/o Aug	RP	Q	RP + Q
master chef can	92.9	87.6	92.6	91.4	92.4	88.6	82.4	91.0	71.5	58.1
cracker box	90.4	88.0	89.4	85.3	86.3	85.5	80.9	84.8	79.1	46.6
sugar box	94.5	91.2	92.4	93.6	93.3	89.4	92.7	88.0	89.5	49.5
tomato soup can	94.0	89.9	92.8	92.8	93.4	88.3	92.4	88.5	85.3	27.3
mustard bottle	94.8	87.4	90.0	89.1	91.6	82.7	88.3	88.2	81.7	59.1
tuna fish can	94.0	90.9	93.1	93.3	95.3	92.0	84.4	91.4	87.5	51.3
pudding box	93.8	93.6	93.0	92.0	94.1	89.2	88.0	84.4	89.9	37.2
gelatin box	92.7	90.6	92.3	93.4	91.0	92.9	84.6	86.9	91.7	40.0
potted meat can	94.1	90.8	92.9	91.8	93.9	89.7	91.6	90.2	83.9	59.3
banana	94.3	91.0	89.4	88.0	86.0	91.5	86.9	80.5	87.0	52.8
pitcher base	94.3	92.6	92.8	91.5	92.9	89.8	92.0	83.9	89.5	70.8
bleach cleanser	92.6	90.1	89.3	91.6	90.5	87.9	84.5	87.7	85.6	58.8
bowl*	92.1	87.6	94.2	91.1	93.2	92.4	93.5	89.9	80.5	17.8
mug	94.1	89.5	92.2	93.9	93.8	92.0	86.4	93.6	82.0	53.7
power drill	94.3	91.9	91.9	91.6	92.0	88.1	90.2	88.5	85.9	45.7
wood block*	92.0	81.4	92.9	86.4	90.5	85.2	75.6	91.2	63.4	31.5
scissors	92.5	90.8	91.8	88.3	92.6	90.5	83.8	68.7	80.6	25.9
large marker	81.6	84.5	84.7	83.1	85.3	82.0	60.8	85.6	77.3	18.5
large clamp*	95.7	93.0	95.6	94.7	96.0	92.6	88.7	93.6	80.5	20.3
extra large clamp*	96.0	88.5	94.8	94.8	94.5	88.6	91.4	91.3	83.6	45.3
foam brick*	89.7	75.2	88.9	87.2	90.3	85.3	79.3	91.0	69.3	11.6
All	92.8	88.9	91.8	90.7	91.8	88.8	85.6	87.6	82.2	42.0

Table 3: Ablation study results of PoET on YCB-V. We report the AUC of ADD scores for each class and average score over all classes. Ablation of class mode, network size, rotation representation, data augmentation and transformer modifications.

Object	Baseline	Agnostic	Small	Jitter	Quat	SilhoQuat	w/o Aug.	RP	Q	RP + Q
master chef can	40.6	28.3	34.9	37.9	37.6	45.5	40.3	30.6	21.5	16.8
cracker box	79.6	72.8	78.7	67.6	67.4	70.0	55.7	64.0	41.2	3.9
sugar box	89.6	83.4	85.0	87.1	86.7	79.3	85.3	74.0	78.4	12.0
tomato soup can	72.5	66.1	67.9	70.1	69.6	64.1	73.5	62.6	53.0	9.7
mustard bottle	82.2	69.4	76.4	59.1	71.8	50.6	53.7	64.0	63.9	12.2
tuna fish can	66.2	64.3	63.8	62.8	68.8	63.4	40.5	65.2	51.9	16.1
pudding box	88.4	88.2	86.8	84.9	88.5	79.1	77.4	69.3	80.5	13.8
gelatin box	87.0	83.4	85.4	87.9	82.7	86.8	68.4	73.9	84.4	11.3
potted meat can	88.3	81.1	85.4	83.5	87.3	79.9	81.0	79.4	66.6	29.1
banana	86.3	81.6	70.5	45.0	72.4	81.2	52.1	41.8	42.7	8.1
pitcher base	87.0	83.4	82.7	80.1	82.3	77.8	80.5	38.8	75.6	23.8
bleach cleanser	82.8	81.0	77.3	76.6	74.4	68.1	58.8	73.2	69.8	26.3
bowl*	76.9	62.4	82.3	66.8	68.7	59.4	60.9	71.9	45.0	1.3
mug	86.9	77.4	82.6	86.1	85.3	82.1	69.3	85.4	62.5	16.1
power drill	88.0	82.7	83.7	81.6	83.4	74.2	77.3	71.2	71.1	13.5
wood block*	83.8	64.2	86.1	73.1	78.6	70.3	50.6	81.4	35.7	2.9
scissors	85.9	82.3	83.7	75.8	85.5	81.3	63.0	33.5	66.2	15.4
large marker	73.3	76.1	76.3	73.0	76.9	72.8	48.4	76.5	66.8	7.7
large clamp*	90.0	83.7	89.4	87.3	90.8	82.5	74.6	85.3	64.7	2.1
extra large clamp*	90.4	71.7	86.7	87.7	86.4	69.5	72.2	80.5	64.3	14.4
foam brick*	72.0	53.7	75.4	70.0	78.6	66.6	44.9	79.3	44.1	0.3
All	80.8	73.2	78.1	73.5	77.3	71.6	63.3	66.8	59.5	12.2

Table 4: Ablation study results of PoET on YCB-V. We report the average translation error in cm for each class and average score over all classes. Ablation of class mode, network size, rotation representation, data augmentation and transformer modifications.

