

Single-Stage Keypoint-Based Category-Level Object Pose Estimation from an RGB Image

Yunzhi Lin^{1,2}, Jonathan Tremblay¹, Stephen Tyree¹, Patricio A. Vela², Stan Birchfield¹

¹NVIDIA: {jtremblay, styree, sbirchfield}@nvidia.com

²Georgia Institute of Technology: {yunzhi.lin, pvela}@gatech.edu

Abstract—Prior work on 6-DoF object pose estimation has largely focused on *instance-level* processing, in which a textured CAD model is available for each object being detected. *Category-level* 6-DoF pose estimation represents an important step toward developing robotic vision systems that operate in unstructured, real-world scenarios. In this work, we propose a single-stage, keypoint-based approach for category-level object pose estimation that operates on unknown object instances within a known category using a single RGB image as input. The proposed network performs 2D object detection, detects 2D keypoints, estimates 6-DoF pose, and regresses relative bounding cuboid dimensions. These quantities are estimated in a sequential fashion, leveraging the recent idea of convGRU for propagating information from easier tasks to those that are more difficult. We favor simplicity in our design choices: generic cuboid vertex coordinates, single-stage network, and monocular RGB input. We conduct extensive experiments on the challenging Objectron benchmark, outperforming state-of-the-art methods on the 3D IoU metric (27.6% higher than the MobilePose single-stage approach and 7.1% higher than the related two-stage approach).

I. INTRODUCTION

Scene awareness is a fundamental skill for robotic manipulators to operate in unconstrained environments. This ability includes locating objects and their poses, also known as the 6-DoF pose estimation problem (*i.e.*, 6 degrees of freedom, from 3D position + orientation). Accurate, real-time pose information of nearby objects in the scene would allow robots to engage in semantic interaction.

The problem of pose estimation is a rich topic in the computer vision community, yet most existing methods have focused on *instance-level* object pose estimation [1], [2], [3]. Such methods suffer from lack of scalability: a detector trained for the ‘cracker box’ in the YCB object dataset [4] may work reliably on instances of that specific object (similar size and texture) but may fail to detect an instance with different textures (*e.g.*, due to seasonal promotional changes) or sizes. Further, the detector is expected to *ignore* all other types of cracker boxes or food-containing cuboids. As a result, the number of *instance-level* detectors required increases rapidly with scene complexity.

To alleviate this challenge, we focus on *category-level* pose estimation. As shown in Figure 1, our goal is to detect and infer the pose and relative size of all objects within a



Fig. 1. Given a single RGB image containing previously unseen instances of known categories (in this case *cereal boxes*, *cups*, and *shoes*), our proposed method detects objects and estimates 6-DoF poses and 3D bounding box dimensions up to a scale factor. We use a separate network for each category.

specific category using a monocular RGB image processed by a single-stage neural network.

A few recent works have considered *category-level* object pose estimation [5], [6], [7], [8], [9], [10], [11]. By removing the requirement of exact CAD models of object instances *at inference time*, these methods promise to scale better for real-world applications. To train the network, one could use a large collection of 3D CAD models (*e.g.*, from ShapeNet [12]) to render synthetic samples with complex annotations, such as pixel-wise segmentation masks or normalized object coordinate spaces (NOCS [6]). Yet the domain gap between synthetic and real data remains an obstacle [13]. In the meantime, many techniques require depth in addition to color (RGB) images [6], [8], [9]. While monocular RGB-based methods [14], [15] have not received much attention, they have great potential for wide applicability and for handling certain material properties, such as transparent or dark surfaces, that are difficult for depth sensors.

In this work, we address the aforementioned challenges and limitations by proposing a simple and efficient RGB-based approach (without depth) that only requires oriented 3D bounding box annotations at training time, and thus does not require CAD models for training. This design decision allows us to take advantage of large collections of real-world images, such as the Objectron dataset [15], which are annotated with *category-level* 3D bounding boxes.

Work was completed while the first author was an intern at NVIDIA.
This work was supported in part by NSF Award #2026611.

Inspired by CenterNet [16], we use a single-stage neural network to regress object locations in the image, 2D keypoint projections of 3D bounding box vertices, and relative cuboid dimensions. The simple design of generic 3D bounding box keypoint regression allows the same method to be applied to a wide variety of categories. To better handle intra-class shape variability, we adopt a two-fold representation of both displacements and heatmaps for keypoint detection. This choice achieves a good balance between accuracy and design complexity, as shown in our experiments. Furthermore, the single-stage design of our network avoids the complexity of multi-stage networks [15] and enables end-to-end learning, and potentially faster training time. To improve the tractability of regressing so many outputs, the network output modalities are grouped by increasing difficulty, and we use a convolutional gated recurrent unit (convGRU) [17], [18] to compute each output group from an underlying sequentially-refined hidden state. This way, the difficulty of predicting the later groups is alleviated by using information stored in the hidden state from previous groups. Once objects have been detected in image space, our approach of estimating relative cuboid dimensions allows us to leverage robust off-the-shelf PnP algorithms for pose estimation.

