

ZebraPose: Coarse to Fine Surface Encoding for 6DoF Object Pose Estimation

Yongzhi Su^{1,2*} Mahdi Saleh^{3*} Torben Fetzner² Jason Rambach¹
Nassir Navab³ Benjamin Busam³ Didier Stricker^{1,2} Federico Tombari^{3,4}

¹ German Research Center for Artificial Intelligence (DFKI) ² TU Kaiserslautern
³Technische Universität München ⁴Google

{yongzhi.su; jason.rambach; torben.fetzner; didier.stricker}@dfki.de

{m.saleh; b.busam; nassir.navab}@tum.de, tombari@in.tum.de

Abstract

Establishing correspondences from image to 3D has been a key task of 6DoF object pose estimation for a long time. To predict pose more accurately, deeply learned dense maps replaced sparse templates. Dense methods also improved pose estimation in the presence of occlusion. More recently researchers have shown improvements by learning object fragments as segmentation. In this work, we present a discrete descriptor, which can represent the object surface densely. By incorporating a hierarchical binary grouping, we can encode the object surface very efficiently. Moreover, we propose a coarse to fine training strategy, which enables fine-grained correspondence prediction. Finally, by matching predicted codes with object surface and using a PnP solver, we estimate the 6DoF pose. Results on the public LM-O and YCB-V datasets show major improvement over the state of the art w.r.t. ADD(-S) metric, even surpassing RGB-D based methods in some cases.

1. Introduction

Augmented reality and robotics are two of the main application fields of 3D computer vision. In many augmented reality applications, the location and pose of an object of interest has to be determined at a high precision [46, 55]. Similarly, object grasping and manipulation is needed for many robotic applications (e.g. automatic manufacturing [51], cooperative assistance [9, 21]), and also demands accurate 6 Degree-of-Freedom (6DoF) object pose information. As the crucial element in both application domains, estimating the 6DoF object pose has received increasing attention from the computer vision research community.

The correspondence problem is a classical problem in

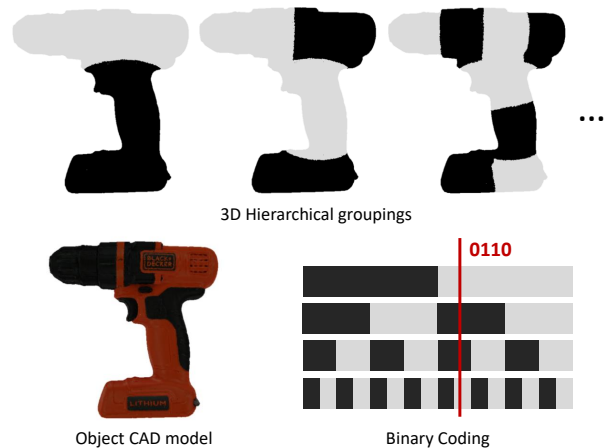


Figure 1. ZebraPose assigns a discrete code to each surface vertex hierarchically. We project the code as binary black and white values (top) and learn them using deep neural networks. Our binary descriptor allows one-to-one correspondence for the problem of 6DoF object pose efficiently.

computer vision. While finding correspondences across the same domain is more straightforward, estimating the 6DoF object pose requires 2D-3D correspondences. In earlier object pose estimation research, depth maps came to help to match image pixels to 3D surface points [28, 65]. Due to cost and setup complications, the detection of 6DoF pose without depth information can be advantageous. However, RGB approaches typically achieve a lower accuracy with respect to their depth-based counterparts [16, 26].

Driven by the recent developments in deep learning and Convolutional Neural Networks (CNNs), various methods were proposed, which make 6DoF pose estimation from a single RGB image feasible [10, 34, 54, 62, 70]. In a correspondence-based setting, to estimate the object pose, Perspective-n-Points (PnP) algorithms require at least 4 2D-3D point matches [39]. Therefore, sparse methods are ap-

*The authors contributed equally to this paper

Code: <https://github.com/suyz526/ZebraPose>

plied to extract points of interest [48, 53]. However such methods might fail to find object landmarks under viewpoint changes, occlusion, or lack of texture. With the success of deep neural networks in image synthesis problems, researchers use such tools to generate dense correspondence maps. For instance, several methods learn UV [72] or UVW [49, 68] values in object local coordinates. Since the network produces dense smooth results, certain low-level geometry is lost. Moreover, neural networks tend to achieve a higher performance in classification tasks [34].

To this aim, we propose a dense correspondence pipeline that combines the concepts of handcrafted features and image segmentation in a hierarchical fashion for RGB-based 6DoF pose estimation. In order to design a descriptor that encodes surfaces efficiently, we use the binary numeral system. Binary-based descriptors are applied in ORB [57] and are still in use in robust SLAM applications [12]. In our work, we split the surface into halves in multiple iterations and define our vertex encoding by stacking the assigned group labels. By leveraging a hierarchical discrete representation, we guarantee a compact mapping and simple learning objective as a multi-label classification problem [29, 34]. Moreover, learning how to encode a full sequence at once might be challenging for neural networks. Therefore, we propose a coarse to fine learning scheme. By design, our encodings on the coarse levels are continuously shared in wider object regions. As the network learns to differentiate coarse splits, we focus on finer encoding positions. With a coarse to fine loss and training strategy, we then manage to predict fine-grained surface correspondences.

In contrast to previous works where there is no guaranteed putative correspondence [49, 50, 68], our encoding promotes direct pixel-to-surface matching just by means of a look-up table. With a simple matching and PnP-RANSAC scheme of Progressive-X [2], we outperform the state of the art in 6DoF pose on the most commonly used benchmarks w.r.t. ADD(-S) metric.

In summary, we propose ZebraPose, a two-stage RGB-based approach that defines the matching of dense 2D-3D correspondence as a hierarchical classification task. We divide the general two-stage approach for 6DoF object pose estimation into three components: 1) assigning a unique descriptor to the 3D vertex; 2) predicting a dense correspondence between the 2D pixels and 3D vertices; 3) solving the object pose using the predicted correspondences. We can summarize our proposed contributions in this paper related to the first two components:

- A novel coarse to fine surface encoding method assigning the dense vertex descriptor in an efficient way, which also fully exploits traditional outlier filters used in computer vision task.
- A novel hierarchical training loss and strategy to auto-

matically adjust the weights of each code position.

Extensive experiments on LM-O [3] and YCB-V [11] datasets show that our proposed approach achieves state of the art results.

2. Related Work

We limit our in-depth discussion of related work to the most relevant methods to our work, i.e. RGB-based 6DoF pose estimation, and object surface encoding techniques.

2.1. RGB-based 6DoF Pose Estimation

Traditional Methods. With the development of the feature descriptor [43], the object pose problem could be solved by feeding estimated 2D-3D correspondences into a RANSAC/PnP framework. However, dealing with textureless objects remained a challenge. To overcome the lack of keypoints, Hinterstoisser *et al.* [25] proposed to utilize the image gradient information and formulate the pose estimation task within a template matching pipeline. Later advances [5] avoided the template searching time by applying a statistical learning-based framework to regress object coordinates and object labels jointly. However, the accuracy that handcrafted methods can achieve is far from that of deep learning methods nowadays.

End-to-End Methods. PoseNet [35] was the first work that attempted to regress the camera viewpoint with a CNN. Following works usually concatenated an object detector with the pose regression, making multi-object pose estimation possible [70]. Finding a suitable rotation representation for pose regression was a problem at that time and typical rotation parametrization did not populate Euclidean spaces [8]. SSD6D [34] avoided complex parameters by discretization of the rotation space thus treating the rotation estimation as a classification problem. Zhou *et al.* [74] proposed a continuous 6-dimensional rotation representation that shows advantages over quaternions [44, 45] or Lie algebra [20, 61] parametrization for neural network training. This representation is utilized in several direct regression works [19, 37, 67].

In parallel, several efforts have been made to integrate RANSAC and PnP modules to pose learning frameworks. [4, 6, 7] propose differentiable RANSAC variants, which are not applicable to object pose estimation as they require a good initialization and complex training strategy. [30] proposes a network to solve the PnP problem, with a loss function reflecting pose metrics. At the same time, a new branch of methods has been developed with the growth of neural renderers [15, 33, 42]. [32] is able to define the loss according to the texture colour on pixel level. [59, 66] used a differentiable depth map and achieved self-supervised network fine-tuning with unlabeled RGB-D data. In an effort to combine correspondence-based methods with direct regression

of 6DoF parameters, [67] used correspondence maps as an intermediate geometric representation to regress the pose. [19] further enhances [67] by employing self-occlusion information that provides richer information to predict the object pose with the predicted 2D-3D correspondences.

