

Regular Article

Semantic keypoint-based pose estimation from single RGB frames

Karl Schmeckpeper¹, Philip R. Osteen², Yufu Wang¹, Georgios Pavlakos^{3,*},
Kenneth Chaney¹, Wyatt Jordan², Xiaowei Zhou^{4,*}, Konstantinos G. Derpanis⁵
and Kostas Daniilidis¹

¹University of Pennsylvania, 3330 Walnut Street, Philadelphia, PA, USA 19104

²DEVCOM Army Research Laboratory, 2800 Powder Mill Road, Adelphi, MD, USA 20783

³University of California, Berkeley, 253 Cory Hall, Berkeley, CA 94720

⁴Zhejiang University, 866 Yuhangtang Rd, Hangzhou, China

⁵Ryerson University, 350 Victoria Street, Toronto, ON, Canada

Abstract: This paper presents an approach to estimating the continuous 6-DoF pose of an object from a single RGB image. The approach combines semantic keypoints predicted by a convolutional network (convnet) with a deformable shape model. Unlike prior investigators, we are agnostic to whether the object is textured or textureless, as the convnet learns the optimal representation from the available training-image data. Furthermore, the approach can be applied to instance- and class-based pose recovery. Additionally, we accompany our main pipeline with a technique for semi-automatic data generation from unlabeled videos. This procedure allows us to train the learnable components of our method with minimal manual intervention in the labeling process. Empirically, we show that our approach can accurately recover the 6-DoF object pose for both instance- and class-based scenarios even against a cluttered background. We apply our approach both to several, existing, large-scale datasets - including PASCAL3D+, LineMOD-Occluded, YCB-Video, and TUD-Light - and, using our labeling pipeline, to a new dataset with novel object classes that we introduce here. Extensive empirical evaluations show that our approach is able to provide pose estimation results comparable to the state of the art.

Keywords: computer vision, perception

1. Introduction

This paper addresses the task of estimating an object’s continuous, six degrees-of-freedom (6-DoF) pose (3D translation and rotation) of an object from a single image. Despite its importance in a variety of applications, e.g., robotic manipulation, navigation, etc., and its intense study, most solutions tend to treat objects on a case-by-case basis. For instance, approaches can be distinguished by whether they apply to “sufficiently” textured objects or apply to textureless objects. In addition,

* Work done while at the University of Pennsylvania

Received: 14 September 2020; revised: 25 May 2021; accepted: 02 September 2021; published: 16 March 2022.

Correspondence: Karl Schmeckpeper, University of Pennsylvania, 3330 Walnut Street, Philadelphia, PA, USA 19104, Email: karls@seas.upenn.edu

This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Copyright © 2022 Schmeckpeper, Osteen, Wang, Pavlakos, Chaney, Jordan, Zhou, Derpanis and Daniilidis



Figure 1. Overview of approach. Given a single RGB image of an object (first and third columns), we estimate its 6-DoF pose. The corresponding CAD model is overlaid on the image (second and fourth columns) using the estimated pose. Our method deals with both instance-based (left pair) and class-based scenarios (right pair).

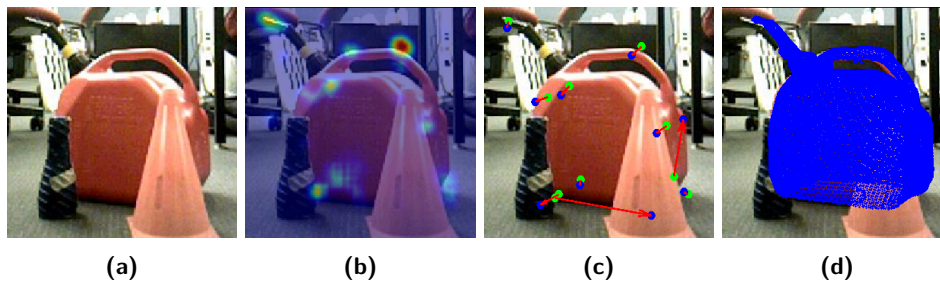


Figure 2. Processing pipeline of our approach. Given a single RGB image of an object (a), we localize a set of class-specific keypoints using a convnet. The output of this step is a set of heatmaps, one for each keypoint, which are combined for visualization in (b), sometimes leading to false detections. In (c), green dots represent the detected keypoints and the corresponding blue dots (connected with an arrow) the groundtruth locations. For robustness against such localization errors, we solve a fitting problem to enforce global consistency of the keypoints, where the response of the heatmaps is used as a measure of certainty for each keypoint. The optimization recovers the full 6-DoF pose of the object (d).

some approaches focus on instance-based object detection, while others address object classes. In this work we strive for an approach where the admissibility of objects considered is as wide as possible (examples in Figure 1).

Our approach combines statistical models of appearance and the 3D shape layout of objects for pose estimation. It consists of two stages that first reason about the 2D projected shape of an object, as captured by a set of 2D semantic keypoints, then estimates the 3D pose consistent with those keypoints (Figure 2). In the first stage, we use a high capacity, convolutional network (convnet) to predict a set of semantic keypoints. Here, the network takes advantage of its ability to aggregate appearance information over a wide receptive field, as compared to localized part models, e.g., (Gu and Ren, 2010), to make reliable predictions of the semantic keypoints. In the second stage, the semantic keypoint predictions are used to reason explicitly about the intra-class shape variability and the camera pose modeled by a weak- or full-perspective camera model. Pose estimates are realized by maximizing the geometric consistency between the parametrized, deformable model and the 2D semantic keypoints. While this work focuses on RGB-based pose estimation, our method can provide a robust way to initialize the iterative closest point (ICP) algorithm (Besl and McKay, 1992) to further refine the pose in the case where a corresponding point cloud is provided with the image.