Object	Baseline	Agnostic	Small	Jitter	Quat	SilhoQuat	w/o Aug.	RP	Q	RP + Q
master chef can	1.37	2.43	1.38	1.68	1.49	2.29	3.47	1.75	5.47	7.60
cracker box	1.48	2.05	1.78	2.50	2.14	2.54	3.12	2.25	3.00	9.28
sugar box	0.94	1.57	1.42	1.16	1.04	2.03	1.33	2.22	1.99	9.38
tomato soup can	1.09	1.87	1.30	1.35	1.09	2.25	1.39	2.04	2.70	11.70
mustard bottle	0.94	2.28	1.77	1.93	1.02	3.46	1.96	1.94	3.23	6.64
tuna fish can	0.95	1.80	1.32	1.30	0.79	1.63	3.12	1.58	2.51	7.83
pudding box	1.01	0.99	1.19	1.45	0.84	2.00	2.02	2.49	1.62	8.79
gelatin box	1.20	1.59	1.42	1.14	1.44	1.28	2.67	2.14	1.40	8.90
potted meat can	1.13	1.82	1.40	1.59	1.06	1.98	1.54	1.83	3.07	7.20
banana	1.06	1.76	2.02	1.60	1.49	1.77	2.30	2.51	1.62	6.63
pitcher base	0.95	1.32	1.10	1.49	1.07	1.88	1.39	2.92	1.71	5.82
bleach cleanser	1.09	1.68	1.75	1.39	1.44	2.22	2.63	1.87	2.50	6.40
bowl*	1.51	2.37	1.05	1.38	1.17	1.46	1.14	1.71	2.85	14.54
mug	1.28	2.24	1.69	1.33	1.31	1.77	2.92	1.40	3.73	8.47
power drill	0.98	1.57	1.52	1.48	1.02	2.05	1.54	1.54	2.60	7.53
wood block*	1.41	3.54	1.29	2.52	1.53	2.81	4.57	1.38	6.69	13.10
scissors	1.38	1.73	1.53	2.16	1.14	1.93	2.16	3.67	3.23	12.31
large marker	2.68	2.30	2.15	2.57	2.12	2.63	4.64	2.11	2.94	11.16
large clamp*	0.98	1.59	1.01	1.22	0.75	1.71	2.42	1.37	3.34	12.40
extra large clamp*	0.91	2.12	1.24	1.19	1.13	2.20	1.83	1.73	3.15	8.34
foam brick*	1.90	4.64	2.31	2.58	1.79	2.80	3.42	1.76	5.29	14.12
All	1.20	1.95	1.48	1.58	1.25	2.11	2.36	1.92	2.99	9.06

Table 5: Ablation study results of PoET on YCB-V. We report the average rotation error in degrees for each class and average score over all classes. Ablation of class mode, network size, rotation representation, data augmentation and transformer modifications.

Object	Baseline	Agnostic	Small	Jitter	Quat	SilhoQuat	w/o Aug.	RP	Q	RP + Q
master chef can	89.25	107.22	106.61	92.34	92.68	72.89	58.75	103.36	70.19	82.52
cracker box	9.68	13.76	8.47	14.82	20.46	11.40	27.55	25.62	52.10	94.01
sugar box	3.95	5.12	4.29	5.15	7.57	3.97	5.77	11.32	8.80	32.04
tomato soup can	50.97	54.62	59.57	56.96	56.63	51.62	43.58	59.95	76.68	81.93
mustard bottle	23.71	31.69	21.20	67.50	45.24	76.25	96.55	51.75	20.42	125.40
tuna fish can	60.30	48.55	61.78	64.54	52.37	57.98	84.49	53.85	73.29	75.59
pudding box	6.36	7.02	7.43	4.92	10.96	8.08	11.35	27.65	15.61	51.83
gelatin box	6.69	7.92	4.84	6.09	13.85	3.80	18.01	26.92	10.27	90.58
potted meat can	5.06	6.80	6.20	6.27	9.30	6.67	14.00	11.89	14.77	41.17
banana	7.90	5.12	26.63	81.88	26.71	5.62	60.20	78.62	82.19	119.98
pitcher base	7.51	8.98	11.71	12.04	11.60	10.08	12.46	71.09	14.26	67.58
bleach cleanser	16.32	8.12	15.61	26.35	27.92	33.37	50.96	20.61	17.86	62.75
bowl*	16.06	29.26	12.23	34.62	26.86	48.01	45.97	21.18	62.19	82.84
mug	3.86	7.29	5.88	5.58	8.46	3.76	17.51	6.39	10.32	51.55
power drill	5.92	5.90	4.90	10.47	12.75	12.20	15.90	29.50	11.40	81.50
wood block*	5.88	4.85	3.94	6.90	15.77	6.45	13.76	9.06	20.51	63.41
scissors	3.19	4.20	7.43	8.53	12.73	5.72	68.54	146.30	10.33	113.19
large marker	24.95	19.85	27.00	27.82	23.43	24.64	49.56	26.12	28.68	90.32
large clamp*	2.61	4.00	3.24	3.33	7.26	2.67	8.42	6.36	11.10	85.03
extra large clamp*	2.38	16.15	3.71	2.58	6.89	26.88	28.00	8.60	13.67	70.12
foam brick*	37.20	11.39	12.89	27.38	19.32	25.14	102.18	18.79	28.20	82.56
All	23.65	24.92	25.64	30.35	28.63	27.54	37.95	37.31	35.39	74.26