Our work makes the following contributions:

- A single-stage keypoint-based network for detecting previously unseen objects from known categories and estimating their 6-DoF poses and relative bounding box dimensions from a monocular RGB input.
- Demonstration of the benefit of directly predicting relative dimensions of the 3D bounding cuboid for category-level pose estimation, as well as the benefit of sequential feature association to improve the accuracy of estimating scale information for difficult cases.
- Experiments showing that the proposed method achieves state-of-the-art performance on the large-scale Objectron dataset [15].

II. RELATED WORK

Instance-level object pose estimation. Assuming that a 3D (possibly textured) CAD model is available for each object class at both training and inference time, these methods aim to infer each object's position and orientation in 3D. Current approaches can be divided into two types: template matching and regression. Template matching techniques align known 3D CAD models to the observed 3D point clouds [19], 2D images [20], [21] or local descriptors [22], [23]. State-of-the-art template-based methods have demonstrated impressive results on public benchmarks like BOP [24]. Regression-based methods directly regress the 6-DoF pose [1] or predict the image coordinates of 2D projected keypoints to establish 2D-3D correspondences for solving the 6-DoF pose using a PnP algorithm [25], [26], [27], [2], [14]. Other works have explored different ways to represent objects, including coordinate maps [1], keypoints [28], and symmetry correspondences [29]. Although our method is inspired by keypoint regression techniques, we

do not require 3D CAD models. Thus, the dimensions of the object have to be estimated in addition to pose.

Category-level object pose estimation. Recently, researchers have begun to explore category-level object pose estimation, which does not require instance-specific 3D object models at test time. Wang et al. [6] propose a normalized object coordinate space (NOCS) to serve as a common reference frame for 6-DoF pose and size estimation of unseen objects. Their proposed network is based on two-stage Mask R-CNN [30], which predicts the NOCS map for a pose fitting algorithm that accepts the depth map as input. However, 3D meshes were still required during training to calculate the NOCS map, requiring a synthetic training dataset. Subsequent RGBD works mainly focus on fusing RGB and depth information. Chen et al. [9] propose a correspondence-free approach by learning a canonical shape space for input RGBD images. Their approach also eases network training by matching pose-dependent and pose-independent features separately. Tian et al. [8] model the deformation from the categorical shape prior to the object model by latent embeddings, then recover 6-DoF pose by estimating a similarity transformation between observed points and NOCS map.

To the best of our knowledge, only a few approaches attempt category-level pose estimation from monocular RGB images. Manhardt et al. [11] propose to regress shape and pose parameters and recover depth, while Chen et al. [10] propose a neural analysis-by-synthesis approach. However, both of these methods still require synthetic CAD models (e.g., ShapeNet [12]) at training time. Hou et al. [14] present a single-stage light-weight model with two heads regressing to the centroid location and the 3D bounding box keypoints, respectively, from an RGB image. Similarly, Ahmadyan et al. [15] introduce a two-stage architecture for 3D bounding box keypoint regression from an RGB image. Both approaches are trained directly on real images from Objectron and thus do not require CAD models or synthetic data.¹ These methods do not take the object dimensions into account when solving for pose. They instead directly lift the 2D predicted keypoints to 3D via a modified EPnP algorithm [31] by fixing the homogeneous barycentric coordinates. In contrast, as we show experimentally, our approach achieves better performance by directly regressing the relative dimensions of the cuboid and using an off-the-shelf PnP algorithm.

III. APPROACH

Our approach to category-level object pose estimation is illustrated in Figure 2. We follow the lead of prior correspondence-based methods [25], [26], [27], [2], [14] by predicting 2D image projections of the corners of the 3D bounding cuboid, followed by PnP to compute pose. Inspired by the recent success of works based on the CenterNet architecture [16], [18], [32], [33], we employ a single-stage network to make all predictions, including the relative

¹Also, both methods report real-time performance on a mobile GPU (36 fps for [14] and 83 fps for [15]), using heavily optimized networks. Such optimization is beyond the scope of our work.