Indirect Methods with Deep Learning. While end-to-end methods have evolved through time by integrating differentiable modules, the performance of such methods are normally below geometrical and indirect methods. Combining learning features and geometrical fitting, [69] uses metric learning to learn an implicit pose representation through triplet loss and finally looks for nearest neighbors in pose space. AAE [62] learns to generate a latent vector based on the visual information of the object in discrete viewpoints. At inference stage, the rotation is obtained by comparing the latent code with the pre-generated rotation-latent code lookup table. The rest of the indirect methods usually estimate the 2D-3D correspondence, and solve the object pose using RANSAC/PnP. BB8 [53] firstly defines the 3D object bounding box corners as the keypoints and PVNet [50] reaches high recall rate in LM [27] dataset by predicting the keypoints with a dense pixel-wise voting for sampled keypoints on the object. The main drawback of such sparse 2D-3D correspondence methods is that the prediction of keypoints in the occluded area lacks in accuracy. HybridPose [60] proposed to leverage multiple geometric information to tackle this issue while other methods [29, 49, 72] predict pixel-wise dense 2D-3D correspondences.

2.2. Surface Encoding

The binary surface encoding technique has been successfully used in the field of structured light reconstruction for many years [47, 52, 58, 64]. For this purpose, a video projector illuminates the scene with several successively refined binary fringe patterns. The composition of the different stripe patterns provides an encoding of the surface points. Surface coding using multiple classification problems has proven to be highly reliable and competitive [22]. Since neural networks are ideally suited for solving classification problems, a transfer of the approach as we presented in this work constitutes a logical step.

In pose estimation domain, to estimate the dense 2D-3D correspondence, each 3D corresponding point must be assigned a unique descriptor. Pix2Pose [49] simply treats the 3D vertex coordinates as this descriptor. DPOD [72] textures the object with a 2-channel UV-map with discrete values to learn the correspondences. EPOS [29] divides the object surface into multiple fragments, and estimates the corresponding points by combining fragment segmentation and local fragments coordinates prediction. Although most of these encodings are limited to local object coordinates, we propose a method that learns dense 2D-3D correspondence through a handcrafted code. Compared to methods

that predict local coordinates space [68, 72] in 2D or 3D grid, we encode the object surface in a coarse to fine manner. Moreover, unlike EPOS [29] that divides the object surface into multiple coarse bins at once, we divide the object surface iteratively until the fragments are fine enough to define the unique 3D corresponding point. This allows for gradual refinement of the correspondence through the hierarchical levels.

3. Method: ZebraPose

In this section, we present our approach for the problem of 6DoF object pose, which involves the entire process from our surface encoding to the final pose estimation.

3.1. Coarse to Fine Surface Encoding

Given a surface CAD model of an object and its vertices $v_i \in \mathbb{R}^3$, where i stands for the vertex id, we want to represent each v_i with vertex code $c_i \in \mathbb{N}^d$, where d is the length of the vertex code. We need to define such encoding based on vertices' position relative to the given 3D object surface to enable coarse to fine learning. To enable this, we construct our codes in a non-decimal numeral system.

Defining our encoding in a numeral system with lower radix makes the representation very efficient and provides easier grounds for coarse to fine grouping of the points. For a code of length d , we perform d iterations of grouping of the vertices. The collection of groups G_j in the j -th iteration ($j \in \{0, \dots, d\} \subset \mathbb{N}$) consist of r^j groups. G_0 defines the initial group, including only one group, i.e. the entire object vertices. G_j with $j > 1$ is obtained by splitting each group in G_{j-1} into r groups. In a grouping iteration, each vertex v_i is assigned with a class id $m_{i,j}$, where $m_{i,j} \in \{0, \dots, r-1\} \subset \mathbb{N}$ based on the group it belongs to in the j -th grouping. Finally, each vertex is assigned to a vertex code with d digits by stacking the class id of each grouping operation. This representation is stored and fixed for every 3D object. The vertices in each group share the same code. We build the lookup table to map a code to the centroid of the each group in G_d , which is further used to build 2D-3D correspondence and solve the pose as described in Sec. 3.6. In this paper, we used k-means for the grouping, more details are in Sec. 4.1. We illustrate this process in Fig. 2 with $r = 2$ and break down the CAD model surface into discrete and equally sized groups.

3.2. Choice of the Radix for Vertex Code

Following our grouping described in Sec. 3.1, we would have K total number of classes, where $K = r^d$. In a classification problem we learn these maps using o logits, where $o = r \cdot d$. To minimize the number of outputs while learning the most number of classes we have:

$$o_{min} = \min_r r \cdot d = \min_r r \cdot \log_r K = 2 \cdot \log_2 K. \quad (1)$$

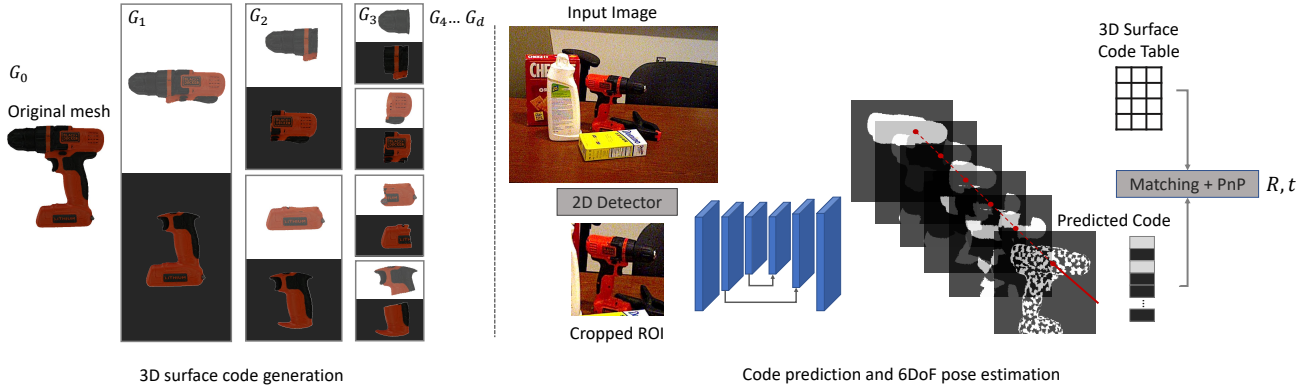


Figure 2. Left: Our hierarchical encoding is defined by grouping surface vertices in several iterations. In each iteration, object vertices are split into equally sized groups. In a binary setting, vertices are classified into two groups, 0 (white) and 1 (black). This process happens offline and the generated mapping between vertex code and the corresponding 3D vertex is stored in a look-up table. Right: Our training framework uses a detector to crop the object ROI and predicts a multi-layer code using a fully convolutional neural network. The predicted code is then matched to the 3D surface vertex and passed to RANSAC and PnP modules for pose estimation.

The best positive integer choice of r to minimize the number of network layers are 2 and 4. Since a value is classified either as positive or negative, we do not need to use the cross entropy loss with 2 explicit output layers for the binary classification. So we can reach $\log_2 K$ as the optimal number of output layers with $r = 2$.

Besides the advantages of reduced GPU memory requirement, we show later in the ablation study (see Sec. 4.2) that using the binary vertex code yields the most accurately predicted pose. Thus we choose a binary base for the vertex code.

3.3. Rendering the Training Labels

Each object pixel in the image corresponds to a 3D object vertex. The network predicts the class id that is assigned to this vertex in each grouping operation. Therefore, we still need to render the class id into the 2D image plane with a given pose for the training. For this purpose, we transfer the class id of vertices into the class id of the mesh faces using the following criteria: if two vertices of a face have the same class id, the face is assigned with this class id. Otherwise, the face has the class id of its first vertex. We repeat this rendering process for d times until the training label class id for each grouping is generated.

3.4. Network Architecture

In Sec. 3.2 we justify our choice of $r = 2$. In this regard, our goal is to classify 2^d regions with only d binary values.

During training, we use the object pose annotations to render the labels as layered black and white maps to image coordinates. This way, our objective learning maps are $d + 1$ binary labels (d for the binary vertex code and 1 for the object mask) for code and visible mask prediction. An

encoder-decoder network generates $d + 1$ outputs with a single decoder. We round the final output probabilities to represent our discrete vertex codes.