5 Backbone Comparison on YCB-V

The results of PoET reported in the main paper for the YCB-V dataset [2] were achieved with a Scaled-YOLOv4 [3] as the object detector. Moreover, the ablation study discusses the influence of the object detector quality on the performance of PoET. However, only the quality influences the performance of PoET and not the actual object detector architecture. In Table 6 we compare PoET trained on top of a Scaled-YOLOv4 (YOLO) to PoET trained on top of a pre-trained Mask R-CNN (R-CNN). We evaluate YOLO and R-CNN with ground-truth bounding boxes as well as the object detectors’ actual predictions. In either case, the actual feature maps of the respective object detector is utilized. Moreover, we provide results of other state-of-the-art RGB-based methods for better comparison.

While the AUC of ADD-S score varies slightly for different classes, PoET achieves the same average score independent of the actual object detector architecture given ground-truth bounding boxes. Using the actual predictions confirms our ablation study’s findings that the quality of the object detector influences PoET’s performance. As reported in the main paper, our YOLO object detector performs better on YCB-V than a Mask R-CNN [5]. This is reflected in the small difference of the average AUC of ADD-S between YOLO and R-CNN. These results are a clear indication that PoET can be used in combination with different pretrained object detectors to achieve state-of-the-art results. Additionally, we provide the AUC of ADD scores in Table 7 and compare them to PoseCNN [2], the only other RGB-based method reporting results for this metric.

Taken together, these results illustrate that PoET can be used with different object detector backbones achieving competitive performance even for older object detector architectures.

Table 6: **AUC of ADD-S.** Comparison of PoET’s performance trained on top of different object detectors for the YCB-V dataset [2]. We report results for PoET trained on top of a Scaled-YOLOv4 [3] (YOLO) and Mask R-CNN [5] (R-CNN) evaluated with ground-truth (*gt*) and predicted bounding boxes. In either case, the actual feature maps of the respective object detector is utilized. Additionally, we compare the results to other RGB-only state-of-the-art methods. The 3D model row indicates how 3D model information is incorporated into the network. 2D indicates that only 2D image information is used.

Object	PoseCNN[2]	SilhoNet[14]	SilhoNet _{gt}	MCN[13]	MCN _{gt}	YOLO _{gt}	YOLO	R-CNN _{gt}	R-CNN
3D Model	Loss	Input + Sym	Input + Sym	2D	2D	2D	2D	2D	2D
master chef can	84.0	84.0	83.6	87.8	91.2	92.9	88.4	92.8	85.2
cracker box	76.9	73.5	88.4	64.3	78.5	90.4	80.5	92.0	87.0
sugar box	84.3	86.6	88.8	82.4	85.1	94.5	92.4	95.1	92.0
tomato soup can	80.9	88.7	89.4	87.9	93.3	94.0	91.4	95.6	91.7
mustard bottle	90.2	89.8	91.0	92.5	91.9	94.8	91.7	92.5	91.0
tuna fish can	87.9	89.5	89.9	84.7	95.2	94.0	90.4	95.3	92.3
pudding box	79.0	60.1	89.1	51.0	84.9	93.8	89.0	88.6	81.1
gelatin box	87.1	92.7	94.6	86.4	92.1	92.7	91.7	93.7	89.3
potted meat can	78.5	78.8	84.8	83.1	90.8	94.1	91.2	93.8	86.9
banana	85.9	80.7	88.7	79.1	70.0	94.3	89.5	84.7	80.3
pitcher base	76.8	91.7	91.8	84.8	91.1	94.3	91.7	94.9	93.1
bleach cleanser	71.9	73.6	72.0	76.0	86.8	92.6	85.4	93.8	88.0
bowl*	69.7	79.6	72.5	76.1	85.0	92.1	90.5	92.9	88.6
mug	78.0	86.8	92.1	91.4	91.9	94.1	91.4	94.2	88.5
power drill	72.8	56.5	82.9	76.0	87.2	94.3	88.8	94.9	91.6
wood block*	65.8	66.2	79.2	54.0	87.2	92.0	75.7	91.6	74.4
scissors	56.2	49.1	78.3	71.6	80.2	92.5	75.2	89.8	63.7
large marker	71.4	75.0	83.1	60.1	66.4	81.6	81.2	90.8	86.5
large clamp*	49.9	69.2	84.5	66.8	86.5	95.7	88.6	96.2	87.9
extra large clamp*	47.0	72.3	88.4	61.1	79.5	96.0	83.5	94.5	88.8
foam brick*	87.8	77.9	88.4	60.9	79.2	89.7	81.3	88.5	79.7
All	75.9	79.6	85.8	75.1	86.9	92.8	87.1	92.7	86.1

Table 7: **AUC of ADD**. Comparison of PoET’s performance trained on top of different object detectors for the YCB-V dataset [2]. We report results for PoET trained on top of a Scaled-YOLOv4 [3] (YOLO) and Mask R-CNN [5] (R-CNN) evaluated with ground-truth and predicted bounding boxes.