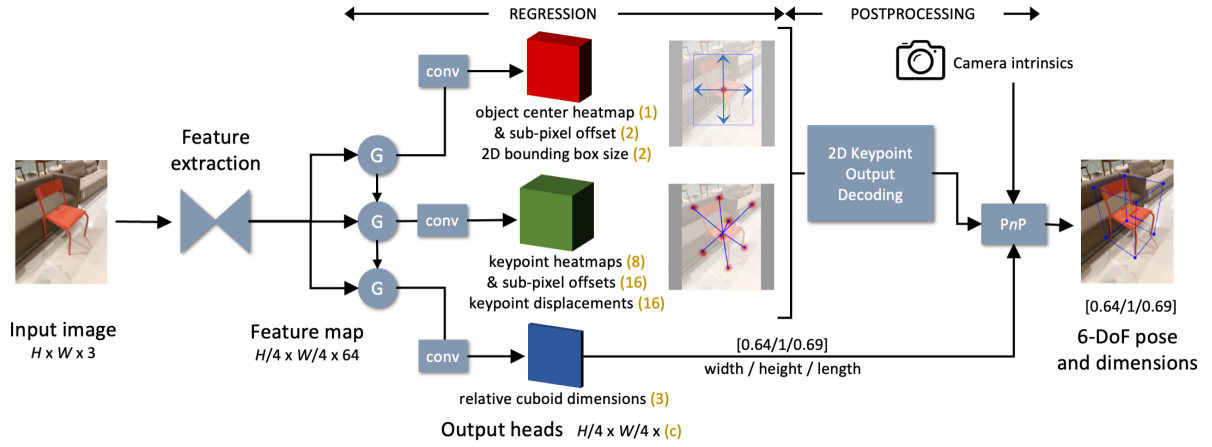


Fig. 2. Overview of our method. For an input image $I \in \mathbb{R}^{H \times W \times 3}$, the backbone network extracts a feature map $\Phi(I) \in \mathbb{R}^{H/4 \times W/4 \times 64}$. A convGRU module associates features into three sequential groups, regressing a total of seven outputs. Each output is in $\mathbb{R}^{H/4 \times W/4 \times c}$, where the number in parentheses (c) represents channels. The corresponding modalities are extracted from the feature map via lookup at the detected center. Finally, we use the decoded 2D keypoints, known camera intrinsics, and the relative cuboid dimensions to obtain a final 6-DoF pose (up to a scale factor) via PnP.

dimensions of the cuboid, which are necessary for PnP. To alleviate the difficulty of inferring 3D structure information from a 2D input, we propose to use a convGRU module [18] to predict outputs grouped in increasing order of difficulty. Further details are provided below.

A. Architecture Design

The network takes an RGB image of resolution $H \times W \times 3$, with source images re-scaled and padded as needed so that $W = H = 512$. We adopt DLA-34 [34] combined with upsampling as the backbone network, where hierarchical aggregation connections are augmented by deformable convolutional layers [35]. The backbone network produces multiple intermediate feature maps of spatial resolutions ranging from $H/4 \times W/4$ to $H/32 \times W/32$, which are aggregated in a single $H/4 \times W/4 \times 64$ output.

The network has a total of seven output heads arranged in three groups. For each output head, a 3×3 convolutional layer with 256 channels followed by a 1×1 convolution layer is used to process the output of the corresponding convGRU module. Outputs are predicted as dense heatmaps or regression maps, but are accessed sparsely in correspondence with detected object centers, as described subsequently.

Object detection branch. The primary output of the entire network, at least conceptually, is the *object center heatmap* whose peaks indicate the centers of the 2D bounding boxes for detected objects.² Other output maps are accessed *w.r.t.* the object center: If a peak is found in the object center heatmap at location (c_x, c_y) , the values at (c_x, c_y) in the remaining outputs are associated with this object. To recover the discretization error resulting from the heatmap output resolution, we regress a local 2D *object center sub-pixel offset* map [16]. Since the Objectron dataset [15] does not provide 2D bounding box annotations, we define it as the smallest axis-aligned rectangle that encloses the extreme points of the projected ground truth 3D bounding box.

²We also tried defining the object center as the projection of the 3D bounding box center, as in [18], [32], but we obtained much better results using the center of the 2D bounding box.

Keypoint detection branch. Our network uses two methods to predict the 2D coordinates of 3D bounding box vertices projected into image space. First, we regress 2D *keypoint displacement* vectors from the bounding box center point. Second, we output a set of 8 *keypoint heatmaps* whose peaks indicate the 2D coordinates of the projected 3D vertices. These peaks are not accessed at the coordinates of the object's center like other outputs. (Further details are given in Section III-C.) Training labels for keypoint heatmaps are generated by a Gaussian kernel centered at the ground truth keypoint coordinates with variance determined by the size of the 2D bounding box. As above, to mitigate discretization error, we also output a local 2D *keypoint sub-pixel offset* for each vertex.

Cuboid dimensions branch. Since category-level pose estimation assumes that we do not have access to the CAD model of the target object instance, we use a final output branch to estimate the *relative dimensions* (width, height, length) of the 3D bounding cuboid. *Relative* values are predicted to avoid the need to implicitly estimate absolute depth from a monocular RGB image, which is a fundamentally ill-posed problem. (For example, we do not know whether we are viewing a full-size actual chair or a toy chair.) Relative values also allow us to apply our network to images obtained with different camera intrinsics without having to retrain the network. Since many target objects in daily life have a canonical orientation when resting on the ground, we choose the up (y) axis as the primary axis. Ground truth scale labels are considered to be $(x/y, 1, z/y)$, with the ratios x/y and z/y estimated by the network. Unlike 3D vehicle detection approaches [18], [32] that use an exponential offset between 3D dimensions and the category-specific dimension template, we directly regress each ratio since the objects we encounter exhibit much more diversity in aspect ratios than are found in vehicles.