The entire process from input images to the predicted pose is presented in Fig. 2. To predict the code per pixel in the frame with fine granularity, we process only a Region of Interest (ROI) around object pixels. Following the pipeline similar to [37, 41, 67], we focus on the object pose and use the available 2D detector predictions to find the ROIs. We crop and resize the ROI from the prediction to a fixed dimension $H \times W$, and apply the exact process to the target vertex code maps during training. Our goal is to predict multiple labels per frame in the ROI.

3.5. Hierarchical Learning

Predicting correspondences directly from object pixels is a fine-grained task. On the other side, deep neural networks are commonly used for coarse level predictions. This means features predicted per pixel are very similar in a small vicinity. As our encoding is also hierarchical by design, we learn the codes in a coarse to fine manner. Therefore, the predictions are learned in different stages, from coarse groupings to fine ones. We use an error histogram for each position on hierarchical level and weight our Hamming-based loss given the error to design this.

Mask loss. Firstly, we predict the visible mask to segment the object area from the background. Here, we simply pass the predicted probability to the sigmoid function and use L1 loss as L_{mask} . It is worth noting, for the binary vertex code prediction in the following, we only calculate loss of the pixels within the predicted object mask.

Hamming distance: The CNN outputs the binary vertex code probabilities $\hat{p} \in \mathbb{R}^d$ for a pixel within a ROI, we

obtain the predicted discrete binary code \hat{b} by rounding \hat{p} . Given \hat{b} and its known ground truth binary vertex code b , the Hamming distance Ham is defined by counting the number of bits \hat{b} which are different from b . This formulation does not favor any of the positions and calculates the error without considering any hierarchical information explicitly. As a common practice in deep learning, we use binary cross-entropy as an activation function for the distance:

$$Ham(b, \hat{p}) = \sum_{j=1}^d b_j \log \hat{p}_j + (1 - b_j) \log(1 - \hat{p}_j), \quad (2)$$

where b_j stands for the j -th bit in b (the j -th bit is generated in the j -th vertices grouping).

Active bits. Lower bits in binary vertex code b hold coarse correspondences, and higher bits define finer estimates. During the initial training phase, the network focuses on learning the coarse splits and has a higher error on fine bits. Therefore we adaptively weight the coarse bits by looking at the histogram of error of all bits. As the training proceeds and coarser predictions become more robust, finer bits are induced with more weights. We define our histogram at training step t by looking at the error at different bits:

$$H_j(t) = avg(\lambda(b_j^t - \hat{b}_j^t) + (1 - \lambda)(b_j^{t-1} - \hat{b}_j^{t-1})), \quad (3)$$

where \hat{b}_j^t defines the predicted binary vertex code \hat{b} at training step t , and λ is a constant. With the *avg* operator, we get the error ratio by calculating the average difference in b_j^t and \hat{b}_j^t of all pixels within the predicted object mask in a mini-batch. During training we update the histogram given the previous histogram in training step $t - 1$ and the current error histogram. We show how to define a hierarchical loss based on the histogram in the following.

Hierarchical loss. We compute a weighting coefficient based on the error histogram, and use it on top of a Hamming distance to form our hierarchical loss with

$$w_j(t) = \exp(\sigma \cdot \min\{H_j(t), 0.5 - H_j(t)\}), \quad (4)$$

where the function w uses an exponential term to softly define a weight for n -th bit at training step t , σ is a constant. All object pixels in the mini-batch share the same weighting coefficients. We normalize the weights across all bits. We then define our hierarchical loss based on the weighting function of active bits and Hamming distance as below:

$$L_{hier} = \sum_{j=1}^d w_j \cdot Ham(b_j, \hat{p}_j). \quad (5)$$

With this loss we focus mainly on active bits which automatically change from coarse to fine during training.

Total loss to train the CNN. We weight the L_{mask} and L_{hier} with a hyper-parameter α (α set as 0 for pixels predicted as background), the total per pixel loss can be mathematically expressed as

$$L_{total} = L_{mask} + \alpha \cdot L_{hier}. \quad (6)$$

3.6. Pose estimation

In previous sections, we discussed how to generate our descriptor and learn to predict them using a fully convolutional neural network. Now we incorporate the predicted code and visible mask and the reference 3D model encoding to match correspondences. Different from common dense correspondences such as [49,67,68], this compact representation also enables a bijective correspondence between the surface vertices and the descriptor space. That means, unlike the regressed 3D point which can be off the object surface, our estimated 3D correspondences always refers to a vertex on the object model, which eases the matching stage for the pose solver. For the matching, we use a look-up table that extracts the corresponding 2D and 3D points. Following that, we use Progressive-X [2] solver to calculate the rotation R and translation t .

4. Experiments

In this section, we firstly introduce the implementation details, the datasets and metrics used for the evaluation. Subsequently, we present ablation study experiments on the LM-O [5] dataset. Finally, we compare our experimental results with state of the art methods on the LM-O [5] and YCB-V [70] datasets. Please refer to supplementary materials for more qualitative results.

4.1. Experiments Setup

Implementation Details. In order to have the same number of classes as DPOD [72] ($K = 256^2$), we firstly upsample the mesh by subdivision of each face using the edge’s midpoint [14] until the mesh has more than 256^2 vertices. Subsequently, we group the 3D vertices of the object model as we described in Sec. 3.1 with $r = 2$ and $d = 16$. After several iterations of the grouping operation, a group could contain fewer points than 2 and cannot be grouped further. To avoid this, we modified the k-means++ clustering algorithm [1], to force both output groups of points to have equal sizes.

We modified Deeplabv3 [13] by adding skip connections and used Resnet34 [23] as the backbone. The input ROI is resized to the shape of $256 \times 256 \times 3$, and the CNN output has a height and width of 128. We applied the same dynamic zoom-in strategy as CDPN [41] to generate the noisy ROI for the training. The parameter λ used in the histogram is 0.05. The parameter σ used in the hierarchical loss is 0.5, and α has been set as 3 to balance the training for mask and

vertex code prediction. The CNN has been trained 380k steps using the Adam optimizer [36] with a batch size of 32 and a fixed learning rate of $2e-4$. During the inference stage, we utilize the detected bounding box with Faster R-CNN [56] and FCOS [63] provided by CDPNv2 [41]. If not specified, we used detected bounding box from Faster R-CNN in the ablation study.

Additionally, by changing any bit in the vertex code, the code refers to another 3D point, possibly even to a vertex on the other side of the object. To maintain the topology presented with the ground truth correspondence map, we disabled the interpolation during the rendering when we generated the ground truth. The resizing of ground truth is also done with nearest neighbourhood interpolation in the training stage.

Datasets. The reported recall rate in LM [28] dataset has lately been higher than 95% and quite saturated, therefore we focus on the more challenging LM-O [5] and YCB-V [70] dataset in this paper. LM-O consist of 1214 images and is only used as test images. LM-O annotated 8 objects poses in the images under partial occlusion, making pose estimation more challenging. About $1.2k$ images per object in LM are used as the real training images for LM-O. Compared to LM-O, YCB-V is a large dataset containing 21 objects. Although YCB-V provides more real training images, the objects are strongly occluded in the scene, and many of the objects are geometrically symmetric.

Since the LM-O dataset includes only a limited number of training images, [34, 50] additionally render a large number of synthetic images for training. However, due to the domain gap between the synthetic and real images, the performance of the methods also heavily depends on the domain randomization and domain adaptation technique [62, 71]. As the physically-based rendering (pbr) training images [18] for both datasets are publicly accessible now, using pbr images to support the training can help us focus on the pose estimation CNN itself. We use the pbr images together with the real images for the training in the same manner as [19, 29, 67].

Error Metrics. We selected the ADD(-S) error metric as the most commonly used metric for the 6DoF pose estimation task. This metric calculates the average distance of model points projected to the camera domain using the predicted pose to the same model points projected using the ground truth pose. For symmetric objects, the metric matches the closest model points projected with the ground truth pose instead of the same model point. In all the experiments in this paper, if ADD(-S) error is smaller than 10% (most commonly used threshold) of the object diameter, the predicted pose is considered to be correct. For YCB-V, we also reported the AUC (area under curve) of ADD(-S) with a maximum threshold of 10 cm [70].