Object	PoseCNN [2]	YOLO _{gt}	YOLO	R-CNN _{gt}	R-CNN
master chef can	50.9	40.6	42.2	38.6	37.0
cracker box	51.7	79.6	57.9	83.2	70.8
sugar box	68.6	89.6	85.0	91.0	84.5
tomato soup can	66.0	72.5	69.5	72.9	68.4
mustard bottle	79.9	82.2	76.9	72.3	67.1
tuna fish can	70.4	66.2	59.7	70.0	65.4
pudding box	62.9	88.4	79.3	80.1	68.2
gelatin box	75.2	87.0	84.9	88.2	79.6
potted meat can	59.6	88.3	81.8	87.4	76.1
banana	72.3	86.3	74.1	57.8	53.7
pitcher base	52.2	87.0	82.1	88.7	84.7
bleach cleanser	50.5	82.8	68.6	88.4	76.0
bowl	6.5	76.9	68.7	69.1	59.7
mug	57.7	86.9	81.1	87.1	74.8
power drill	55.1	88.0	75.5	89.6	82.1
wood block	31.8	83.8	49.3	82.6	52.2
scissors	35.8	85.9	62.5	80.6	47.7
large marker	58.0	73.3	72.6	83.8	77.7
large clamp	25.0	90.0	75.9	91.4	75.2
extra large clamp	15.8	90.4	65.1	86.5	75.0
foam brick	40.4	72.0	60.0	75.1	61.9
All	53.7	80.8	70.1	79.3	68.5

6 Qualitative Results

6.1 6D Pose Estimation



Figure 2: Qualitative results for selected frames of the YCB-V dataset. Given the relative 6D object poses predicted by PoET we project the 3D model into the image.

In Fig. 2 we present PoET’s qualitative results for the 6D pose estimation task on the YCB-V dataset. The results range from almost perfect estimation in the top left image to failure cases in the bottom right image. The qualitative results show that PoET is able to handle occlusion due to its feature maps containing global information.

6.2 Localization Trajectories

In this subsection we provide further examples for the localization task discussed in the main paper. We introduced three different approaches to the problem of localizing a camera given landmark world positions and using PoET’s relative object pose estimates. Furthermore, we investigated the best possible estimate by choosing the hypothesis that is closest to the current ground-truth camera pose. In Fig. 3 we compare the estimated position and attitude for the two outlier rejection approaches to the best trajectory approach across one trajectory. In comparison to the best possible estimate, simple outlier rejection performs only slightly worse for this specific trajectory. We visualize the 3D trajectories in Fig. 4.

Additionally, we provide PoET’s performance for the localization task for a different trajectory in Fig. 5. There we compare between PoET being provided either ground-truth or backbone bounding box information. In both cases, it is possible to estimate good trajectories based on PoET’s relative object poses. However, PoET with backbone information performs slightly worse and is more prone to wrong estimates resulting in a noisier estimated trajectory.

Summarizing, these results support our claim that PoET can be used as a pose sensor for the localization task in robotics applications.