In this paper, we additionally focus on efficiently collecting the data to train the learnable components of our approach. More specifically, although the convnet used for the detection of the semantic keypoints requires a relatively small number of images for training, annotation time can become excessive as all relevant keypoints must be identified. To bypass the manual labor, we propose a semi-automatic procedure that can greatly accelerate deployment of our algorithm, particularly in the instance-based scenario. Our method relies on 3D reconstruction of the object from an unlabeled

input video, and selection of a small number of keypoints on the final 3D surface. Projecting the 3D keypoints on all the images of the video allows us to generate a large number of training data for the keypoint localization network. This procedure can be repeated for multiple videos to significantly increase the quantity of data available for training, capturing multiple backgrounds, poses, and lighting conditions.

In summary, this paper extends our previous work published in ICRA 2017 (Pavlakos et al., 2017). Our main contribution lies in presenting an accurate and robust approach for 6-DoF object pose estimation, both for the class-based and the instance-based cases. Here we extend our previous work in the following ways.

- We present an approach for semi-automatic data generation which is appropriate to train the learnable components of the original method
- We release a dataset of outdoor pose estimation challenges labeled with this pipeline
- We carefully evaluate the accuracy of the data labeling and the effect it has on the rest of the pipeline
- We provide additional comparisons on new datasets between our pose estimation algorithm and those developed since the release of the original paper, demonstrating that our approach has withstood the test of time.

2. Related Works

Estimating the 6-DoF pose of an object from a single image has attracted significant study. Given a rigid 3D object model and a set of 2D-to-3D point correspondences, various solutions have been explored, e.g., (Fischler and Bolles, 1981; Lepetit et al., 2009). This is commonly referred to as the Perspective-n-Point problem (PnP). To relax the assumption of known 2D landmarks, a number of approaches have considered the detection of discriminative image keypoints (Collet et al., 2009; Collet et al., 2011; Xie et al., 2013), such as SIFT (Lowe, 2004), with highly textured objects. A drawback with these approaches is that they are inadequate for addressing textureless objects and their performance is susceptible to scene clutter. An alternative to sparse discriminative keypoints is offered by dense methods (Brachmann et al., 2014; Brachmann et al., 2016; Dumanoglou et al., 2016; Michel et al., 2017), where every pixel or patch is voting for the object pose. These approaches are also applicable for textureless objects, however, the assumption that a corresponding instance-specific 3D model is available for each object limits their general applicability.

Holistic template-based approaches are one of the earliest approaches considered in the object detection literature. To accommodate appearance variation due to camera capture viewpoint, a set of template images of the object instance are captured about the view sphere and are compared to the input image at runtime. In recent years, template-based methods have received renewed interest due to the advent of accelerated matching schemes and their ability to detect textureless objects by way of focusing their model description on the object shape (Muja et al., 2011; Hinterstoisser et al., 2012; Rios-Cabrera and Tuytelaars, 2013; Xie et al., 2013; Cao et al., 2016). While impressive results in terms of accuracy and speed have been demonstrated, holistic template-based approaches are limited to instance-based object detection and are not robust to occlusions. To address class variability and viewpoint, various approaches have used a collection of 2D appearance-based part templates trained separately on discretized views (Gu and Ren, 2010; Fidler et al., 2012; Pepik et al., 2012; Xiang et al., 2014; Zhu et al., 2014).

Convolutional networks or convnets (LeCun et al., 1989; Krizhevsky et al., 2012) have emerged as the method of choice for a variety of vision problems. Closest to the current work is their application in camera viewpoint and keypoint prediction. Convnets have been used to predict the camera's viewpoint with respect to the object by way of direct regression or casting the problem as classification into a set discrete views (Massa et al., 2014; Tulsiani and Malik, 2015; Su et al., 2015a). While these approaches allow for object category pose estimation they do not provide fine-grained information about the 3D layout of the object. Convnet-based keypoint prediction for human pose

estimation (e.g., (Toshev and Szegedy, 2014; Zhou et al., 2015b; Newell et al., 2016; Wei et al., 2016)) has attracted considerable study, while limited attention has been given to their application with generic object categories (Long et al., 2014; Tulsiani and Malik, 2015). Their success is due in part to the high discriminative capacity of the network. Furthermore, their ability to aggregate information over a wide field of view allows for the resolution of ambiguities (e.g., symmetry) and for localizing occluding joints.

Statistical shape-based models tackle recognition by aligning a shape subspace model to image features. While originally proposed in the context of 2D shape (Cootes et al., 1995) they have proven useful for modelling the 3D shape of a host of object classes, e.g., faces (Cao et al., 2013), cars (Zia et al., 2013; Murthy et al., 2017), and human pose (Ramakrishna et al., 2012). In (Zhu et al., 2015), data-driven discriminative landmark hypotheses were combined with a 3D deformable shape model and a weak perspective camera model in a convex optimization framework to globally recover the shape and pose of an object in a single image. Here, we adapt this approach and extend it with a perspective camera model, in cases where the camera intrinsics are known.