B. convGRU Feature Association

We expect that some network outputs are more difficult to learn than the others. Heuristically, we divided them into the

three groups discussed previously, as shown in Figure 2: 1) *object center heatmap, object center sub-pixel offset, and 2D bounding box size*; 2) *x-y displacements to keypoint, keypoint heatmaps, and keypoint sub-pixel offsets*; and 3) *relative cuboid dimensions*. The last group is the most difficult to estimate since 3D structure has to be implicitly deduced from 2D appearance. We hypothesize that keypoints are more easily found once the object centroid and 2D bounding box are estimated, and similarly that bounding box dimensions are more easily predicted after keypoints have been found.

Inspired by Gao et al. [18], this grouping strategy and sequential output construction is naturally formulated by assigning different output groups to different “timesteps” in a recurrent neural network.³ Given an input image I , the i^{th} output ($i = 1, \dots, 7$) is represented as:

$$y_i = \Psi_i(G_t(\Phi(I), h_{t-1})), \quad (1)$$

where $\Phi(I)$ denotes the feature map from the backbone network, $G_t(\cdot)$ represents the GRU at timestep t , $h_{t-1} = G_{t-1}(\Phi(I), h_{t-2})$ denotes the hidden state produced by the GRU cell at the previous timestep, $h_0 = 0$, and Ψ_i is the fully convolutional network for the i^{th} output. The timesteps $t = 1, 2, 3$ correspond to the three groups in our method.

We adopt a single-layer convolutional GRU network where all the convolution layers in the convGRU are set with stride = 1, kernel size = 3, and output channels = 64. The output from a later timestep will have access to the hidden states flowing from the previous timestep, which implements the idea of output grouping and sequential feature association.

C. 2D Keypoint Output Decoding

The outputs of the network are decoded and assembled in the following manner. First, a 3×3 max pooling operation is applied on the heatmap of the 2D object center, which serves as an efficient alternative to non-maximum suppression [16]. For each detected center point, displacement-based keypoint locations are then given by the 2D x - y displacements under the center point. Next, heatmap-based keypoint locations are extracted by finding high confidence peaks in the corresponding heatmaps that are within a margin of the 2D object bounding box. Both estimates of keypoint locations are adjusted according to the sub-pixel offsets, then input, along with the estimated relative cuboid dimensions, to the Levenberg-Marquardt version of PnP [36].

D. Loss Function

We employ penalty-reduced focal losses [37] $\mathcal{L}_{p_{cen}}$ and $\mathcal{L}_{p_{key}}$ in a point-wise manner for the center point and keypoint heatmaps, respectively. The center sub-pixel offset loss, \mathcal{L}_{off} , is computed using an L1 loss. The keypoint sub-pixel offset loss, \mathcal{L}_{offkey} , is computed similarly. The 2D bounding box size, \mathcal{L}_{bbox} , the keypoint displacement loss, \mathcal{L}_{dis} , and the relative cuboid dimensions loss, \mathcal{L}_{dim} , are also computed using an L1 loss w.r.t. their label values. The

³Timesteps here simply refer to recurrent iterations; there is no temporal aspect to the input data.

overall training objective is the weighted combination of the seven loss terms mentioned above.

E. Implementation Details

The network was trained with a batch size of 32 on 4 NVIDIA V-100 GPUs for 140 epochs, starting with pre-trained weights from ImageNet. Data augmentation included random flip, scaling, cropping, and color jittering. We chose Adam as the optimizer with an initial learning rate of $2.5e-4$, dropping 10x at both 90 and 120 epochs. An average of 36 hours was required to train one category (using between 8k to 32k training images depending on the category). Inference speed is around 15 fps on a NVIDIA GTX 1080Ti GPU.

IV. EXPERIMENTAL RESULTS

A. Dataset

The Objectron dataset [15] is a newly proposed benchmark for monocular RGB category-level 6-DoF object pose estimation. It consists of 15k annotated video clips with over 4M annotated frames. Objects are from nine categories: bikes, books, bottles, cameras, cereal boxes, chairs, cups, laptops, and shoes. Each object is annotated with a 3D bounding cuboid, which describes the object’s position and orientation with respect to the camera, as well as the cuboid dimensions. Additional metadata includes camera poses, sparse point clouds, and surface planes, with the latter assuming that the object rests on the ground plane, which yields an absolute scale factor. The cup category contains both cups and mugs, where the former do not have handles. Therefore, we manually differentiate these by training a separate network for each. We also noticed ambiguities in mug instances where the handle is not consistently oriented. To solve this problem, we manually checked all videos and rotated some of the ground truth bounding boxes by 180 degrees to ensure consistent orientation. The cup/mug split is released along with our code.

For training, we extract frames by temporally downsampling the original videos at 15 fps. For testing, we evaluate all the test samples in each category from the official release of the dataset for straightforward comparison with other methods. For symmetric objects like cups, we follow the idea of Wang et al. [6] to generate multiple ground truth labels $\{\mathbf{y}_1, \dots, \mathbf{y}_{|\theta|}\}$ during the training phase, rotating $|\theta|=12$ times around the symmetry axis. The symmetric loss is then computed as $\mathcal{L}_{\text{sym}} = \min_{i=1, \dots, |\theta|} \mathcal{L}(\mathbf{y}_i, \hat{\mathbf{y}})$, where $\hat{\mathbf{y}}$ denotes the prediction, and \mathcal{L} is the asymmetric loss.

