

OBJECT DETECTION AND AUTOENCODER-BASED 6D POSE ESTIMATION FOR HIGHLY CLUTTERED BIN PICKING

Timon Höfer, Faranak Shamsafar, Nuri Benbarka, Andreas Zell

Cognitive Systems Group, CS Department
Eberhard Karls Universität Tübingen, Germany

ABSTRACT

Bin picking is a core problem in industrial environments and robotics, with its main module as 6D pose estimation. However, industrial depth sensors have a lack of accuracy when it comes to small objects. Therefore, we propose a framework for pose estimation in highly cluttered scenes with small objects, which mainly relies on RGB data and makes use of depth information only for pose refinement. In this work, we compare synthetic data generation approaches for object detection and pose estimation and introduce a pose filtering algorithm that determines the most accurate estimated poses. We will make our real dataset for object detection available with the paper.

Index Terms— 6D pose estimation, object detection, synthetic dataset, bin picking.

1. INTRODUCTION

Bin picking is a major automation task with various applications in industrial sectors. The core starting problem of this work is the 6D pose estimation of instances. To tackle this problem, an RGB-D or depth camera is usually installed on top of the bin. There are existing solutions to bin picking of large objects, mostly using local invariant features [1, 2] or template-matching algorithms [3], which rely on the computationally expensive evaluation of many pose hypotheses. Moreover, local features do not perform well for texture-less objects, and thus, template-matching often fails in heavily cluttered scenes with severely occluded objects. Additionally, depth sensors are often more sensitive to lighting variations than RGB cameras [4]. Most importantly, for small objects, the depth information is often insufficient to get accurate pose estimates. Therefore, in this work, we focus on RGB-based convolutional neural networks, which make use of depth information only for pose refinement.

Accordingly, one of the important issues for training a deep network is labeling the training dataset, which requires high effort for tasks like 6D pose estimation [5]. Given a

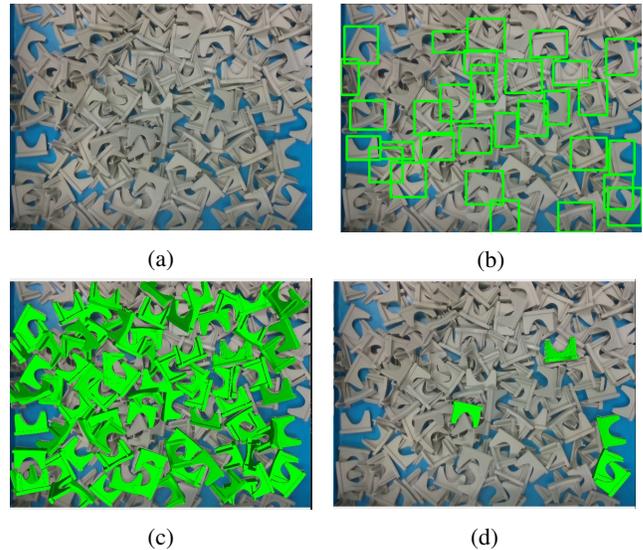


Fig. 1: (a) A sample image of a cluttered bin captured by a Microsoft Azure Kinect RGB-D camera. (b) Detection results, limited to 30 objects, to be visually recognizable. (c) Pose estimation results for all the detected objects. (d) The best five selected poses based on the filtering algorithm.

CAD model of the object, which is usually available in the industry, generating a synthetic dataset is possible. However, training on only synthetic 2D images of the CAD models does not generalize well to real data. Hence, more insightful techniques are required to bridge the gap between simulation and reality [4].

Generally, a state-of-the-art object detector is first used to recognize individual objects, and the resultant cropped images are passed to the pose estimator. Following [6, 7], we use Mask R-CNN [8] for object detection. As for the task of pose estimation, we consider an augmented autoencoder [4], since it has demonstrated good performance in bin picking of deformable products [7]. Sample results of our proposed method are displayed in Fig. 1. The main contributions of this work are as follows:

1. We present a comprehensive framework from creating a synthetic dataset to the prediction of the 6D pose es-

Financial support for this work is provided by the Federal Ministry for Economic Affairs and Energy of Germany (BMWi) in the project iBinPick (ZF4076504DB8).

timates in bin picking scenarios, where no real labeling is needed.

2. We show that a more realistic renderer for data generation significantly improves the performance on heavily cluttered piles.
3. We present a pose filtering scheme to select the best pose predictions.
4. We give an analysis of how the performance of the auto-encoder can be improved in bin picking scenarios.