Since our initial work on pose estimation from semantic keypoints (Pavlakos et al., 2017), several other works have been inspired by or have used our method. Wang et al. demonstrate the effectiveness of using learned keypoints for object tracking (Wang et al., 2019). Qin et al. show a pipeline that learns task specific keypoints for manipulation (Qin et al., 2019). Manuelli et al. use semantic keypoints as their object representation for manipulation (Manuelli et al., 2019). Peng et al. propose a refinement of the keypoint detection for pose estimation using pixel-wise voting (Peng et al., 2019). Hu et al. use object segmentations rather than object bounding boxes to input detected objects into their semantic keypoint based pose estimation pipeline (Hu et al., 2019). Zhou et al. learn a set of category agnostic keypoints for object pose estimation (Zhou et al., 2018).

Other works directly use our work or similar approaches. Zuo et al. show that a keypoint based pose estimation pipeline can be used to track the position of a low cost robotic arm that does not have its own sensors, enabling cheap manipulation (Zuo et al., 2019). The keypoint-based pose estimation pipeline was used by (Bowman et al., 2017) as the first step in their semantic SLAM pipeline. Our pipeline has been used by (Vasilopoulos and Koditschek, 2018; Vasilopoulos et al., 2020b; Vasilopoulos et al., 2020a) to detect and localize objects for reactive planning for navigation, and was used by the Robotics Collaborative Technology Alliance (RCTA) program to rapidly collect and annotate data for pose estimation for mobile manipulation (Narayanan et al., 2020; Kessens et al., 2020).

RGB-based instance-level pose estimation has also made improvements since our method was initially proposed, largely due to deep network based pipelines (Kehl et al., 2017; Xiang et al., 2017; Mousavian et al., 2017; Kehl et al., 2017; Tekin et al., 2018; Tremblay et al., 2018; Do et al., 2018; Li et al., 2018; Montserrat et al., 2019; Labbé et al., 2020; Ke et al., 2020; Hou et al., 2020) and higher quality training data (Zeng et al., 2017; Hodaň et al., 2018; Hodaň et al., 2020; Sundermeyer et al., 2020; Wang et al., 2020). Interestingly, while several single-shot methods have focused on speed (Kehl et al., 2017; Tekin et al., 2018; Hou et al., 2020), many state of the art approaches still follow the two-stage approach, where 2D-3D correspondences are detected with a deep network in the first stage, and 6D pose is solved in the second stage with PnP methods. Different intermediate representations have been explored to establish correspondences, such as 3D bounding box corners (Rad and Lepetit, 2017), dense object coordinates (Park et al., 2019), dense fragments (Hodaň et al., 2020), pixel-wise voting (Peng et al., 2019) and their hybrids (Song et al., 2020). We will show through new experiments in Section 6.2 that our simple semantic keypoint based representation is still effective and competitive, even in challenging occluded cases.

3. Pose estimation from Semantic Keypoints

Our pipeline includes object detection, keypoint localization, and pose optimization. As object detection has been a well studied problem, we assume that a bounding box around the object has

been provided by an off-the-shelf object detector, e.g., Faster R-CNN (Ren et al., 2015), and focus on the later two processes.

3.1. Pose optimization

Given the keypoint locations on the 3D model as well as their correspondences in the 2D image, one naive approach is to simply apply an existing PnP algorithm to solve for the 6-DoF pose. This approach is problematic because the keypoint predictions can be rendered imprecise due to occlusions and false detections in the background. Moreover, the exact 3D model of the object instance in the testing image is often unavailable. To address these difficulties, we fit a deformable shape model to the 2D detections while considering the uncertainty in keypoint predictions. This approach additionally allows the optimization to include the uncertainty of the detections, making it better able to handle false or imprecise detections.

We build a deformable shape model for each object category using 3D CAD models with annotated keypoints. More specifically, the p keypoint locations on a 3D object model are denoted by $\mathbf{S} \in \mathbb{R}^{3 \times p}$ and

$$\mathbf{S} = \mathbf{B}_0 + \sum_{i=1}^k c_i \mathbf{B}_i, \quad (1)$$

where \mathbf{B}_0 is the mean shape of the given 3D model and $\mathbf{B}_1, \dots, \mathbf{B}_k$ are several modes of possible shape variability computed by Principal Component Analysis (PCA).

Given detected keypoints in an image, which are denoted by $\mathbf{W} \in \mathbb{R}^{2 \times p}$, the goal is to estimate the rotation $\mathbf{R} \in \mathbb{R}^{3 \times 3}$ and translation $\mathbf{T} \in \mathbb{R}^{3 \times 1}$ between the object and camera frames as well as the coefficients of the shape deformation $\mathbf{c} = [c_1, \dots, c_k]^\top$.

The inference is formulated as the following optimization problem:

$$\min_{\theta} \frac{1}{2} \left\| \xi(\theta) \mathbf{D}^{\frac{1}{2}} \right\|_F^2 + \frac{\lambda}{2} \|\mathbf{c}\|_2^2, \quad (2)$$

where θ is the set of unknowns, $\xi(\theta)$ denotes the fitting residuals dependent on θ , and the Tikhonov regularizer $\|\mathbf{c}\|_2^2$ is introduced to penalize large deviations from the mean shape.

To incorporate the uncertainty in 2D keypoint predictions, a diagonal weighting matrix $\mathbf{D} \in \mathbb{R}^{p \times p}$ is introduced:

$$\mathbf{D} = \begin{bmatrix} d_1 & 0 & \cdots & 0 \\ 0 & d_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & d_p \end{bmatrix}, \quad (3)$$

where d_i indicates the localization confidence of the i th keypoint in the image. In our implementation, d_i is assigned the peak value in the heatmap corresponding to the i th keypoint. As shown previously (Newell et al., 2016), the peak intensity of the heatmap provides a good indicator for the visibility of a keypoint in the image.