4.2. Ablation Study on LM-O

In this section, we present the results of several ablation studies as follows:

Length of Binary Vertex Code. The object 3D surface is encoded through iterative k-means++ clustering until the size of the segmented cluster is small enough so that we can map the vertex code to the centroid of each cluster. We used the same total number of classes as DPOD [72], which means each binary vertex code has 16 bits. However, if the objects are small or the distance of the object to the camera is too large, different clusters in the fine level could be rendered into the same pixel when we generate the ground truth data. This makes the binary code in the fine levels (the last few bits) redundant. Due to the distance variation of the object to the camera, we can not determine which bits are redundant.

In this ablation study, we research which bits are the redundant bits. The models are trained without the hierarchical training strategy, and we use Progressive-X [2] to solve the pose. We ignore the last few bits of the predicted binary code in the inference stage. The new binary vertex code with fewer bits refers to a larger point cloud group (see Fig. 2 left). We calculated the new centroid of the group and reassigned the centroid as the corresponding 3D points for the binary vertex code with fewer bits. From Fig. 3 b) we can see that using 10-bit code is already sufficient to yield an accurate prediction for the objects in LM-O, indicating the last 6 bits are redundant for those objects. The results fluctuated a bit when we applied the redundant bits, which indicates that for some objects the best results is achieved by not using the full 16-bit code. However, in the following experiment, we always report the results with the full predicted vertex code.

Radix used in Vertex Code. The number of the clusters in each iteration decides the radix of the generated vertex code that describes the 3D vertex. Since our CNN predicts the vertex code, it is meaningful to compare which radix in vertex code suits the representation better. We do not need to generate all the vertex codes used in this ablation study from scratch. More specifically, by merging every $\log_2 r$ bit of a vertex code, we get a code with a radix r . For instance, a vertex with a binary code $\{11111110\ 11111111\}$ can be transformed to $\{254\ 255\}$ using 256 as the radix. We will get exactly the same code for this vertex if we split the object into 256 groups and split each group again into 256 groups. We use a fix $w_j = 1$ (see eq. 4) for all positions so the loss is essentially a binary cross entropy when $r = 2$, and cross-entropy loss for other radices.

We present the comparison results in Tab. 1. If RANSAC/PnP is used to solve the pose, the results with different radices are quite similar. There is no clear indication whether using the small or large radix is better. If we switch the pose solver to Progressive-X [2], the code

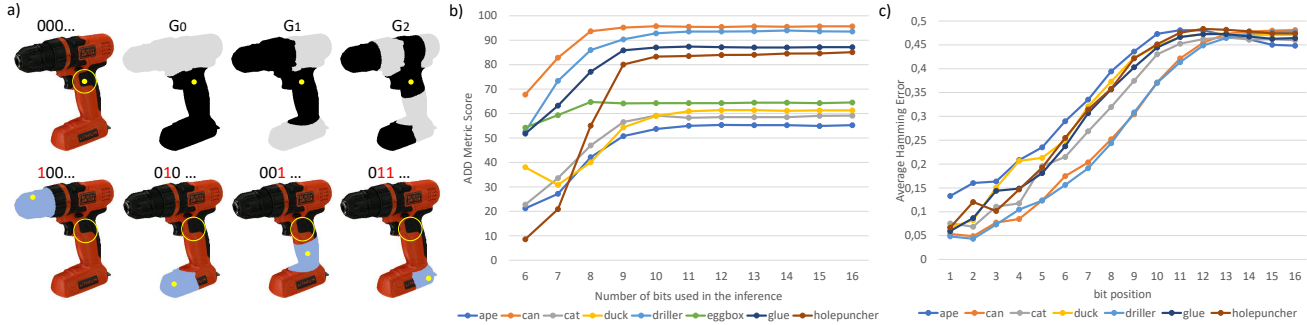


Figure 3. a) In the first row, a yellow dot has the ground truth binary vertex code beginning with 000, and the yellow circle refers to the neighborhood vertex of this yellow dot. If the vertex code has been predicted as 100 (the first bit is wrong, marked as red in the figure) during the inference stage, the estimated yellow dot lies somewhere on the head of the drill (marked in blue). The estimated 3D vertex is far away from its original neighborhood and can be easily found by checking the spatial coherency. We show four similar cases in this figure. b) We calculated ADD pose metrics only on the first j bits of the predicted code to build the 2D-3D correspondence. Here we observe from which bit the predictions are stable. c) We present the average error rate at different bit positions on the LM-O dataset [5].

with small radix improves the most and yields the best accuracy. Progressive-X solver includes a spatial coherence filter that checks neighboring 3D points with respect to its assigned 2D correspondences based on label cost energy minimization as introduced in [17]. It can therefore deal particularly well with the type of outliers, from our method as we mentioned in Fig. 3. We show in the Fig. 3 a), if the CNN predicts the first few bits wrong for the yellow dot, the estimated corresponding 3D points is far away from the ground truth position and totally incoherent with its original neighbourhood. Assuming that most predictions are correct, most neighbourhood vertices are posed within the yellow circle in the figure. In this case, this false estimated 3D corresponding can be easily filtered by calculating the coherency with its neighbourhood. This spatial coherency filter can not detect outliers well if the wrong prediction is in the last few bits, and the same for 256 as radix, as divide the vertices into 256 groups is already a fine grouping. Nevertheless, the Fig. 3 b) already shows that the last few bits do not affect the solved pose. So we argue that the binary vertex code suits this task the best. Moreover, the prediction of binary codes requires the least RAM in GPU, as we discussed in Sec. 3.1.

Effectiveness of Hierarchical Training. According to the first ablation study, the last few bits are redundant and may not be trainable (see Fig. 3c)). During the training, we can recognize redundant bits based on the error histogram and focus on the decisive bits as described in Sec. 3.5. Tab. 2 shows that the results are further improved by our proposed hierarchical training.

Influence of 2D detection. The CNN estimates the pose with the cropped ROI from the detected bounding box. The object pose estimation is meaningless with a false-positive detection, also the pose is not even estimated in the case of false-negative detection. By leveraging the detected bound-

Method	RANSAC/ Pnp [39]	Progressive-X [2]
2 as radix	73.06	75.23 (+2.17)
4 as radix	72.94	74.59 (+1.65)
16 as radix	73.04	74.98 (+1.94)
256 as radix	73.25	74.52 (+1.27)

Table 1. **Ablation study on LM-O [5].** We tested the use of different radices to encode the vertices, and using different solvers to calculate the pose. The results are presented in terms of average recall of ADD(-S) in %.

Method	ADD
2 as radix	75.23
2 as radix + Hierarchical Learning	75.86
2 as radix + Hierarchical Learning + Faster R-CNN [56] → FCOS [63]	76.91

Table 2. **Ablation study on LM-O [5].** We compare the result with and w/o applying our hierarchical loss, as well as the impact of the prior object detector. The results are presented with average recall of ADD(-S) in %.

ing box with FCOS [63] instead of the one from Faster R-CNN [56], the recall rate improved 1.05%.

4.3. Comparison to State of the Art

We use 2 as radix, i.e. binary vertex code and apply the hierarchical training strategy and Progressive-X pose solver [2] in our proposed ZebraPose to compare to state of the art on LM-O [5] and YCB-V [70] datasets. The detected bounding box of FCOS [63] detector are provided by CDPNv2 [41].

Results on LM-O. We report the recall of ADD(-S) met-

Method	RGB Input				RGB-D Input		
	HybridPose [60]	RePose [32]	GDR-Net [67]	SO-Pose [19]	Ours	PR-GCN [73]	FFB6D [24]
ape	20.9	31.1	46.8	48.4	57.9	40.2	47.2
can	75.3	80.0	90.8	85.8	95.0	76.2	85.2
cat	24.9	25.6	40.5	32.7	60.6	57.0	45.7
driller	70.2	73.1	82.6	77.4	94.8	82.3	81.4
duck	27.9	43.0	46.9	48.9	64.5	30.0	53.9
eggbox*	52.4	51.7	54.2	52.4	70.9	68.2	70.2
glue*	53.8	54.3	75.8	78.3	88.7	67.0	60.1
holepuncher	54.2	53.6	60.1	75.3	83.0	97.2	85.9
mean	47.5	51.6	62.2	62.3	76.9	65	66.2

Table 3. **Comparison with State of the Art on LM-O [5]**. We report the Recall of ADD(-S) in % and compare with state of the art. (*) denotes symmetric objects.

ric in Tab. 3. We ordered the methods according to the input modality. HybridPose [60] and RePose [32] have been trained with synthetic and real images. GDR-Net [67] also reported their recall of 53% when trained with synthetic and real images. Therefore, GDR-Net outperforms HybridPose and RePose. In our Tab.3, we report the best results that GDR-Net and SO-Pose [19] presented, which are also trained with pbr and real images. GDR-Net used faster R-CNN [56] as the detector, ZebraPose yields a recall of 75.86% with faster R-CNN (see Tab. 2), which can be seen as a more fair comparison with GDR-Net.