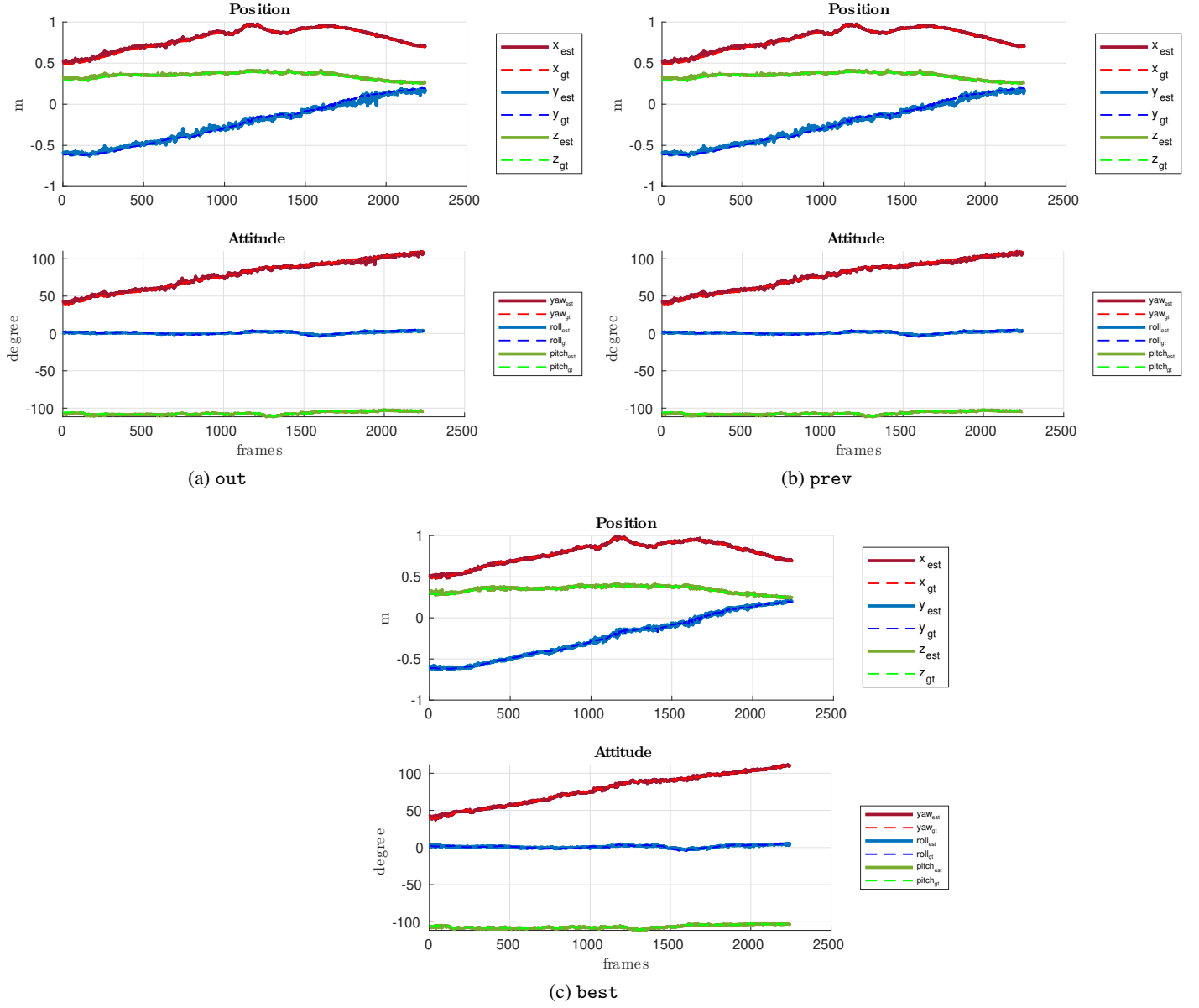
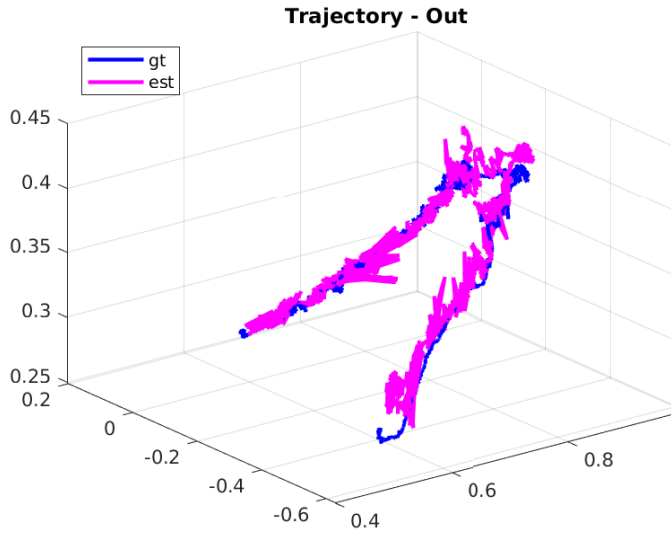
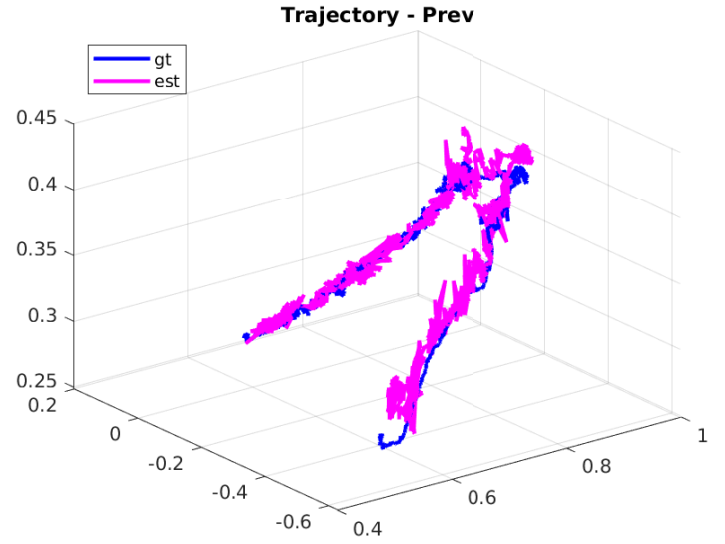


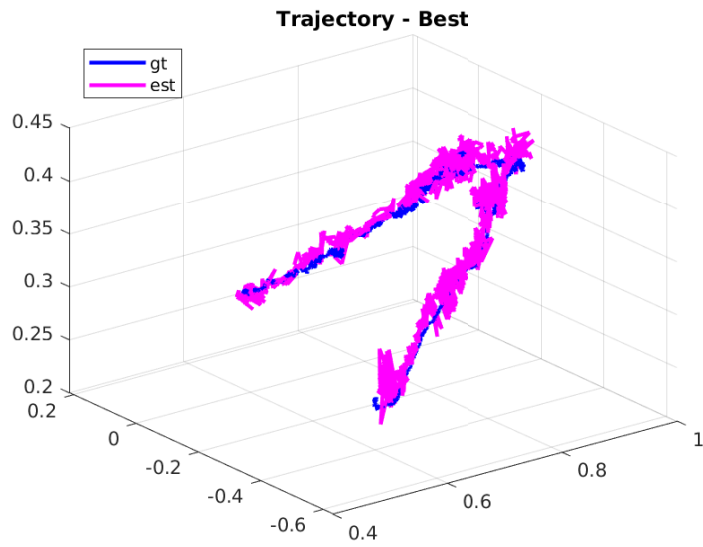
Figure 3: Estimated trajectories (est) for the same trajectory as in the main paper.