The fitting residuals, $\xi(\theta)$, measure the differences between the given 2D keypoints, provided by the previous processing stage, and the projections of 3D keypoints. Two camera models are next considered.

3.1.1. Weak perspective model

If the camera intrinsic parameters are unknown, the weak perspective camera model is adopted, which is usually a good approximation to the full perspective case when the camera is relatively far

away from the object. In this case, the reprojection error is written as

$$\xi(\theta) = \mathbf{W} - s\bar{\mathbf{R}} \left(\mathbf{B}_0 + \sum_{i=1}^k c_i \mathbf{B}_i \right) - \bar{\mathbf{T}} \mathbf{1}^\top, \quad (4)$$

where s is a scalar, $\bar{\mathbf{R}} \in \mathbb{R}^{2 \times 3}$ and $\bar{\mathbf{T}} \in \mathbb{R}^2$ denote the first two rows of \mathbf{R} and \mathbf{T} , respectively, and $\theta = \{s, \mathbf{c}, \bar{\mathbf{R}}, \bar{\mathbf{T}}\}$.

The problem in (2) is continuous and in principal can be locally solved by any gradient-based method. We solve it with a block coordinate descent scheme because of its fast convergence and the simplicity in implementation. We alternately update each of the variables while fixing the others. The updates of s , \mathbf{c} and $\bar{\mathbf{T}}$ are simply solved using closed-form least squares solutions. The update of $\bar{\mathbf{R}}$ should consider the $SO(3)$ constraint. Here, the Manopt toolbox (Boumal et al., 2014) is used to optimize $\bar{\mathbf{R}}$ over the Stiefel manifold. As the problem in (2) is non-convex, we further adopt a convex relaxation approach (Zhou et al., 2015a) to initialize the optimization. More specifically, we only estimate the pose parameters while fixing the 3D model as the mean shape in the initialization stage. By setting $\mathbf{c} = \mathbf{0}$ and replacing the orthogonality constraint on $\bar{\mathbf{R}}$ by the spectral norm regularizer, the problem in (2) can be converted to a convex program and solved with global optimality (Zhou et al., 2015a).

3.1.2. Full perspective model

If the camera intrinsic parameters are known, the full perspective camera model is used, and the residuals are defined as

$$\xi(\theta) = \tilde{\mathbf{W}} \mathbf{Z} - \mathbf{R} \left(\mathbf{B}_0 + \sum_{i=1}^k c_i \mathbf{B}_i \right) - \mathbf{T} \mathbf{1}^\top, \quad (5)$$

where $\tilde{\mathbf{W}} \in \mathbb{R}^{3 \times p}$ represents the normalized homogeneous coordinates of the 2D keypoints and \mathbf{Z} is a diagonal matrix:

$$\mathbf{Z} = \begin{bmatrix} z_1 & 0 & \cdots & 0 \\ 0 & z_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & z_p \end{bmatrix}, \quad (6)$$

where z_i is the depth for the i th keypoint in 3D. Intuitively, the distances from the 3D points to the rays crossing the corresponding 2D points are minimized. In this case, the unknown parameter set θ is given by $\{\mathbf{Z}, \mathbf{c}, \mathbf{R}, \mathbf{T}\}$.

The optimization here is similar to the alternating scheme in the weak perspective case. The update of \mathbf{Z} also admits a closed-form solution and the update of \mathbf{R} can be analytically solved by the orthogonal Procrustes analysis. To avoid local minima, the optimization is initialized by the weak perspective solution.

3.2. Keypoint localization

The keypoint localization step employs the “stacked hourglass” network architecture (Newell et al., 2016) that has been shown to be particularly effective for 2D human pose estimation. Motivated by this success, we use the same network design and train the network for object keypoint localization.

Network architecture. A high level overview of the main network components is presented in Figure 3. The network takes as input an RGB image, and outputs a set of heatmaps, one per keypoint, with the intensity of the heatmap indicating the confidence of the respective keypoint to be located at this position. The network consists of two hourglass components, where each component can be further subdivided into two main processing stages. In the first stage, a series

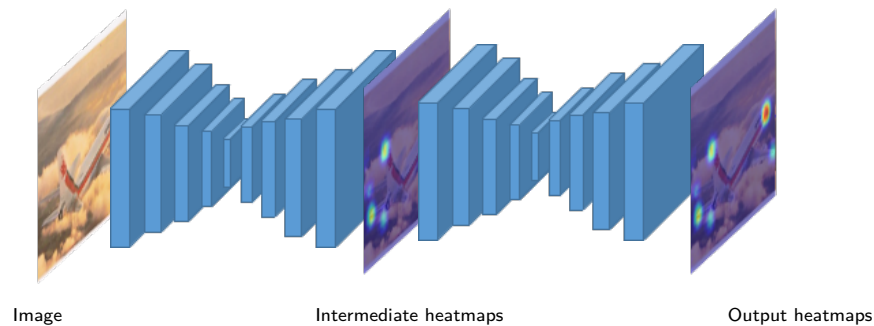


Figure 3. Overview of the stacked hourglass architecture (Newell et al., 2016). Here, two hourglass modules are stacked together. The symmetric nature of the design allows for bottom-up processing (from high to low resolution) in the first half of the module, and top-down processing (from low to high resolution) in the second half. Intermediate supervision is applied after the first module. The heatmap responses of the second module represent the final output of the network that is used for keypoint localization.