To summarize, our ZebraPose outperforms state of the art RGB based methods with a large margin on this dataset. Additionally, we found that our ZebraPose also outperforms state of the art RGB-D based methods [24,73]. Most objects in the LM-O dataset are texture-less, meaning that RGB-D based methods should have more advantage in feature extraction on the objects with the help of depth image. Even in this case, our results still exceed theirs.

Results on YCB-V. We compare ZebraPose with other approaches in the YCB-V dataset in Tab. 4. The AUC reported in Tab. 4 has been calculated using all-points interpolation. Tab. 4 shows that ZebraPose is still better than state of the art w.r.t. ADD(-S) and AUC of ADD(-S) metrics and comparable to them w.r.t. the AUC of ADD-S metric.

4.4. Runtime Analysis

We tested the runtime on a desktop with an Intel 3.50GHz CPU and an Nvidia 2080Ti GPU. The CNN runtime plus the time to build the 2D-3D correspondence is about 52 ms. The FCOS detector [63] takes 55 ms. RANSAC/PnP [39] needs only 4 ms to solve the pose, while Progressive-X [2] requires 150 ms to obtain the pose. So for ZebraPose used in Sec.4.3, it totally needs about 250 ms to estimate the object pose. If we use RANSAC/PnP to solve the pose, the runtime reduces to 110 ms, while with about 2.6% recall drop on LM-O dataset.

Method	ADD(-S)	AUC of ADD-S	AUC of ADD(-S)
SegDriven [31]	39.0	-	-
SingleStage [30]	53.9	-	-
CosyPose [37]	-	89.8	84.5
RePose [32]	62.1	88.5	82.0
GDR-Net [67]	60.1	91.6	84.4
SO-Pose [19]	56.8	90.9	83.9
Ours	80.5	90.1	85.3

Table 4. **Comparison with State of the Art on YCB-V [70]**. We compare our ZebraPose with state of the art w.r.t ADD(-S), AUC of ADD(-S) and AUC of ADD-S in %. (-) denotes results missing from the original paper.

5. Conclusion

In this work, we proposed a novel coarse to fine surface encoding technique to provide 2D-3D correspondences for 6DoF object pose estimation. We also designed a specific hierarchical training strategy that maximizes the prediction accuracy for our proposed binary vertex code. Solving the object pose using a PnP solver based on our vertex code surpasses the state of the art on different benchmarks, proving our approach’s effectiveness. In the future, we would like to extend our vertex code solution to the problem of category-level object pose [40].

Acknowledgements

This work was partially funded by the *Federal Ministry of Education and Research* of the Federal Republic of Germany (BMBF), under grant agreements 16SV8732 (GreifbAR) and 01IW21001 (DECODE). We are thankful to Rene Schuster, Fangwen Shu, Yaxu Xie and Ghazal Ghazaei for proofreading the paper.

References

- [1] David Arthur and Sergei Vassilvitskii. k-means++: The advantages of careful seeding. Technical report, Stanford, 2006. [5](#)
- [2] Daniel Barath and Jiri Matas. Progressive-x: Efficient, anytime, multi-model fitting algorithm. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3780–3788, 2019. [2](#), [5](#), [6](#), [7](#), [8](#), [12](#)
- [3] Eric Brachmann, Alexander Krull, Frank Michel, Stefan Gumhold, Jamie Shotton, and Carsten Rother. Learning 6d object pose estimation using 3d object coordinates. In *European conference on computer vision*, pages 536–551. Springer, 2014. [2](#)
- [4] Eric Brachmann, Alexander Krull, Sebastian Nowozin, Jamie Shotton, Frank Michel, Stefan Gumhold, and Carsten Rother. Dsac-differentiable ransac for camera localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6684–6692, 2017. [2](#)
- [5] Eric Brachmann, Frank Michel, Alexander Krull, Michael Ying Yang, Stefan Gumhold, and others. Uncertainty-driven 6d pose estimation of objects and scenes from a single rgb image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3364–3372, 2016. [2](#), [5](#), [6](#), [7](#), [8](#), [12](#), [13](#), [14](#)
- [6] Eric Brachmann and Carsten Rother. Learning less is more-6d camera localization via 3d surface regression. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4654–4662, 2018. [2](#)
- [7] Eric Brachmann and Carsten Rother. Neural-guided ransac: Learning where to sample model hypotheses. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4322–4331, 2019. [2](#)
- [8] Benjamin Busam, Tolga Birdal, and Nassir Navab. Camera pose filtering with local regression geodesics on the riemannian manifold of dual quaternions. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 2436–2445, 2017. [2](#)
- [9] Benjamin Busam, Marco Esposito, Simon Che’Rose, Nassir Navab, and Benjamin Frisch. A stereo vision approach for cooperative robotic movement therapy. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 127–135, 2015. [1](#)
- [10] Benjamin Busam, Hyun Jun Jung, and Nassir Navab. I like to move it: 6d pose estimation as an action decision process. *arXiv preprint arXiv:2009.12678*, 2020. [1](#)
- [11] Berk Calli, Arjun Singh, Aaron Walsman, Siddhartha Srinivasa, Pieter Abbeel, and Aaron M Dollar. The ycb object and model set: Towards common benchmarks for manipulation research. In *Advanced Robotics (ICAR), 2015 International Conference on*, pages 510–517. IEEE, 2015. [2](#)
- [12] Carlos Campos, Richard Elvira, Juan J Gómez Rodríguez, José MM Montiel, and Juan D Tardós. Orb-slam3: An accurate open-source library for visual, visual-inertial, and multi-map slam. *IEEE Transactions on Robotics*, 2021. [2](#)
- [13] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017. [5](#)
- [14] Qi Chen and Hartmut Prutzsch. General midpoint subdivision. *arXiv preprint arXiv:1208.3794*, 2012. [5](#)
- [15] Wenzheng Chen, Huan Ling, Jun Gao, Edward Smith, Jaakko Lehtinen, Alec Jacobson, and Sanja Fidler. Learning to predict 3d objects with an interpolation-based differentiable renderer. *Advances in Neural Information Processing Systems*, 32:9609–9619, 2019. [2](#)
- [16] Alvaro Collet, Manuel Martinez, and Siddhartha S Srinivasa. The moped framework: Object recognition and pose estimation for manipulation. *The international journal of robotics research*, 30(10):1284–1306, 2011. [1](#)
- [17] Andrew DeLong, Anton Osokin, Hossam N Isack, and Yuri Boykov. Fast approximate energy minimization with label costs. *International journal of computer vision*, 96(1):1–27, 2012. [7](#)
- [18] Maximilian Denninger, Martin Sundermeyer, Dominik Winkelbauer, Youssef Zidan, Dmitry Olefir, Mohamad Elbadrawy, Ahsan Lodhi, and Harinandan Katam. Blender-proc. *arXiv preprint arXiv:1911.01911*, 2019. [6](#)
- [19] Yan Di, Fabian Manhardt, Gu Wang, Xiangyang Ji, Nassir Navab, and Federico Tombari. So-pose: Exploiting self-occlusion for direct 6d pose estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12396–12405, 2021. [2](#), [3](#), [6](#), [8](#)
- [20] T Do, Trung Pham, Ming Cai, and Ian Reid. Real-time monocular object instance 6d pose estimation. 2019. [2](#)
- [21] Ghazal Ghazaei, Iro Laina, Christian Rupprecht, Federico Tombari, Nassir Navab, and Kianoush Nazarpour. Dealing with ambiguity in robotic grasping via multiple predictions. In *Asian Conference on Computer Vision*, pages 38–55. Springer, 2018. [1](#)
- [22] Silvio Giancola, Matteo Valenti, and Remo Sala. *A survey on 3D cameras: Metrological comparison of time-of-flight, structured-light and active stereoscopy technologies*. Springer, 2018. [3](#)
- [23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [5](#)
- [24] Yisheng He, Haibin Huang, Haoqiang Fan, Qifeng Chen, and Jian Sun. Ffb6d: A full flow bidirectional fusion network for 6d pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3003–3013, 2021. [8](#)
- [25] Stefan Hinterstoisser, Cedric Cagniart, Slobodan Ilic, Peter Sturm, Nassir Navab, Pascal Fua, and Vincent Lepetit. Gradient response maps for real-time detection of textureless objects. *IEEE transactions on pattern analysis and machine intelligence*, 34(5):876–888, 2011. [2](#)
- [26] Stefan Hinterstoisser, Cedric Cagniart, Slobodan Ilic, Peter Sturm, Nassir Navab, Pascal Fua, and Vincent Lepetit. Gradient response maps for real-time detection of textureless objects. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(5):876–888, 2012. [1](#)
- [27] Stefan Hinterstoisser, Stefan Holzer, Cedric Cagniart, Slobodan Ilic, Kurt Konolige, Nassir Navab, and Vincent Lepetit.