(a) out

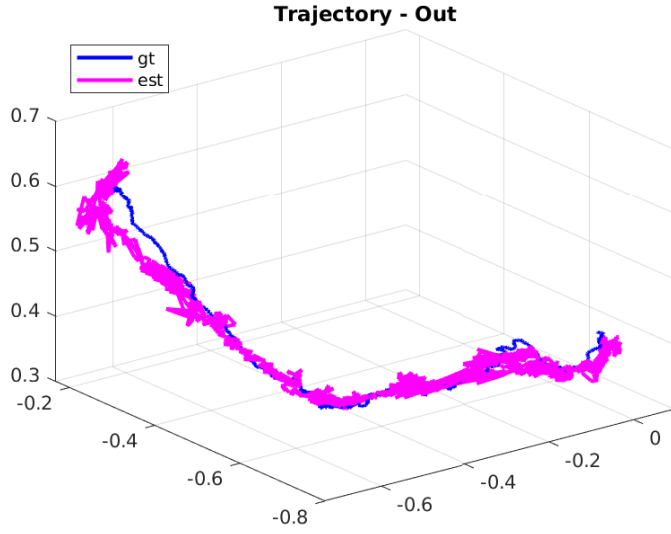


(b) prev

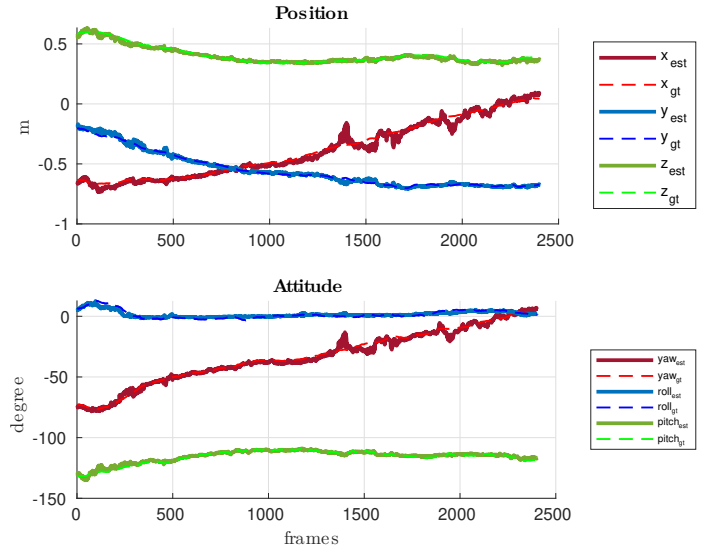


(c) best

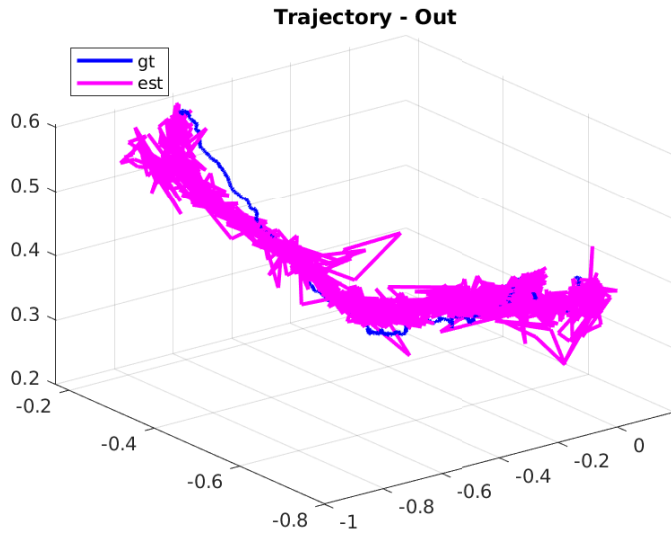
Figure 4: Visualization of the estimated trajectory for the same YCB-V scene shown in the main paper. We compare the estimated trajectory (est) to the ground-truth trajectory (gt).



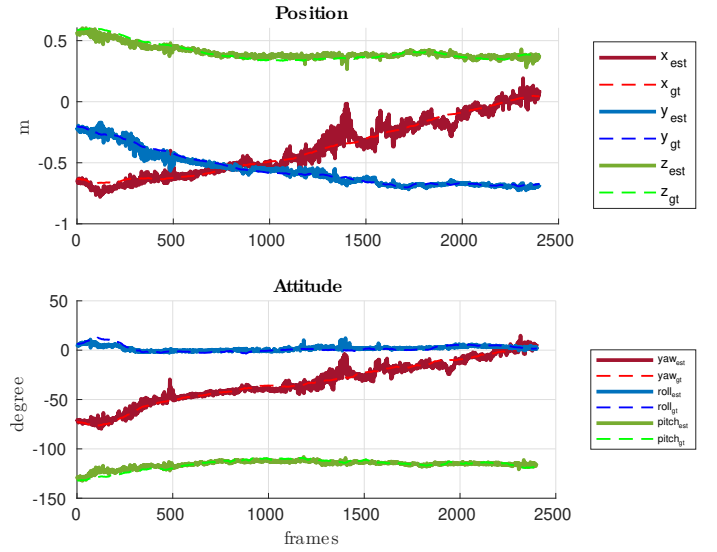
(a) gt + out



(b) gt + out



(c) bb + out



(d) bb + out

Figure 5: Comparison of estimated trajectories between PoET with ground-truth (gt) or predicted (bb) bounding box information. This is a different trajectory than the one used for discussion in the main paper.