B. Metrics

Following the Objectron dataset [15], we evaluate 3D detection and object dimension estimation using the average precision (AP) of 3D IoU metric proposed by [14] with a threshold of 50%. The mean pixel error metric computes the mean normalized distance between the 2D projections of 3D bounding box keypoints given the estimated and ground truth pose. For viewpoint estimation, we report the AP of azimuth and elevation with a threshold of 15° and 10° , respectively. For symmetric object categories (bottle*

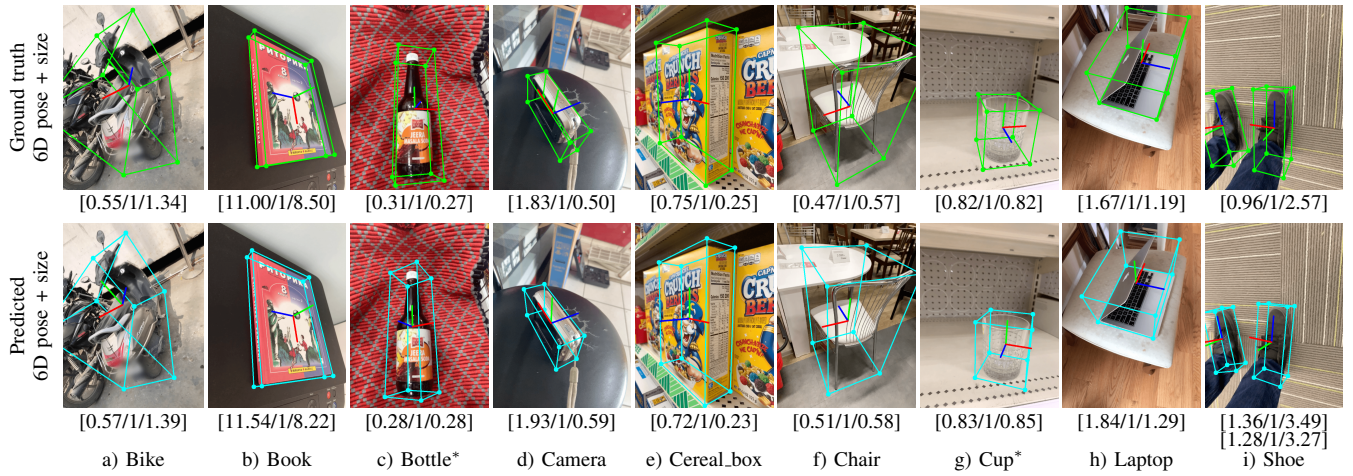


Fig. 3. Sample results of our method on the Objectron dataset [15], where the number below represents the relative dimensions of the 3D bounding box. The proposed method handles large intra-class shape variance, diverse viewpoints, noisy backgrounds, and transparent surfaces.

TABLE I
AVERAGE PRECISION AT 0.5 3D IoU (\uparrow) ON THE OBJECTRON TEST SET [15].

Stages	Method	Bike	Book	Bottle*	Camera	Cereal_box	Chair	Cup*	Laptop	Shoe	Mean
One	MobilePose [14]	0.3109	0.1797	0.5433	0.4483	0.5419	0.6847	0.3665	0.5225	0.4171	0.4461
Two	Two-stage [15]	0.6127	0.5218	0.5744	0.8016	0.6272	0.8505	0.5388	0.6735	0.6606	0.6512
One	Ours	0.6419	0.5565	0.8021	0.7188	0.8211	0.8471	0.7704	0.6766	0.6618	0.7218

TABLE II
POSE ESTIMATION RESULTS ON THE OBJECTRON TEST SET [15],
AVERAGED ACROSS CLASSES.

Method	[14]	[15]	Ours
AP@0.5 3D IoU (\uparrow)	0.4461	0.6512	0.7218
AP@15° azimuth error (\uparrow)	0.4894	0.7844	0.8398
AP@10° elevation error (\uparrow)	0.6414	0.8584	0.8560
Mean pixel error of vertices (\downarrow)	0.1001	0.0467	0.0520

and cup*), we rotate the estimated bounding box along the symmetry axis N times ($N = 100$ following [15]) and evaluate the prediction w.r.t. each rotated instance. The reported number is the instance that maximizes 3D IoU or minimizes pixel projection error, respectively. Although the cup category includes mug instances which are asymmetric, we still treat them as symmetric for comparison with [15]. For the comparison on relative dimension prediction, we use mean relative dimension error, which computes the relative error of the relative dimension across all predictions $\frac{1}{n} \sum_{i=1}^n \frac{|\hat{y}_i - y_i|}{y_i}$, where \hat{y}_i denotes the prediction and y_i denotes the ground truth.

C. Category-Level 6-DoF Pose and Size Estimation

We compare our proposed method with two state-of-the-art methods: single-stage MobilePose [14] and a two-stage network [15]. To the best of our knowledge, these are the only methods available for the Objectron dataset. Results on 3D IoU for each class are shown in Table I. Results averaged across classes (to conserve space) are shown in Table II for pixel projection error, azimuth, and elevation. We present qualitative results in Figure 3. In addition, our supplementary video depicts a robot arranging a pair of shoes by placing the

one shoe next to its match using the position and orientation estimated by our system.