of convolutional and max-pooling layers are applied to the input. After each max-pooling layer, the resolution of the feature maps decreases by a factor of two, allowing the next convolutional layer to process the features at a coarser scale. This sequence of processing continues until reaching the lowest resolution (4×4 feature maps), which is illustrated by the smallest layer in the middle of each module in Figure 3. Following these downsampling layers, the processing continues with a series of convolutional and nearest-neighbor upsampling layers. Each upsampling layer increases the resolution by a factor of two. This process culminates with a set of heatmaps at the same resolution as the input of the hourglass module. A second hourglass component is stacked at the end of the first one to refine the output heatmaps. The groundtruth labels used to supervise the training are synthesized heatmaps based on a 2D Gaussian centered at each keypoint with a standard deviation set to one. The ℓ_2 loss is minimized during training. Additionally, intermediate supervision is applied at the end of the first module, which provides a richer gradient signal to the network and guides the learning procedure towards a better optimum (Lee et al., 2015). We evenly weight the the loss from the intermediate heatmaps with the heatmaps from the last module. We train the model using the RMSProp optimizer (Tieleman and Hinton, 2012). The heatmap responses of the last module are considered as the final output of the network and the peak in each heatmap indicates the most likely location for the corresponding keypoint. For more details of our architecture, please see our official implementation.¹

Design benefits. The most critical design element of the hourglass network is the symmetric combination of bottom-up and top-down processing that each hourglass module performs. Given the large appearance changes of objects due to in-class and viewpoint variation, both local and global cues are needed to effectively decide the locations of the keypoints in the image. The consolidation of features across different scales in the hourglass architecture allows the network to successfully integrate both local and global appearance information, and commit to a keypoint location only after this information has been made available to the network. Moreover, the stacking of hourglass modules provides a form of iterative processing that has been shown to be effective with several other recent network designs (Carreira et al., 2016; Wei et al., 2016) and offers additional refinement of the network estimates. Additionally, the application of intermediate supervision at the end of each module has been validated as an effective training strategy, particularly ameliorating the practical

¹ Our implementation of the keypoint detector network is available at <https://github.com/kschmeckpeper/keypoint-detection>

issue of vanishing gradients when training a deep neural network (Lee et al., 2015). Finally, residual layers are introduced (He et al., 2016), which have achieved state-of-the-art results for many visual tasks, including object classification (He et al., 2016), instance segmentation (Dai et al., 2016), and 2D human pose estimation (Newell et al., 2016).

4. Data Collection and Annotation

With most modern machine learning frameworks, the data quantity and quality is as important, if not more so, than details of a particular learning algorithm. While unlabeled data can be rapidly obtained, high quality labeled data is time consuming and costly to acquire. Recent efforts (Castrejon et al., 2017; Sohn et al., 2020; Xie et al., 2019) leverage machine learning to assist with manual annotation or to provide semi-supervised labels from unlabeled data. Similarly, (Hua et al., 2016) and (Marion et al., 2018) leverage multi-view scene reconstructions to efficiently annotate 3D objects, then reproject the annotation masks to individual frames based on estimated camera viewpoints. (Osteen et al., 2019) combined geometric reconstructions with hierarchical metric segmentations as well as learned object-segmentation proposals to annotate 3D scene objects, reprojecting label masks to individual frames. Here, we extend this work to add rapid keypoint annotation with depth-based keypoint refinement to enable downstream applications such as object pose estimation from color images. The work is most closely related to (Marion et al., 2018), which also uses object reconstructions to annotate objects and their poses. While their work requires 3D meshes of the objects as a prerequisite to using the annotation tool, we assume no existing models and emphasize the scenario of needing to rapidly annotate a completely unknown object. In addition, our tool specifically produces annotations for training keypoint detectors, while the alternative tool does not. As a result, we can leverage the structure of the object given by the 3D keypoint annotations to perform automatic keypoint refinement with no additional human work beyond keypoint labeling.

To label a new object class with keypoint annotations, we collect sequences of streaming data from different viewpoints, environments, and lighting conditions using an RGB-D sensor. We annotate the resulting images by performing 3D reconstruction, manually annotating the resulting models, projecting the keypoints into the image frame, and refining the resulting keypoints.

4.1. 3D Reconstruction

We leverage state-of-the-art scene reconstruction from RGB-D sequences (Whelan et al., 2016) as well as surface generation algorithms (Kazhdan and Hoppe, 2013) to pre-process the collected data. In previous work (Osteen et al., 2019), we evaluated the performance of modern scene reconstruction algorithms, specifically Elastic Fusion (Whelan et al., 2016), Elastic Reconstruction (Choi et al., 2015), and Kintinuous (Whelan et al., 2015). Kintinuous is an extension to the original KinectFusion reconstruction algorithm (Izadi et al., 2011) which supports operation over larger scales with loop closure. Kintinuous was designed for corridor-like motion with infrequent loop closures over large distances, rather than the potentially frequent loop-closures that can arise from the context of overlapping camera motion for object reconstruction (Whelan et al., 2016). Elastic Reconstruction is an offline reconstruction algorithm that constructs sets of scene fragments (or sub-models), each consisting of a small number (50) of integrated sequential RGB-D frames, and then globally registers the fragments and optimizes the associated sensor poses for a final reconstruction.