The remainder of the paper is as follows: after presenting the related work in section 2, we explain our methodology in section 3. Experimental results are discussed in section 4 and we conclude the paper in section 5.

2. RELATED WORK

The object pose estimation problem in bin picking scenarios has been investigated using local invariant features (e.g., point pair features [1, 2]) and template-matching [3]. However, these approaches do not show acceptable performance in bin picking from a cluttered pile of textureless small objects. Recently, convolutional neural networks [4, 6, 9, 10] have proven to be promising in the BOP2020 challenge [11] on 6D pose estimation, even surpassing the depth-based methods. They also tend to be faster, mostly with a runtime of less than one second. These methods mainly depend on an object detection phase, which is achieved using a state-of-the-art object detector (Mask R-CNN [8], RetinaNet [12], Faster R-CNN [13]).

The most important requirements for these deep networks are labeled datasets. As the effort of labeling 6D poses in cluttered scenes is high and demands a complex setup [5], some works [14, 4] have proposed training on synthetic images rendered from a 3D model. To bridge the gap to reality, random augmentations and domain randomization techniques have been applied [15, 16]. As a different solution to this problem, a more realistic data generation using a physics engine has been proposed in [17]. In our work, we benefit from this approach to generate photorealistic cluttered piles and propose the full framework for object detection and pose estimation.

Moreover, once the general pipeline is given, the predicted poses can be further refined using a pose refinement method. Previously, this step has been achieved [6, 4] by the ICP algorithm [18]. As the ICP-based methods show slow performance, we show that incorporating the depth information into the pose estimation procedure, achieves comparable results. Besides, a filtering algorithm is applied to choose the best poses among the estimated ones.

3. METHODOLOGY

In this work, we consider heavily cluttered and occluded scenes of small industrial objects. Here, we explain the meth-

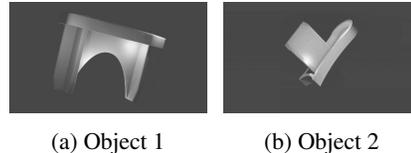


Fig. 2: The objects of interest are grey plastic pieces with sizes of $2.3 \times 3.6 \times 0.8\text{cm}^3$ and $1.5 \times 2.7 \times 0.9\text{cm}^3$, respectively.

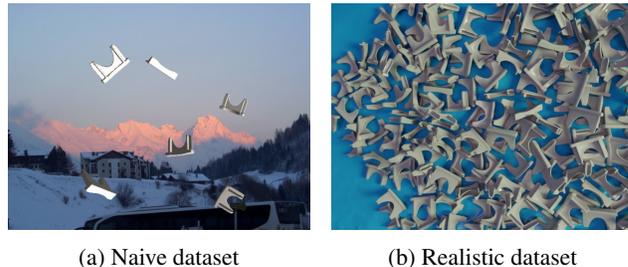


Fig. 3: The synthetic datasets generated with two different pipelines for object 1.

ods for dataset generation, followed by the full framework for object detection and pose estimation.

3.1. Dataset generation

Creating a synthetic dataset can be achieved by using the CAD models of the objects. Since our pipeline has two main tasks, we need to create a dataset for both, object detection and pose estimation.

3.1.1. Dataset for object detection

As the first approach, we make use of the pipeline in [4] to generate synthetic images to train the object detector. In particular, a CAD model is used to render the object on a black background. Then, random images from the Pascal VOC dataset [19] are added as a background, followed by random augmentations strategies. We create 60K images per object with 5-20 instances per scene. Due to their simplicity, we call these images “naive dataset”.

For the second approach, we employ BlenderProc [17] to generate more realistic synthetic images. BlenderProc utilizes a physics engine to make the synthetic data look more realistic. Furthermore, it uses different lighting effects, object materials and applies physics and collision checking as well. With BlenderProc, we generate for each object type 20K images with 30 instances and 5K images with 300 objects. The camera configuration is sampled in a range of 20° around the top of the scene with a height between 27-33cm. We call this the “realistic dataset”.

In this work, we consider two industrial objects (see Fig. 2). Sample images of the naive and realistic images for object 1 are depicted in Fig. 3.



(a) Real test dataset 1 (b) Real test dataset 2

Fig. 4: Examples of our labeled test dataset in real scenarios.