- Multimodal templates for real-time detection of texture-less objects in heavily cluttered scenes. In *2011 international conference on computer vision*, pages 858–865. IEEE, 2011. [3](#)
- [28] Stefan Hinterstoisser, Vincent Lepetit, Slobodan Ilic, Stefan Holzer, Gary Bradski, Kurt Konolige, and Nassir 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. [1](#), [6](#)
- [29] Tomas Hodan, Daniel Barath, and Jiri 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. [2](#), [3](#), [6](#)
- [30] Yinlin Hu, Pascal Fua, Wei Wang, and Mathieu Salzmann. Single-stage 6d object pose estimation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2930–2939, 2020. [2](#), [8](#), [15](#)
- [31] Yinlin Hu, Joachim Hugonot, Pascal Fua, and Mathieu Salzmann. Segmentation-driven 6d object pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3385–3394, 2019. [8](#), [15](#)
- [32] Shun Iwase, Xingyu Liu, Rawal Khirodkar, Rio Yokota, and Kris M Kitani. Repose: Fast 6d object pose refinement via deep texture rendering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3303–3312, 2021. [2](#), [8](#), [15](#)
- [33] Hiroharu Kato, Yoshitaka Ushiku, and Tatsuya Harada. Neural 3d mesh renderer. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3907–3916, 2018. [2](#)
- [34] Wadim Kehl, Fabian Manhardt, Federico Tombari, Slobodan Ilic, and Nassir Navab. SSD-6D: Making RGB-based 3D detection and 6D pose estimation great again. In *Proceedings of the International Conference on Computer Vision (ICCV 2017), Venice, Italy*, pages 22–29, 2017. [1](#), [2](#), [6](#)
- [35] Alex Kendall, Matthew Grimes, and Roberto Cipolla. Posenet: A convolutional network for real-time 6-dof camera relocalization. In *Proceedings of the IEEE international conference on computer vision*, pages 2938–2946, 2015. [2](#)
- [36] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. [6](#)
- [37] Yann Labbé, Justin Carpentier, Mathieu Aubry, and Josef Sivic. Cosypose: Consistent multi-view multi-object 6d pose estimation. In *European Conference on Computer Vision*, pages 574–591. Springer, 2020. [2](#), [4](#), [8](#), [16](#)
- [38] Vincent Lepetit, Francesc Moreno-Noguer, and Pascal Fua. Epnp: An accurate o(n) solution to the pnp problem. *International journal of computer vision*, 81(2):155, 2009. [12](#)
- [39] Vincent Lepetit, Francesc Moreno-Noguer, and Pascal Fua. EPnP: An Accurate O(n) Solution to the PnP Problem. *International Journal of Computer Vision*, 81(2):155–166, 2009. [1](#), [7](#), [8](#)
- [40] Xiaolong Li, He Wang, Li Yi, Leonidas J Guibas, A Lynn Abbott, and Shuran Song. Category-level articulated object pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3706–3715, 2020. [8](#)
- [41] Zhigang Li, Gu Wang, and Xiangyang 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. [4](#), [5](#), [6](#), [7](#)
- [42] Matthew M Loper and Michael J Black. Opendr: An approximate differentiable renderer. In *European Conference on Computer Vision*, pages 154–169. Springer, 2014. [2](#)
- [43] David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60(2):91–110, 2004. [2](#)
- [44] Fabian Manhardt, Diego Martin Arroyo, Christian Rupprecht, Benjamin Busam, Tolga Birdal, Nassir Navab, and Federico Tombari. Explaining the ambiguity of object detection and 6d pose from visual data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6841–6850, 2019. [2](#)
- [45] Fabian Manhardt, Wadim Kehl, Nassir Navab, and Federico Tombari. Deep model-based 6d pose refinement in rgb. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 800–815, 2018. [2](#)
- [46] Eric Marchand, Hideaki Uchiyama, and Fabien Spindler. Pose estimation for augmented reality: a hands-on survey. *IEEE transactions on visualization and computer graphics*, 22(12):2633–2651, 2016. [1](#)
- [47] Michihiko MIMOU, Takeo Kanade, and Toshiyuki SAKAI. A method of time-coded parallel planes of light for depth measurement. *IEICE TRANSACTIONS (1976-1990)*, 64(8):521–528, 1981. [3](#)
- [48] Markus Oberweger, Mahdi Rad, and Vincent Lepetit. Making Deep Heatmaps Robust to Partial Occlusions for 3D Object Pose Estimation. pages 2–4. [2](#)
- [49] Kiru Park, Timothy Patten, and Markus 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. [2](#), [3](#), [5](#)
- [50] Sida Peng, Yuan Liu, Qixing Huang, Xiaowei Zhou, and Hujun 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. [2](#), [3](#), [6](#)
- [51] Luis Pérez, Íñigo Rodríguez, Nuria Rodríguez, Rubén Usamentiaga, and Daniel García. Robot guidance using machine vision techniques in industrial environments: A comparative review. *Sensors*, 16(3):335, 2016. [1](#)
- [52] Jeffrey L Posdamer and MD Altschuler. Surface measurement by space-encoded projected beam systems. *Computer graphics and image processing*, 18(1):1–17, 1982. [3](#)
- [53] Mahdi Rad and Vincent Lepetit. BB8: A scalable, accurate, robust to partial occlusion method for predicting the 3D poses of challenging objects without using depth. In *ICCV*, 2017. [2](#), [3](#)
- [54] Jason Rambach, Chengbiao Deng, Alain Pagani, and Didier Stricker. Learning 6dof object poses from synthetic single