In contrast to these algorithms, Elastic Fusion does not perform pose-graph optimization. Instead, Elastic Fusion is designed with a focus on the quality of the reconstructed model, formulating the problem using a deformation graph. A surfel map is used to model the environment, where surfels store properties such as color, normal, radius, etc. Frame-to-model sensor tracking is performed using the portion of the model that has been recently observed (*active*), while global loop-closure candidates are determined by attempting to register new data with portions of the model that have not been recently observed (*inactive*). During loop closure, the model is non-rigidly deformed to

align surfels in the deformation graph, in contrast to pose-graph optimization approaches which do not perform non-rigid deformation of the data, instead optimizing the sensor trajectory given the measured data and associated uncertainty. After deforming sets of matching surfels according to surface correspondences, the camera poses are updated by applying the relative transform which brings the surfaces into alignment.

The evaluation from (Osteen et al., 2019) on sequences from five classes of the Redwood dataset (Choi et al., 2016) found that the reconstructions of Elastic Fusion were comparable in terms of accuracy to the offline Elastic Reconstruction approach at a fraction of the runtime. Therefore, for this work we use Elastic Fusion to create scene models, though the annotation system is agnostic to the choice of reconstruction algorithm.

Once reconstruction has taken place, the resulting 3D model as well as the estimated camera poses from each frame are input to the annotation tool. The only phase of the process that requires human input is the placement of keypoints on the surface of the 3D reconstruction, after which keypoint annotations are automatically projected back to the original input image stream.

4.2. Annotation Refinement

An ideal reconstruction system would correctly estimate the camera poses for each frame of a sequence, such that the generated depth images match corresponding real depth images for all frames. In practice, pose-graph based SLAM techniques such as Elastic Reconstruction as well as deformation techniques such as Elastic Fusion attempt to minimize accumulated pose or model errors across the sequence. While this is optimal for the purposes of estimating sensor trajectory over a sequence, the global pose estimate for a given frame depends the accuracy of the estimates of other poses in the sequence.

While we only annotate sequences for which qualitatively good reconstructions are produced, we observe the presence of drift in the camera pose estimates over the course of a sequence. In many cases, this drift is mitigated by loop closure in the reconstruction algorithm, but we still develop techniques to account for potentially inaccurate camera pose estimates at any given frame.

Specifically, we use the 3D model and camera trajectory to generate depth images at each camera pose, then perform data association between the generated depth images and the corresponding true depth images from the sensor. This technique ensures that the reconstruction results, integrated over the course of an entire sequence, are updated to align with the real data at each individual frame. While we could also compare input color images with their generated counterparts, we choose to focus on depth comparisons on the assumption that 3D reconstructions preserve object shape better than object appearance for any given viewpoint and lighting condition. Therefore we use the depth images to refine annotations, and we compare two alternative approaches to refining annotations: those that update keypoints individually versus those that update all keypoints together.

4.2.1. Keypoint-level Refinement

The most noticeable artifact of keypoint projection with camera trajectory drift is pixels which project to the background rather than the object of interest. More subtly, visible keypoints may project to the correct object but be a few pixels away from the correct location, which may also be obvious to the observer depending on the choice of keypoint location.

To address the issues of keypoints projecting to background objects ("jump-edge" pixels), as well as keypoints that are visible but not correctly located, we create an approach to refine individual keypoint locations. In order to determine if a keypoint is eligible for refinement, we first determine whether it is occluded or not for the current camera pose \mathbf{T}_C^F in fixed frame F . Given N keypoints specified for an object, the set of object keypoints is defined as $K = \{^F\mathbf{k}_0, ^F\mathbf{k}_1, \dots, ^F\mathbf{k}_N\}$.

For each $^C\mathbf{k}_i \in \mathbb{R}^3$ in K , transformed from fixed frame F to camera frame C by $^C\mathbf{k}_i = \mathbf{T}_C^F \mathbf{k}_i$, the occlusion test is performed by comparing the depth of the predicted 3D keypoint location $^C\mathbf{k}_i$ with the depth value of the generated depth image at the pixel which the keypoint projects to, effectively performing a z-buffer test to identify which keypoints are occluded. Also, with knowledge of the

3D object geometry, local normals for each keypoint ${}^C\mathbf{n}_i$ are compared to the camera viewpoint ${}^C\mathbf{v}$, where ${}^C\mathbf{v} \in \mathbb{R}^3$ represents a unit vector along the principal axis of the camera, and keypoints are identified as occluded if their local normals are nearly orthogonal to or face away from the camera. More specifically, a keypoint is determined to be occluded if ${}^C\mathbf{v} \cdot {}^C\mathbf{n}_i > \tau$, where τ is an approximate orthogonality threshold that is set to -0.15 radians for this evaluation.

The following keypoint-level refinement methods are mutually exclusive, such that no two refinements will be applied to the same keypoint.

Jump-edge adjustment. Often, corners and edges of objects are good candidates for keypoint locations. For even small camera pose estimation errors, the projected keypoint location can jump past the object of interest to the background. This is addressed by an adjustment approach that takes advantage of the fact that our 3D keypoint positions are known for each estimated camera pose, and that we can easily compare the depth values of the generated depth image and real depth image at the keypoint’s projected pixel location to identify jump-edges. Specifically, we use our reconstructed model to project the keypoints into the estimated camera frame. Any keypoint where the projected depth is less than the measured depth at its projected pixel coordinates is deemed a jump edge. If a jump edge is detected, the refined keypoint location is updated to be the projection of the closest point in the real depth image to the estimated 3D keypoint location.