3.1.2. Dataset for pose estimation

Similar to the naive dataset in the previous section, we generate images with corresponding 6D pose annotations, with an additional step of image cropping around the objects. We also compare the original pipeline results against a new dataset, where we render multiple objects in the image crops. In Fig. 5, samples of different training data are displayed on the left side. While the top image shows one single object in each image crop, the bottom one includes multiple objects.

3.1.3. Test dataset for object detection

To evaluate the object detector, we captured 50 real images per object model with more than 100 instances in the bin. The images were taken using a Microsoft Azure Kinect camera mounted at a height of 30cm over the bin (see Fig. 4).

3.1.4. Test dataset for pose estimation

While we show qualitative results in real-world scenarios, the quantitative results are reported on a synthetic dataset, because only for this we have full ground truth data. The autoencoder is trained on synthetic data following the pipeline in [4], and we create the test dataset with BlenderProc [17]. To be more precise, BlenderProc generates 3D scenes, whereas the pipeline in [4] creates augmented 2D images. Since the data distributions of these synthetic datasets are different, we will show that training the pose estimator on the naive dataset and testing on the photorealistic images leads to suitable performance.

As such, for each object, we created one dataset consisting of 1K images with 300 objects. The camera is located at 30cm on top of the bin ground. We choose a blue background to make it visually comparable to our real settings.

3.2. Object detection

In general, any state-of-the-art object detector can be used (Faster R-CNN [13], RetinaNet [12], SSD [20]) for object detection. However, these methods only predict the bounding boxes. Therefore, we choose Mask R-CNN, which has the advantage of predicting the segmentation masks of the object as well. This can be used for pose refinement [21] by segmenting the point clouds of the objects. In addition, we

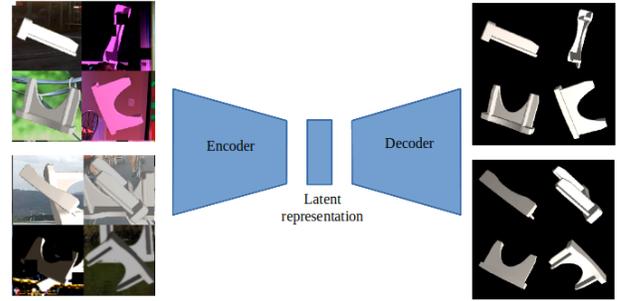


Fig. 5: The architecture of the autoencoder. The autoencoder is trained to map the augmented images to the original image. On the left, there are the two different types of training data.

can compare the pose estimation results when, instead of the whole bounding box, only the pixels visible in the segmentation mask are given to the pose estimation module. Given an image, we predict a set \mathcal{D} of object detections.

3.3. Pose estimation

To receive pose estimates from the set \mathcal{D} of the detections, we resort to the autoencoder network presented in [4]. As the autoencoder’s training procedure is based on only synthetic data, it is more applicable to new industrial settings, where no labeled data exists. In addition, the autoencoder has demonstrated good performance in pick-and-place tasks [22, 7]. Its method of operation is as follows: the autoencoder is a dimensionality reduction technique trained to extract a 3D object from image crops (see Fig. 5). After training, we create a codebook to determine the rotation of the object. The codebook is the set of all latent representations of the discretized 3D rotations that cover the whole $SO(3)$.

At test time, the image crops from the set \mathcal{D} are fed into the encoder. The resulting latent representation z_{test} is then compared with all the latent representations (z_i) from the codebook via a k-NN search, with the similarity function as:

$$cos_i = \frac{z_i z_{test}}{\|z_i\| \|z_{test}\|}. \quad (1)$$

We then choose the rotations with the highest cosine similarities.

3.4. Selecting the best pose estimates

In cluttered scenes, we have pose estimations of several hundred objects. These will be used for the picking task, and as the robot will only pick one object at a time, we are only interested in the k-top pose estimates, and the question arises, how to select the best k-pose estimates. While one can sort the pose estimates regarding the scores given by Mask R-CNN, or by the highest cosine similarities (1), we define a new selecting method that compares the depth of the original image with

Object	dataset	AP _{50:95} (%)	AP ₅₀ (%)	AR ^{max=100} (%)
Object 1	naive	9.3	12.5	9.6
	realistic	66.8	82.8	80.3
Object 2	naive	0.9	1.0	0.1
	realistic	50.4	68.7	59.2