- channel images. In *2018 IEEE International Symposium on Mixed and Augmented Reality Adjunct (ISMAR-Adjunct)*, pages 164–169. IEEE, 2018. 1
- [55] Jason Rambach, Alain Pagani, Michael Schneider, Oleksandr Artemenko, and Didier Stricker. 6dof object tracking based on 3d scans for augmented reality remote live support. *Computers*, 7(1):6, 2018. 1
- [56] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015. 6, 7, 8
- [57] Ethan Rublee, Vincent Rabaud, Kurt Konolige, and Gary Bradski. ORB: An efficient alternative to SIFT or SURF. In *2011 International Conference on Computer Vision*, pages 2564–2571. IEEE, 11 2011. 2
- [58] Joaquim Salvi, Jordi Pages, and Joan Batlle. Pattern codification strategies in structured light systems. *Pattern recognition*, 37(4):827–849, 2004. 3
- [59] Juil Sock, Guillermo Garcia-Hernando, Anil Armagan, and Tae-Kyun Kim. Introducing pose consistency and warp-alignment for self-supervised 6d object pose estimation in color images. In *2020 International Conference on 3D Vision (3DV)*, pages 291–300. IEEE, 2020. 2
- [60] Chen Song, Jiaru Song, and Qixing Huang. Hybridpose: 6d object pose estimation under hybrid representations. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 431–440, 2020. 3, 8
- [61] Yongzhi Su, Jason Rambach, Alain Pagani, and Didier Stricker. Synpo-net—accurate and fast cnn-based 6dof object pose estimation using synthetic training. *Sensors*, 21(1):300, 2021. 2
- [62] Martin Sundermeyer, Zoltan-Csaba Marton, Maximilian Durner, Manuel Brucker, and Rudolph Triebel. Implicit 3d orientation learning for 6d object detection from rgb images. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 699–715, 2018. 1, 3, 6
- [63] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9627–9636, 2019. 6, 7, 8, 12, 14, 17
- [64] Marjan Trobina. Error model of a coded-light range sensor. *Technical report*, 1995. 3
- [65] Joel Vidal, Chyi-Yeu Lin, and Robert Martini. 6D pose estimation using an improved method based on point pair features. In *2018 4th International Conference on Control, Automation and Robotics (ICCAR)*, pages 405–409. IEEE, 2018. 1
- [66] Gu Wang, Fabian Manhardt, Jianzhun Shao, Xiangyang Ji, Nassir Navab, and Federico Tombari. Self6d: Self-supervised monocular 6d object pose estimation. In *European Conference on Computer Vision*, pages 108–125. Springer, 2020. 2
- [67] Gu Wang, Fabian Manhardt, Federico Tombari, and Xiangyang 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. 2, 3, 4, 5, 6, 8, 15, 16
- [68] He Wang, Srinath Sridhar, Jingwei Huang, Julien Valentin, Shuran Song, and Leonidas J Guibas. Normalized object coordinate space for category-level 6d object pose and size estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2642–2651, 2019. 2, 3, 5
- [69] Paul Wohlhart and Vincent Lepetit. Learning descriptors for object recognition and 3d pose estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3109–3118, 2015. 3
- [70] Yu Xiang, Tanner Schmidt, Venkatraman Narayanan, and Dieter Fox. Posecnn: A convolutional neural network for 6d object pose estimation in cluttered scenes. *arXiv preprint arXiv:1711.00199*, 2017. 1, 2, 5, 6, 7, 8, 12, 16, 17
- [71] Sergey Zakharov, Wadim Kehl, and Slobodan Ilic. Deceptionnet: Network-driven domain randomization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 532–541, 2019. 6
- [72] Sergey Zakharov, Ivan S. Shugurov, and Slobodan Ilic. Dpod: 6d pose object detector and refiner. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1941–1950, 2019. 2, 3, 5, 6
- [73] Guangyuan Zhou, Huiqun Wang, Jiaxin Chen, and Di Huang. Pr-gen: A deep graph convolutional network with point refinement for 6d pose estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2793–2802, 2021. 8
- [74] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao 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. 2

6. Supplementary Material

6.1. Hyper-parameters in the Pose Solver

For RANSAC/PnP [38], we set the threshold value for reprojection error as 2 pixels, and execute 150 iterations. For Progressive-X [2], we also set the threshold value for the reprojection error as 2 pixels, and execute 400 iterations. The additional parameters for Progressive-X are "neighborhood_ball_radius=20", "spatial_coherence_weight=0.1", "maximum_tanimoto_similarity=0.9".

6.2. BOP Challenge

We submitted the results on 4 datasets of the BOP challenge and will test our method on the rest 3 datasets. The results are online in [BOP Leaderboards](#) with the submission name "zebrapose".

6.3. YCB-V Evaluation per Object

We present a more detailed result on the YCB-V dataset [70] in Tab. 5 and Tab. 6. As the Tab. 5 shows, in the evaluation of the estimate pose w.r.t ADD(-S) metric, we show major improvement over the state of the art.

In Tab. 6, we carefully calculated the AUC with all-points interpolation algorithm with the maximum threshold of 10 cm. If we calculate the AUC with 11-points interpolation, we will reach AUC of ADD-S of 94%, and AUC of ADD(-S) of 89.8%.

6.4. Qualitative Results

6.4.1 Vertex Code Prediction LM-O

We visualized the predicted binary code of the "duck" object in LM-O dataset [5] with a few examples in Fig. 4. Due to the size limits, we only show the predicted binary code till the 11-th bits. We render the object with the predicted pose on top of the original input ROI. To make the predicted pose more visible in the figure, we set the colour of the object model as red just for this figure. So the duck appears with the orange colour (red + yellow) in the last row. We can see that the rendered object overlapped the object in the original image quite well, indicating that our predicted pose is very accurate.

6.4.2 Pose Prediction LM-O

Qualitative Results on LM-O [5] can be found in Fig. 5. We render the objects with estimated pose on top of the original images. The presented confidence scores are from the 2D object detection with FCOS detector [63].

6.4.3 Pose Prediction YCB-V

Qualitative Results on YCB-V [70] are available in Fig. 6. We render the objects with estimated pose on top of the

original images. The presented confidence scores are from the 2D object detection with FCOS detector [63].

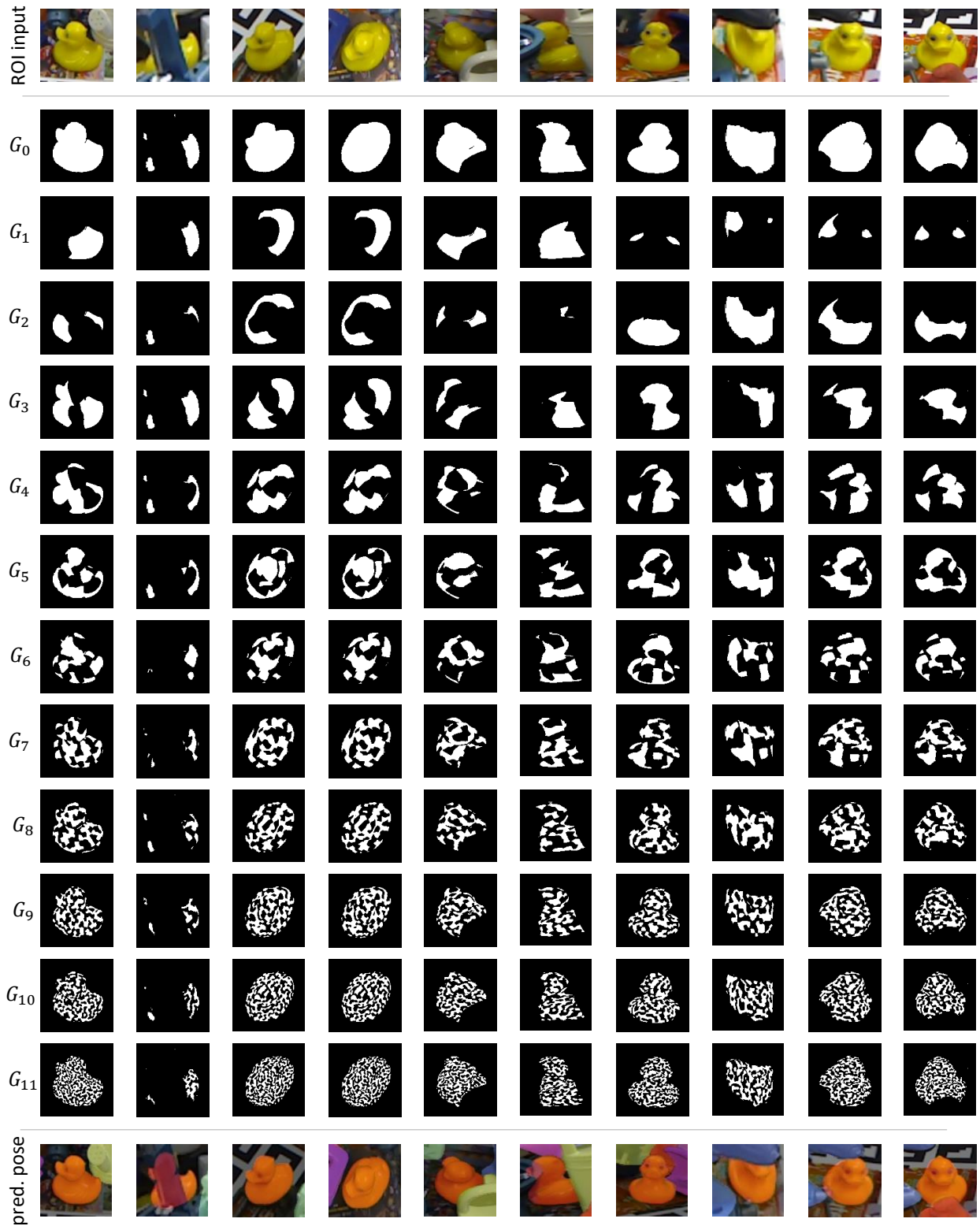


Figure 4. We visualized the predicted binary code of the "duck" in LM-O dataset [5] with a few examples. Due to the size limits, we only show the predicted binary code till the 11-th bit. We set the colour of the object model as red and render the object with the predicted pose on the top of the input ROI. We can see that the rendered object overlaps the object in the image quite well.