7 Results on LM-O

In Table 8 we report the results of PoET for the ADD(-S) metric as described by Hinterstoisser et al. [15]. In contrast to other state-of-the-art approaches, we only train on pbr data as dictated by the BOP challenge [10]. Moreover, we only train a single PoET model, instead of training an individual network for each object class. Similarly as for YCB-V, methods that utilize either PnP or iterative refinement and thus 3D object models achieve state-of-the-art results. PoET outperforms other purely RGB-based approaches by more than 10 points.

Furthermore, we conduct a comparison between PoET trained on top of a pre-trained Mask R-CNN [5] for the LM-O dataset [1, 15] by evaluating it with ground-truth and predicted bounding boxes. Once again, evaluating with ground-truth bounding boxes results in better scores. Nevertheless, PoET still achieves state-of-the-art results for RGB-only methods when evaluated on actual predictions.

Table 8: Comparison with state-of-the-art on LM-O. We report the ADD(-S). PE indicates the number of pose estimators trained. N means that one PE was trained for each class. (*) denotes symmetric objects. real, synt and pbr, respectively, refer to real data, synthetically generated data by projecting the 3D models onto a black image background and photorealistic simulated images.

Method				PnP			IR
	PoseCNN [2]	PoET	PoET _{gt}	Pix2Pose [19]	PVNet [20]	GDR-Net [16]	DeepIM [11]
PE	1	1	1	N	N	N	1
Data	real + syn	pbr	pbr	real	real + syn	real + pbr	real + syn
Ape	9.6	10.2	12.7	22.0	15.8	46.8	59.2
Can	45.2	31.8	51.0	44.7	63.3	90.8	63.5
Cat	0.9	9.0	10.9	22.7	16.7	40.5	26.2
Driller	41.4	33.9	53.3	44.7	65.7	82.6	55.6
Duck	19.6	15.4	22.6	15.0	25.2	46.9	52.4
Eggbox*	22.0	44.7	50.4	25.2	50.2	54.2	63.0
Glue*	38.5	58.7	63.9	32.4	49.6	75.8	71.7
Holep.	22.1	24.7	29.8	49.5	36.1	60.1	52.5
MEAN	24.9	28.5	36.8	32.0	40.8	62.2	55.5

8 BOP Results

In the main paper we presented the results on YCB-V [2] following the most commonly used evaluation metrics [2, 11, 12, 13, 14, 16]. In recent years the evaluation protocol of the BOP challenge [10, 21] has become more popular and thus, we present the results of PoET for YCB-V on the BOP metrics. We refer the reader to [10] for a detailed explanation of the evaluation metrics. We report the average recall (AR) score by calculating the mean of the three main metrics: $AR = (AR_{MSPD} + AR_{MSSD} + AR_{VSD}) / 3$. The final results are summarized in Table 9. To the best of our knowledge, we are the only approach to report their results on the BOP challenge that relies solely on 2D image information. In contrast to the AUC of ADD-S metric, the metrics employed by the BOP challenge are stricter with respect to the objects final estimated rotation. Therefore, it is no surprise that state-of-the-art results are achieved by methods that either use 3D object model keypoints to predict the final pose (PnP) or that perform iterative refinement utilizing the 3D model given an initial estimate (IR). Nevertheless, PoET achieves competitive results and even outperforms the two PnP-based methods Pix2Pose [19] and CDPNv2 [22].

Table 9: Comparison of state-of-the-art methods on YCB-V for BOP metrics [10]. The results of other approaches are either obtained from the corresponding work or from the official leaderboard: <https://bop.felk.cvut.cz/leaderboards/>. PnP and IR stand for methods that either use PnP or iterative refinement respectively. 2D indicates that only RGB-image information was utilized.

Method	3D Model	AR_{MSPD}	AR_{MSSD}	AR_{VSD}	AR
PoET	2D	54.2	57.2	49.4	53.6
PoET _{gt}	2D	66.4	71.9	66.5	68.3
Pix2Pose [19]	PnP	57.1	42.9	37.2	45.7
EPOS [23]	PnP	78.3	67.7	62.6	69.6
CDPNv2 [22]	PnP	63.1	57.0	39.6	53.2
GDR-Net [16]	PnP	84.2	75.6	66.8	75.5
CosyPose [12]	IR	85.0	84.2	77.2	82.1