Our method significantly outperforms MobilePose on all metrics, while the two-stage method [15] achieves better performance on the metric of mean pixel error of 2D projection but falls behind on 3D IoU metric. Their two-stage structure allows the keypoint detector to operate at a higher image resolution for better keypoint location performance, but also limits its ability for end-to-end training and fast scale-up to more categories (since the two networks have to be trained independently). Moreover, they do not take the object dimensions into account but rather rely on a modified EPnP algorithm by fixing homogeneous barycentric coordinates across all the cases, which leads to an unstable solution of the 2D–3D correspondence equation [31].

D. Different strategies for 2D Keypoint Output Decoding

Most existing keypoint-based object pose estimation methods adopt either a heatmap [27], [2] or displacement [39], [14] representation for 2D keypoint detection. But large intra-class shape variance poses a key challenge for the keypoint representation. Thus we designed an experiment to compare five different ways for post-processing the 2D keypoint output: 1) *Displacement* ignores the heatmap. 2) *Heatmap* ignores the displacement. 3) *Distance* implements a heuristic similar to [16] that tries to select the more reliable point to use from the displacement or heatmap. 4) *Sampling*, inspired by [38], fits a Gaussian mixture model to the heatmap peak estimate and the displacement prediction for each keypoint and then samples N points ($N = 20$) to obtain a distribution of possible poses. 5) *Our* proposed method keeps both displacement and heatmap. As shown in Table III, our pro-

TABLE III
DIFFERENT STRATEGIES FOR 2D KEYPOINT OUTPUT DECODING (AVERAGE PRECISION AT 0.5 3D IOU METRIC (\uparrow)).

Strategy	w/o add. proc.	Bike	Book	Bottle*	Camera	Cereal_box	Chair	Cup*	Laptop	Shoe	Mean
Displacement	✓	0.6254	0.5263	0.7917	0.7191	0.8115	0.8492	0.7553	0.6737	0.6688	0.7134
Heatmap	✓	0.5788	0.5539	0.7970	0.7035	0.8138	0.8260	0.7626	0.6124	0.6079	0.6951
Distance [16]	✗	0.6305	0.5436	0.7837	0.7111	0.8044	0.8460	0.7640	0.6692	0.6529	0.7117
Sampling [38]	✗	0.6279	0.5516	0.7873	0.7182	0.8134	0.8466	0.7687	0.6751	0.6641	0.7170
Disp. + Heatmap	✓	0.6419	0.5565	0.8021	0.7188	0.8211	0.8471	0.7704	0.6766	0.6618	0.7218

TABLE IV
DIFFERENT STRATEGIES FOR COMPUTING CUBOID DIMENSIONS.

Method	Mean cuboid dimension error (\downarrow)				Average precision at 0.5 3D IoU (\uparrow)			
	Book	Laptop	Others	Mean	Book	Laptop	Others	Mean
Keypoint lifting [14] (no dim. pred.)	-	-	-	-	0.3999	0.5159	0.6540	0.6104
Estimated dim. (w/o convGRU)	0.8474	0.9124	0.2434	0.3849	0.5401	0.6378	0.7528	0.7164
Estimated dim. (w/ convGRU)	0.7440	0.6799	0.2475	0.3507	0.5565	0.6766	0.7519	0.7218
Ground truth dim. (oracle)	0	0	0	0	0.6955	0.6942	0.7907	0.7694

posed combined method (*Displacement + Heatmap*) is better than either of the single representations (*Displacement* or *Heatmap*), and it also does not require additional processing (as in *Distance* [16] or *Sampling* [38]). Thus, to balance accuracy and efficiency, we use this combined representation in all the other experiments.

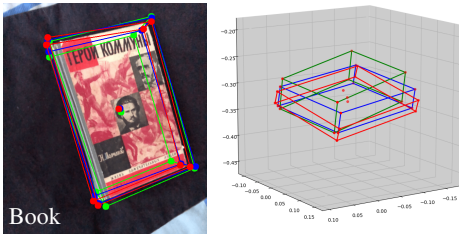


Fig. 4. Improvement due to the convGRU feature association module, where green is ground truth, red is our proposed method without feature association module, and blue is our proposed method with the module. When viewing a thin object from a certain perspective (azimuth close to 90°), it is challenging to estimate the thickness of the target.

E. Different Strategies for Cuboid Dimension Prediction

In this section, we present an experiment on different strategies for cuboid dimension prediction, which further reveals the importance of accurate scale prediction and demonstrates the value of the sequential feature association module (convGRU) for hard cases. We tested the following variants of our system: 1) *Keypoint lifting*, where we reimplemented the postprocessing part proposed by [14] to retrieve the final pose using only the 2D projected cuboid keypoints; 2) *No convGRU*, in which the convGRU layers were removed from our method (See Figure 2); 3) *with convGRU*, our proposed method; 4) *Oracle*, which has access to the ground truth 3D aspect ratio (relative dimensions).