Table 1: Object detection results after training Mask R-CNN on different synthetic datasets

the depth of the rendered image. Let A_1^i be the predicted segmentation mask of object i . Given a predicted 6D pose (\bar{t}, \bar{R}) , we render the depth image $\bar{\Omega}$ and we define the set of pixel $A_2^i := \{(p, q) : |\Omega(p, q) - \bar{\Omega}(p, q)| < m\}$, for some margin m and where Ω represents the real depth image. A_3^i is the segmentation mask of the rendered image. With these definitions, we can build the intersection: $\mathcal{A}^i := A_1^i \cap A_2^i \cap A_3^i$. For each pose estimate we calculate the depth error as follows:

$$e_i = \sum_{(p,q) \in \mathcal{A}^i} |\Omega(p, q) - \bar{\Omega}(p, q)| \quad (2)$$

In the next section, we compare the different approaches.

4. EXPERIMENTAL RESULTS

4.1. Experiments on object detection

We fine-tuned a pretrained Mask R-CNN with a ResNet-50 backbone for 15 epochs with an initial learning rate of 0.001 and a mini-batch size of 4 images. The learning rate was reduced by a factor of 10 at epochs 3, 6, 9 and 12. Stochastic gradient descent (SGD) with momentum (0.9) and weight decay (0.0005) was used for optimization. Our work is based on the torchvision implementation of Mask R-CNN [23]. In Table. 1, the accuracies of object detection in terms of AP_{50} (the average precision with IoU thresholded at 0.50), $AP_{50:95}$ and $AR^{max=100}$ (the average recall with 100 detections per image) are tabulated. While training on the naive dataset does not generalize well to our heavily cluttered real scenarios, Mask R-CNN trained on the realistic dataset considerably boosts the performance.

4.2. Experiments on pose estimation

To this goal, we trained the autoencoder with a latent space size of 128. We chose the L2 loss function, a learning rate of 0.0001 and used the Adam optimizer with a batch size of 32 and trained it for 40K iterations. The pose error metrics used for evaluation are the Visible Surface Discrepancy (VSD), the Maximum Symmetry-aware Surface Distance (MSSD) and the Maximum Symmetry-aware Projection Distance (MSPD), that are being used in the BOP2020 challenge [11]. An estimated pose is considered as correct w.r.t. the pose-error function e if $e < \theta_e$, where $e \in \{e_{VSD}, e_{MSSD}, e_{MSPD}\}$ and θ_e is the threshold of correctness. We used the same values for θ_e as

Object	sorted by Mask R-CNN scores	sorted by max cosine sim.	sorted by depth differences (e_i)
Object 1	0.509	0.394	0.812
Object 2	0.533	0.449	0.633

Table 2: Average recall of top 5 pose estimates sorted by three different approaches for selecting the best estimates

Object	RGB	RGB + depth		
		+ ICP	normal	mult.-obj. mask
Object 1	0.691	0.829	0.812	0.790
Object 2	0.348	0.703	0.633	0.627
time (s)	0.699	18.39	0.697	

Table 3: Average recall of top 5 pose estimates using ICP and depth measurements and combinations.

in [11] to calculate the average recall rates AR_{VSD} , AR_{MSSD} and AR_{MSPD} . The performance of the method on a dataset is measured by the Average Recall $AR = (AR_{VSD} + AR_{MSSD} + AR_{MSPD})/3$.

In Table 2, we compare the results of the different methods on selecting the best five pose estimates. In these experiments, we did not make use of ICP, but we took the depth measurement at the center of the object. It shows that sorting according to the vector (e_i) in (2), results in superior performance compared to other approaches.

The experiments in Table 3 are conducted using the selection method defined by the vector (e_i) . We compare the results, when testing with only RGB information, using the depth measurement at the object center and the improvement through ICP refinement. While ICP refinement increases the performance slightly, the refinement of several hundred poses per image takes time, making it impractical for the usage in real-time settings. The experiments to reduce the noise in cluttered scenes, like feeding only the pixels visible in the mask to the autoencoder, or training the autoencoder with multiple objects, have shown, against our intuition, a worse performance than the normal pipeline. Note how incorporating the depth has improved the accuracy. We have shown qualitative results in the supplementary materials.

5. CONCLUSION

In this paper, we investigated the task of bin picking from piles of crowded, small-sized and identical objects. In particular, we explored the main required vision modules for this challenging problem, i.e. object detection and pose estimation. For each task, we employed convolutional neural networks, which are trained on two types of generated synthetic datasets. Experimental results on synthetic and real images show that the proposed comprehensive framework, from dataset generation to pose estimation, is promising for industrial bin picking.