9 Quaternion Loss

In this section we compare different quaternion loss functions and their loss landscapes to each other. In the main paper we achieve state-of-the-art results on the YCB-V dataset [2] using a 6D rotation representation. PoET achieves similar results by using a quaternion representation and the loss function

$$L_{rot} = -\log(\langle q, \tilde{q} \rangle^2 + \epsilon), \quad (2)$$

where $\langle \cdot, \cdot \rangle$ represents the regular vector dot product and ϵ is a small number for numerical stability. However, as shown in the ablation study, using the same quaternion loss function as SilhoNet [14] results in a worse performance. This loss function is given by

$$L_{rot} = \log(\epsilon + 1 - |\langle q, \tilde{q} \rangle|), \quad (3)$$

where $|\cdot|$ denotes the absolute value. Comparing the loss function landscape from Eq. (2) with the one from Eq. (3) in Fig. 6, it is observable that for small errors our loss converges to 0, while the loss from [14] goes to $-\infty$. The former will result in a more stable training towards the end due to smaller gradients. On the other hand, for large errors our loss function is close to ∞ and thus yielding stronger gradients, while the other loss function is close to 0. Having those smaller gradients towards the end of the training benefits the multi-task loss as the translation loss is not overshadowed by the rotation loss. Thus, the network can focus on further improving its performance with respect to the translation estimation as can be seen in Table 4. On the contrary, Table 5 shows that the large gradients of Eq. (3) for small errors influence PoET to achieve better average rotational error. In the end, being able to better estimate the object translation benefits the pose estimation task in terms of the AUC of ADD-S metric.

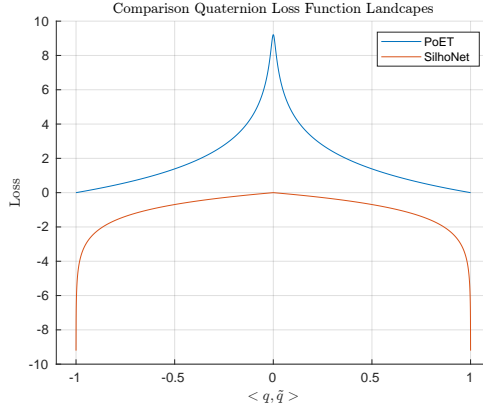


Figure 6: Comparison of the loss landscapes between our proposed quaternion loss function in Eq. (2) and the one used by [14] as shown in Eq. (3). For small errors, that is $\langle q, \tilde{q} \rangle$ being close to -1 or 1, our loss function results in smaller gradients and thus allows our network to better learn. For both loss functions ϵ is set to $1e - 4$.

References

- [1] A. Krull, E. Brachmann, F. Michel, M. Y. Yang, S. Gumhold, and C. Rother. Learning analysis-by-synthesis for 6d pose estimation in rgb-d images. In *Proceedings of the IEEE international conference on computer vision*, pages 954–962, 2015.
- [2] 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)*, 2018.
- [3] C.-Y. Wang, A. Bochkovskiy, and H.-Y. M. Liao. Scaled-yolov4: Scaling cross stage partial network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13029–13038, 2021.
- [4] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, pages 740–755. Springer, 2014.
- [5] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2961–2969, 2017.
- [6] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, and J. Dai. Deformable detr: Deformable transformers for end-to-end object detection. In *International Conference on Learning Representations (ICLR)*, 2021.
- [7] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- [8] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. *International Conference on Learning Representations (ICLR)*, 2019.
- [9] B. Calli, A. Singh, A. Walsman, S. Srinivasa, P. Abbeel, and A. M. Dollar. The ycb object and model set: Towards common benchmarks for manipulation research. In *2015 International Conference on Advanced Robotics (ICAR)*, pages 510–517. IEEE, 2015.
- [10] 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*, pages 577–594. Springer, 2020.
- [11] Y. Li, G. Wang, X. Ji, Y. Xiang, and D. Fox. Deepim: Deep iterative matching for 6d pose estimation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 683–698, 2018.
- [12] Y. Labbé, J. Carpentier, M. Aubry, and J. Sivic. Cosypose: Consistent multi-view multi-object 6d pose estimation. In *European Conference on Computer Vision*, pages 574–591. Springer, 2020.
- [13] C. Li, J. Bai, and G. D. Hager. A unified framework for multi-view multi-class object pose estimation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 254–269, 2018.
- [14] G. Billings and M. Johnson-Roberson. Silhonet: An rgb method for 6d object pose estimation. *IEEE Robotics and Automation Letters*, 4(4):3727–3734, 2019.
- [15] 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*, pages 548–562. Springer, 2012.
- [16] G. Wang, F. Manhardt, F. Tombari, and X. Ji. Gdr-net: Geometry-guided direct regression network for monocular 6d object pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16611–16621, 2021.

- [17] Z. Tian, C. Shen, H. Chen, and T. He. Fcos: Fully convolutional one-stage object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9627–9636, 2019.
- [18] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in Neural Information Processing Systems*, 28, 2015.
- [19] K. Park, T. Patten, and M. Vincze. Pix2pose: Pixel-wise coordinate regression of objects for 6d pose estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7668–7677, 2019.
- [20] S. Peng, Y. Liu, Q. Huang, X. Zhou, and H. Bao. Pvnet: Pixel-wise voting network for 6dof pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4561–4570, 2019.
- [21] T. Hodan, F. Michel, E. Brachmann, W. Kehl, A. GlentBuch, D. Kraft, B. Drost, J. Vidal, S. Ihrke, X. Zabulis, et al. Bop: Benchmark for 6d object pose estimation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 19–34, 2018.
- [22] Z. Li, G. Wang, and X. Ji. Cdpn: Coordinates-based disentangled pose network for real-time rgb-based 6-dof object pose estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7678–7687, 2019.
- [23] T. Hodan, D. Barath, and J. Matas. Epos: Estimating 6d pose of objects with symmetries. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11703–11712, 2020.