Results are shown in Table IV, where we isolated two specific categories (book and laptop) since they have the greatest difference. The results indicate a strong relationship between 3D IoU result and the corresponding mean cuboid dimension error. For many categories (“Others” in

Table IV), the performance does not differ much, as well as their cuboid dimension prediction. Those instances are of similar aspect ratios and easier to estimate, *e.g.*, bottles. On the other hand, the book and laptop categories are more challenging as the thickness of a book varies greatly while the laptop exhibits different modes (whether the lid is open or closed). The proposed convGRU module improves the prediction of their cuboid dimensions and leads to a better 3D IoU result. Moreover, we found oracle with ground truth dimension achieved the best result while performance degraded using the simplified EPnP variant from [14], which suggests that predicting relative dimensions for category-level pose estimation from monocular RGB input is crucial to solving this problem. Figure 4 uses a particular example to show the ability of convGRU (blue) to retrieve the object 3D aspect ratio (relative dimensions) when compared without convGRU (red). Even though both 2D keypoints look accurate, their scale predictions are different, leading to 3D IoU (\uparrow) improvement (0.5059 with convGRU vs. 0.3204 without convGRU).

V. CONCLUSION

We have presented a single-stage method for category-level 6-DoF pose prediction of previously unseen object instances from RGB input. Unlike previous approaches, CAD models of instances are not needed at training nor test time, and complex annotations are not required for training. We also show the importance of accurate cuboid dimension prediction, and we demonstrate the use of convGRU sequential feature association to further improve accuracy for challenging cases with varied aspect ratios. We demonstrate state-of-the-art performance on the large-scale real-world Objectron dataset, along with a robotic experiment indicating the potential of our proposed method to serve real-world applications. Future work will aim to improve the results by incorporating shape geometry embeddings, exploring lightweight backbone networks, and leveraging iterative post refinement.