Figure 5. **Qualitative Results on LM-O [5]:** We render the objects with estimated pose on top of the original images. The presented confidence score are from the 2D object detection with FCOS detector [63].

Method	SegDriven [31]	Single-Stage [30]	RePose [32]	GDR-Net [67]	Ours
002_master_chef_can	33.0	-	-	41.5	62.6
003_cracker_box	44.6	-	-	83.2	98.5
004_sugar_box	75.6	-	-	91.5	96.3
005_tomato_soup_can	40.8	-	-	65.9	80.5
006_mustard_bottle	70.6	-	-	90.2	100.0
007_tuna_fish_can	18.1	-	-	44.2	70.5
008_pudding_box	12.2	-	-	2.8	99.5
009_gelatin_box	59.4	-	-	61.7	97.2
010_potted_meat_can	33.3	-	-	64.9	76.9
011_banana	16.6	-	-	64.1	71.2
019_pitcher_base	90.0	-	-	99.0	100.0
021_bleach_cleanser	70.9	-	-	73.8	75.9
024_bowl*	30.5	-	-	37.7	18.5
025_mug	40.7	-	-	61.5	77.5
035_power_drill	63.5	-	-	78.5	97.4
036_wood_block*	27.7	-	-	59.5	87.6
037_scissors	17.1	-	-	3.9	71.8
040_large_marker	4.8	-	-	7.4	23.3
051_large_clamp*	25.6	-	-	69.8	87.6
052_extra_large_clamp*	8.8	-	-	90.0	98.0
061_foam_brick*	34.7	-	-	71.9	99.3
mean	39.0	53.9	62.1	60.1	80.5

Table 5. **Comparison with State of the Art on YCB-V.** We report the Average Recall of ADD(-S) in % and compare with state of the art. (*) denotes symmetric objects, (-) denotes the results missing from the original paper.

Method	PoseCNN [70]		CosyPose [37]		GDR-Net [67]		Ours	
Metric	AUC of ADD-S	AUC of ADD(-S)	AUC of ADD-S	AUC of ADD(-S)	AUC of ADD-S	AUC of ADD(-S)	AUC of ADD-S	AUC of ADD(-S)
002_master_chef_can	84.0	50.9	-	-	96.3	65.2	93.7	75.4
003_cracker_box	76.9	51.7	-	-	97.0	88.8	93.0	87.8
004_sugar_box	84.3	68.6	-	-	98.9	95.0	95.1	90.9
005_tomato_soup_can	80.9	66.0	-	-	96.5	91.9	94.4	90.1
006_mustard_bottle	90.2	79.9	-	-	100	92.8	96.0	92.6
007_tuna_fish_can	87.9	70.4	-	-	99.4	94.2	96.9	92.6
008_pudding_box	79.0	62.9	-	-	64.6	44.7	97.2	95.3
009_gelatin_box	87.1	75.2	-	-	97.1	92.5	96.8	94.8
010_potted_meat_can	78.5	59.6	-	-	86.0	80.2	91.7	83.6
011_banana	85.9	72.3	-	-	96.3	85.8	92.6	84.6
019_pitcher_base	76.8	52.5	-	-	99.9	98.5	96.4	93.4
021_bleach_cleanser	71.9	50.5	-	-	94.2	84.3	89.5	80.0
024_bowl*	69.7	69.7	-	-	85.7	85.7	37.1	37.1
025_mug	78.0	57.7	-	-	99.6	94.0	96.1	90.8
035_power_drill	72.8	55.1	-	-	97.5	90.1	95.0	89.7
036_wood_block*	65.8	65.8	-	-	82.5	82.5	84.5	84.5
037_scissors	56.2	35.8	-	-	63.8	49.5	92.5	84.5
040_large_marker	71.4	58.0	-	-	88.0	76.1	80.4	69.5
051_large_clamp*	49.9	49.9	-	-	89.3	89.3	85.6	85.6
052_extra_large_clamp*	47.0	47.0	-	-	93.5	93.5	92.5	92.5
061_foam_brick*	87.8	87.8	-	-	96.9	96.9	95.3	95.3
mean	75.9	61.3	89.8	84.5	91.6	84.3	90.1	85.3

Table 6. **Comparison with State of the Art on YCB-V.** We report the Average Recall w.r.t AUC of ADD(-S) and AUC of ADD-S in % and compare with state of the art. (*) denotes symmetric objects, (-) denotes the results missing from the original paper.

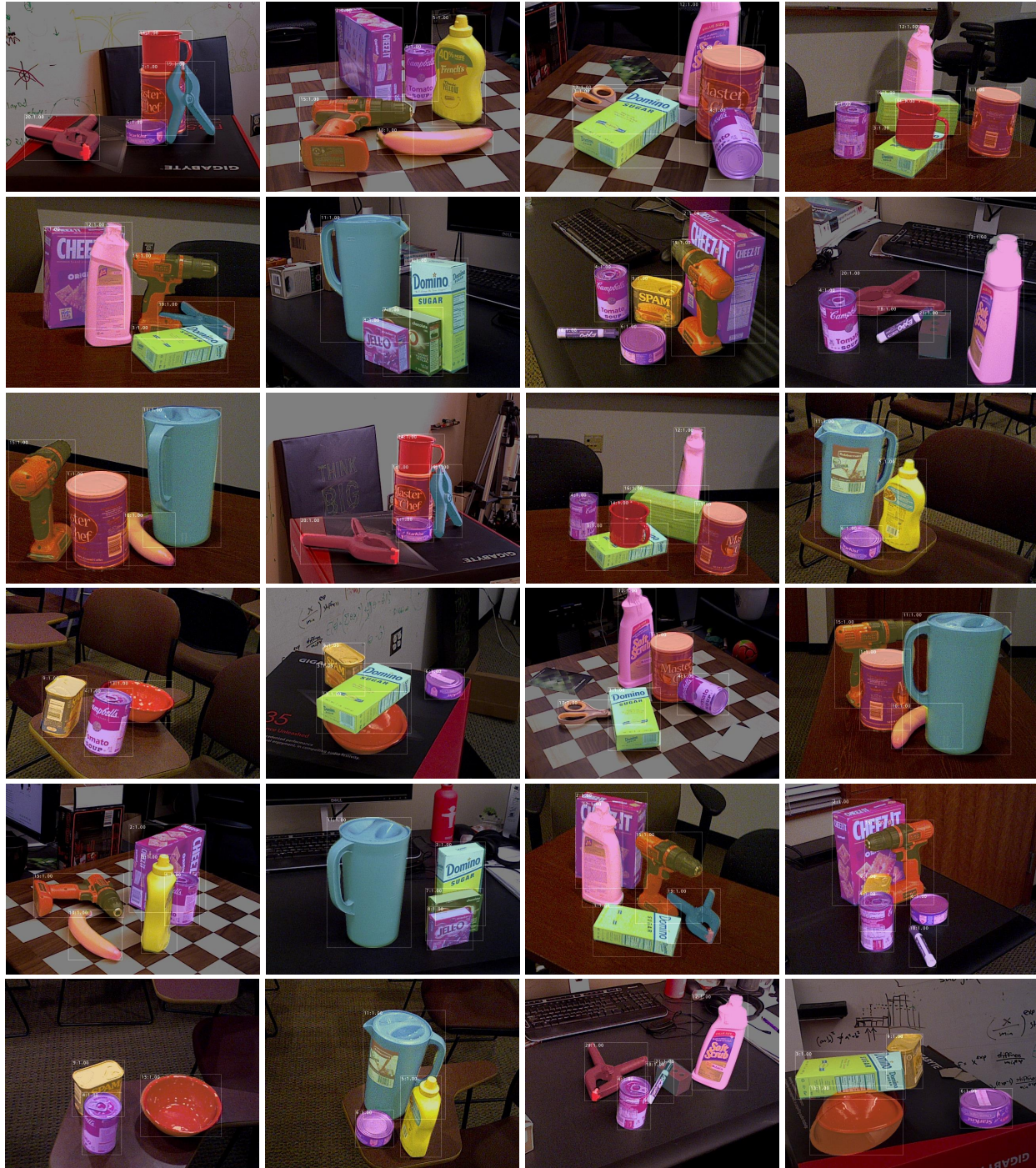


Figure 6. **Qualitative Results on YCB-V [70]:** We render the objects with estimated pose on top of the original images. The presented confidence score are from the 2D object detection with FCOS detector [63].