Feature matching. For keypoints that correctly project to the desired object, 3D feature matching is employed to improve the keypoint location estimate. If a keypoint is clearly visible for the given camera pose, then a 3D descriptor is extracted from the generated depth image and compared to descriptors extracted in the true depth image around the estimated keypoint location. For this study, we use the FPFH feature descriptor (Rusu et al., 2009), which was shown by (Choi et al., 2015) to be highly effective at geometric pairwise registration, but have also tested with the SHOT descriptor (Salti et al., 2014). The updated keypoint is defined by a weighted sum of the Euclidean distance and the matched feature vector distances (with weight values of 0.7 and 0.3, respectively) between the real and generated depth images. The closest match in the real depth image is then chosen as the refined keypoint location.

4.2.2. Object-level Refinement

Since the goal of annotation refinement is to maximize the alignment of the annotated model and an individual frame, regardless of any other camera pose estimates in the trajectory, we perform iterative closest point (ICP) on the generated and real depth images to ensure the model is best aligned with each depth frame. For this work, our primary concern is to align points belonging to the object of interest, so we limit our ICP refinement to the bounding volumes that contain the annotated objects. We use the volume that minimally bounds the annotated keypoints to define the volume in the generated image, and since there may be error in the pose estimate, the volume is dilated for the real depth image. Then, point-to-plane ICP (Chen and Medioni, 1992) is performed on the cropped points from the generated and real point clouds derived from the respective depth images. The camera pose estimate is adjusted based on the ICP alignment, and all keypoints are transformed accordingly. This means that all keypoints undergo the same rigid transformation, which preserves the relative poses of the keypoints, and allows us to refine the positions of both visible and occluded keypoints.

5. RCTA Object Keypoints Dataset

While large scale datasets such as PASCAL3D+ contain many types of objects, there are still many objects and environments for which there is no corresponding dataset. This is especially true for the Robotics Collaborative Technology Alliance (RCTA) program, which is a consortium of government, industry, and academic researchers focusing on fielding embodied autonomous systems in challenging

Table 1. Number of keypoints for each class in the RCTA Object Keypoints Dataset

Class	Barrel	Barrier	Crate	Gas Can	Hedgehog	Robot
Number of Keypoints	6	9	12	10	6	11

Table 2. Number of keypoints for each class in PASCAL3D+

Class	aero	bike	bottle	bus	car	chair	sofa	train	TV monitor	boat
Number of Keypoints	8	11	7	12	12	10	10	8	4	6

and unstructured environments. The scope of the work includes many unique objects that are not well-represented in existing datasets, such as the Czech Hedgehog.

As part of this work, we have therefore created a dataset containing 101 sequences split across six of the objects used by the RCTA (Narayanan et al., 2020) in both indoor and outdoor environments: *barrel*, *barrier*, *crate*, *gas can*, *Czech Hedgehog*, and *robot*. The number of keypoints for each object class are shown in Table 1.

We collected many data sequences for each object, not knowing ahead of time how many would yield accurate 3D reconstructions. Indeed, many reconstructions did fail or were not accurate enough to include in the analysis. To avoid adding the burden to a user of filtering poor reconstructions manually, we defer the determination of reconstruction viability to the actual annotation process and allow the annotator to trivially skip a poor reconstruction. In total, we collected 283 sequences across the six object classes, and determined that 101 were suitable for annotation, some of which are shown in Figure 4. For each object class, a random subset of sequences and frames were annotated by hand. The total number of frames analyzed for this work is 10281; of those, 951 were manually annotated.

With the public release of this dataset, we hope to increase diversity in the types of environments and objects used for keypoint annotation. Furthermore, we release the raw data collection sequences as well as the reconstructed models, with the hope of testing new reconstruction algorithms as they are released.

The full dataset is available for download at <https://sites.google.com/view/rcta-object-keypoints-dataset>.

6. Results

We present an analysis of our architecture on existing benchmark datasets as well as custom data collections. In both qualitative and quantitative output, we show the effectiveness of the approach, and demonstrate the challenges of acquiring accurate manual annotations for occluded keypoints.

6.1. Class-based pose recovery: PASCAL3D+

We demonstrate the full strength of our approach using the large-scale PASCAL3D+ dataset (Xiang et al., 2014). The stacked hourglass network was trained from scratch with the training set of PASCAL3D+. Instead of training separate models for different object classes, a single network was trained to output heatmap predictions for all of the 88 keypoints from all classes. The number of keypoints for each class are shown in Table 2. Using a single network for all keypoints allows us to share features across the available classes and significantly decreases the number of parameters needed for the network. At test time, given the class of the test object, the heatmaps corresponding to the keypoints belonging to this class were extracted. For pose optimization, two cases were tested: (i) the CAD model for the test image was known; and (ii) the CAD model was unknown and the pose was estimated with a deformable model whose basis was learned by PCA on all CAD models for each class in the dataset. Two principal components were used ($k = 2$) for each class, which was sufficient to explain greater than 95% of the shape variation. The 3D model was fit to the 2D keypoints with a weak-perspective model, as the camera intrinsic parameters were not available.



Figure 4. Example 3D reconstructions of the objects in the RCTA dataset.

Semantic correspondences. A crucial component of our approach is the powerful learning procedure that is particularly successful at establishing correspondences between the semantically related keypoints across different instances of an object class. To demonstrate this network property, in Figure 5 we present a subset of the keypoints for each class, along with the localizations of these keypoints, in a randomly-selected set of images among those with the top 50 responses. It is interesting to note that, despite the large appearance differences due to extreme viewpoint and intra-class variability, the predictions are very consistent and preserve the semantic relation across various class instances.