REFERENCES

- [1] Y. Xiang, T. Schmidt, V. Narayanan, and D. Fox, "PoseCNN: A convolutional neural network for 6D object pose estimation in cluttered scenes," in *RSS*, 2018.
- [2] J. Tremblay, T. To, B. Sundaralingam, Y. Xiang, D. Fox, and S. Birchfield, "Deep object pose estimation for semantic robotic grasping of household objects," in *CoRL*, 2018, pp. 306–316.
- [3] C. Wang, D. Xu, Y. Zhu, R. Martín-Martín, C. Lu, L. Fei-Fei, and S. Savarese, "DenseFusion: 6D object pose estimation by iterative dense fusion," in *CVPR*, 2019, pp. 3343–3352.
- [4] B. Calli, A. Walsman, A. Singh, S. Srinivasa, P. Abbeel, and A. M. Dollar, "Benchmarking in manipulation research: Using the Yale-CMU-Berkeley object and model set," *IEEE Robotics and Automation Magazine*, vol. 22, no. 3, Sep. 2015.
- [5] G. Pavlakos, X. Zhou, A. Chan, K. G. Derpanis, and K. Daniilidis, "6-DoF object pose from semantic keypoints," in *ICRA*, 2017, pp. 2011–2018.
- [6] H. Wang, S. Sridhar, J. Huang, J. Valentin, S. Song, and L. J. Guibas, "Normalized object coordinate space for category-level 6D object pose and size estimation," in *CVPR*, 2019, pp. 2642–2651.
- [7] L. Ke, S. Li, Y. Sun, Y.-W. Tai, and C.-K. Tang, "GSNet: Joint vehicle pose and shape reconstruction with geometrical and scene-aware supervision," in *ECCV*, 2020, pp. 515–532.
- [8] M. Tian, M. H. Ang, and G. H. Lee, "Shape prior deformation for categorical 6D object pose and size estimation," in *ECCV*, 2020, pp. 530–546.
- [9] D. Chen, J. Li, Z. Wang, and K. Xu, "Learning canonical shape space for category-level 6D object pose and size estimation," in *CVPR*, 2020, pp. 11973–11982.
- [10] X. Chen, Z. Dong, J. Song, A. Geiger, and O. Hilliges, "Category level object pose estimation via neural analysis-by-synthesis," in *ECCV*, 2020, pp. 139–156.
- [11] F. Manhardt, G. Wang, B. Busam, M. Nickel, S. Meier, L. Minciullo, X. Ji, and N. Navab, "CPS++: Improving class-level 6D pose and shape estimation from monocular images with self-supervised learning," *arXiv preprint arXiv:2003.05848*, 2020.
- [12] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, J. Xiao, L. Yi, and F. Yu, "ShapeNet: An Information-Rich 3D Model Repository," *arXiv preprint arXiv:1512.03012*, 2015.
- [13] A. Kortylewski, A. Schneider, T. Gerig, B. Egger, A. Morel-Forster, and T. Vetter, "Training deep face recognition systems with synthetic data," *arXiv preprint arXiv:1802.05891*, 2018.
- [14] T. Hou, A. Ahmadyan, L. Zhang, J. Wei, and M. Grundmann, "MobilePose: Real-time pose estimation for unseen objects with weak shape supervision," *arXiv preprint arXiv:2003.03522*, 2020.
- [15] A. Ahmadyan, L. Zhang, A. Ablavatski, J. Wei, and M. Grundmann, "Objectron: A large scale dataset of object-centric videos in the wild with pose annotations," in *CVPR*, 2021, pp. 7822–7831.
- [16] X. Zhou, D. Wang, and P. Krähenbühl, "Objects as points," *arXiv preprint arXiv:1904.07850*, 2019.
- [17] N. Ballas, L. Yao, C. Pal, and A. C. Courville, "Delving deeper into convolutional networks for learning video representations," in *ICLR*, 2016.
- [18] T. Gao, H. Pan, and H. Gao, "Monocular 3D object detection with sequential feature association and depth hint augmentation," *arXiv preprint arXiv:2011.14589*, 2020.
- [19] A. Zeng, K.-T. Yu, S. Song, D. Suo, E. Walker, A. Rodriguez, and J. Xiao, "Multi-view self-supervised deep learning for 6D pose estimation in the amazon picking challenge," in *ICRA*, 2017, pp. 1386–1383.
- [20] Y. Li, G. Wang, X. Ji, Y. Xiang, and D. Fox, "DeepIM: Deep iterative matching for 6D pose estimation," in *ECCV*, 2018, pp. 683–698.
- [21] M. Sundermeyer, Z.-C. Marton, M. Durner, M. Brucker, and R. Triebel, "Implicit 3D orientation learning for 6D object detection from RGB images," in *ECCV*, 2018, pp. 699–715.
- [22] C. Choi and H. I. Christensen, "3D textureless object detection and tracking: An edge-based approach," in *IROS*, 2012, pp. 3877–3884.
- [23] T. Birdal and S. Ilic, "Point pair features based object detection and pose estimation revisited," in *3DV*, 2015, pp. 527–535.
- [24] T. Hodaň, F. Michel, E. Brachmann, W. Kehl, A. Glent Buch, D. Kraft, B. Drost, J. Vidal, S. Ihrke, X. Zabulis, C. Sahin, F. Manhardt, F. Tombari, T.-K. Kim, J. Matas, and C. Rother, "BOP: Benchmark for 6D object pose estimation," *ECCV*, 2018.
- [25] 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 *ICCV*, 2017, pp. 3828–3836.
- [26] B. Tekin, S. N. Sinha, and P. Fua, "Real-time seamless single shot 6D object pose prediction," in *CVPR*, 2018, pp. 292–301.
- [27] M. Oberweger, M. Rad, and V. Lepetit, "Making deep heatmaps robust to partial occlusions for 3D object pose estimation," in *ECCV*, 2018, pp. 119–134.
- [28] S. Peng, Y. Liu, Q. Huang, X. Zhou, and H. Bao, "PVNet: Pixel-wise voting network for 6DoF pose estimation," in *CVPR*, 2019, pp. 4561–4570.
- [29] C. Song, J. Song, and Q. Huang, "HybridPose: 6D object pose estimation under hybrid representations," in *CVPR*, 2020, pp. 431–440.
- [30] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask R-CNN," in *ICCV*, 2017, pp. 2961–2969.
- [31] V. Lepetit, F. Moreno-Noguer, and P. Fua, "EPnP: An accurate $O(n)$ solution to the PnP problem," *IJCV*, vol. 81, no. 2, p. 155, 2009.
- [32] Z. Liu, Z. Wu, and R. Tóth, "SMOKE: Single-stage monocular 3D object detection via keypoint estimation," in *CVPR Workshops*, 2020, pp. 996–997.
- [33] Y. Wang, Z. Xu, H. Shen, B. Cheng, and L. Yang, "CenterMask: Single shot instance segmentation with point representation," in *CVPR*, 2020, pp. 9313–9321.
- [34] F. Yu, D. Wang, E. Shelhamer, and T. Darrell, "Deep layer aggregation," in *CVPR*, 2018, pp. 2403–2412.
- [35] X. Zhu, H. Hu, S. Lin, and J. Dai, "Deformable Convnets v2: More deformable, better results," in *CVPR*, 2019, pp. 9308–9316.
- [36] Y. I. Abdel-Aziz, H. M. Karara, and M. Hauck, "Direct linear transformation from comparator coordinates into object space coordinates in close-range photogrammetry," *Photogrammetric Engineering & Remote Sensing*, vol. 81, no. 2, pp. 103–107, 2015.
- [37] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *ICCV*, 2017, pp. 2980–2988.
- [38] M. A. Lee, C. Florensa, J. Tremblay, N. Ratliff, A. Garg, F. Ramos, and D. Fox, "Guided Uncertainty-Aware Policy Optimization: Combining learning and model-based strategies for sample-efficient policy learning," in *ICRA*, 2020, pp. 7505–7512.
- [39] Y. Hu, J. Hugonot, P. Fua, and M. Salzmann, "Segmentation-driven 6D object pose estimation," in *CVPR*, 2019, pp. 3385–3394.