6. REFERENCES

- [1] W. Abbeloos and T. Goedemé, “Point pair feature based object detection for random bin picking,” in *2016 13th Conference on Computer and Robot Vision (CRV)*. IEEE, 2016, pp. 432–439.
- [2] D. Liu, S. Arai, J. Miao, et al., “Point pair feature-based pose estimation with multiple edge appearance models (ppf-meam) for robotic bin picking,” *Sensors*, vol. 18, no. 8, pp. 2719, 2018.
- [3] S. Hinterstoisser, S. Holzer, et al., “Multimodal templates for real-time detection of texture-less objects in heavily cluttered scenes,” in *2011 international conference on computer vision*. IEEE, 2011, pp. 858–865.
- [4] M. Sundermeyer, Z. Marton, M. Durner, and R. Triebel, “Augmented autoencoders: Implicit 3d orientation learning for 6d object detection,” *International Journal of Computer Vision*, vol. 128, no. 3, pp. 714–729, 2020.
- [5] T. Hodan, P. Haluza, et al., “T-less: An rgb-d dataset for 6d pose estimation of texture-less objects,” in *2017 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2017, pp. 880–888.
- [6] Y. Labbé, J. Carpentier, M. Aubry, and J. Sivic, “Cosypose: Consistent multi-view multi-object 6d pose estimation,” in *European Conference on Computer Vision*. Springer, 2020, pp. 574–591.
- [7] B. Joffe, T. Walker, R. Gourdon, and K. Ahlin, “Pose estimation and bin picking for deformable products,” *IFAC-PapersOnLine*, vol. 52, no. 30, pp. 361–366, 2019.
- [8] K. He, G. Gkioxari, P. Dollár, and R. Girshick, “Mask r-cnn,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2961–2969.
- [9] K. Park, T. Patten, and M. Vincze, “Pix2pose: Pixel-wise coordinate regression of objects for 6d pose estimation,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 7668–7677.
- [10] 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 International Conference on Computer Vision*, 2019, pp. 7678–7687.
- [11] T. Hodaň, M. Sundermeyer, B. Drost, Y. Labbé, et al., “Bop challenge 2020 on 6d object localization,” in *European Conference on Computer Vision*. Springer, 2020, pp. 577–594.
- [12] T. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, “Focal loss for dense object detection,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2980–2988.
- [13] S. Ren, K. He, R. Girshick, and J. Sun, “Faster r-cnn: Towards real-time object detection with region proposal networks,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 6, pp. 1137–1149, 2016.
- [14] 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*, 2017, pp. 1521–1529.
- [15] A. Pashevich, R. Strudel, I. Kalevatykh, I. Laptev, and C. Schmid, “Learning to augment synthetic images for sim2real policy transfer,” *arXiv preprint arXiv:1903.07740*, 2019.
- [16] J. Tobin, R. Fong, A. Ray, et al., “Domain randomization for transferring deep neural networks from simulation to the real world,” in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2017, pp. 23–30.
- [17] M. Denninger, M. Sundermeyer, D. Winkelbauer, Y. Zidan, and D. Olefir, “Blenderproc,” *arXiv preprint arXiv:1911.01911*, 2019.
- [18] Z. Zhang, “Iterative point matching for registration of free-form curves and surfaces,” *International journal of computer vision*, vol. 13, no. 2, pp. 119–152, 1994.
- [19] M. Everingham, S. M. A. Eslami, and others., “The pascal visual object classes challenge: A retrospective,” *International Journal of Computer Vision*, vol. 111, no. 1, pp. 98–136, Jan. 2015.
- [20] W. Liu, D. Anguelov, D. Erhan, et al., “Ssd: Single shot multibox detector,” in *European conference on computer vision*. Springer, 2016, pp. 21–37.
- [21] J. M Wong, V. Kee, T. Le, et al., “Segicp: Integrated deep semantic segmentation and pose estimation,” in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2017, pp. 5784–5789.
- [22] X. Deng, Y. Xiang, A. Mousavian, et al., “Self-supervised 6d object pose estimation for robot manipulation,” in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 3665–3671.
- [23] S. Marcel and Y. Rodriguez, “Torchvision the machine-vision package of torch,” in *Proceedings of the 18th ACM international conference on Multimedia*, 2010, pp. 1485–1488.

- [24] J. Vidal, C. Lin, and R. Martí, ,” in *6D pose estimation using an improved method based on point pair features*. IEEE, 2018, pp. 405–409.