Pose estimation. The quantitative evaluation for pose estimation on PASCAL3D+ is presented in Table 3. We only report rotation errors, as the 3D translation cannot be determined in the weak perspective case nor is the ground truth available. Following work from (Tulsiani and Malik, 2015), we use the geodesic distance in Equation 7 to calculate the rotational error between a pose estimate,

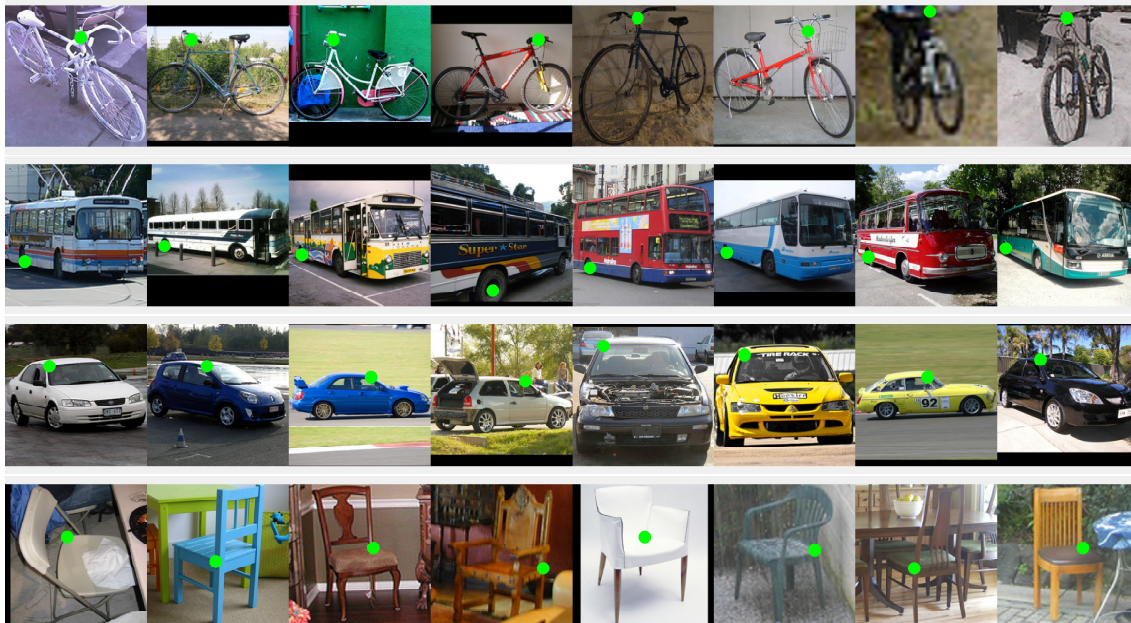


Figure 5. Localization results for diverse keypoint categories. We visualize eight images per category in each row selected randomly from the top 50 responses for each keypoint. The keypoint localization network is particularly successful at establishing semantic correspondences across the instances of a class, despite the significant intra-class variation and wide ranging camera viewpoints.

Table 3. Viewpoint estimation median error (degrees) on PASCAL3D+.

Approach	aero	bike	bottle	bus	car	chair	sofa	train	TV monitor	boat
Tulsiani and Malik (Tulsiani and Malik, 2015)	13.8	17.7	12.9	5.8	9.1	14.8	15.2	8.7	15.4	21.3
Deep3DBox (Mousavian et al., 2017)	13.6	12.5	8.3	3.1	5.8	11.9	12.8	6.3	11.9	22.8
RenderForCNN (Su et al., 2015b)	15.4	14.8	9.3	3.6	6.0	9.7	9.5	6.1	12.6	25.6
ours - PCA basis	11.2	15.2	13.1	4.7	6.9	12.7	21.7	9.1	38.5	37.9
ours - CAD basis	8.0	13.4	11.7	2.0	5.5	10.4	9.6	8.3	32.9	40.7
ours - uniform weights	16.3	17.8	14.1	11.7	30.7	17.6	32.4	20.8	25.0	72.0

\mathbf{R}_1 , and the groundtruth, \mathbf{R}_2 .

$$\Delta(\mathbf{R}_1, \mathbf{R}_2) = \frac{\|\log(\mathbf{R}_1^T \mathbf{R}_2)\|_F}{\sqrt{2}}. \quad (7)$$

where $\|x\|_F$ denotes the Frobenius norm of x . Our method shows improvement across several categories with respect to the state-of-the-art. The best results are achieved in the case where the instance subclass for the object is known and there exists an accurate CAD model correspondence. Our method with uniform weights for all keypoints is also compared as a baseline, which is worse than considering the confidences during model fitting. The worst results occur on object classes that were labeled with few semantic keypoints per instance. The numbers of keypoints for each class in the dataset are shown in Table 2. The two classes which performed worst for our method, the TV Monitor and the boat, had four and six keypoints respectively, while the remaining classes averaged just under ten keypoints each. A subset of results of our method are visualized in Figure 6. While our approach applies to both instance- and category-level pose estimation, it yields comparable results to other leading, but more limited methods.

Category-level pose is useful for robotic applications. In (Kessens et al., 2020), the authors demonstrated a system for mobile manipulation, wherein our keypoint-based pose estimation algorithm

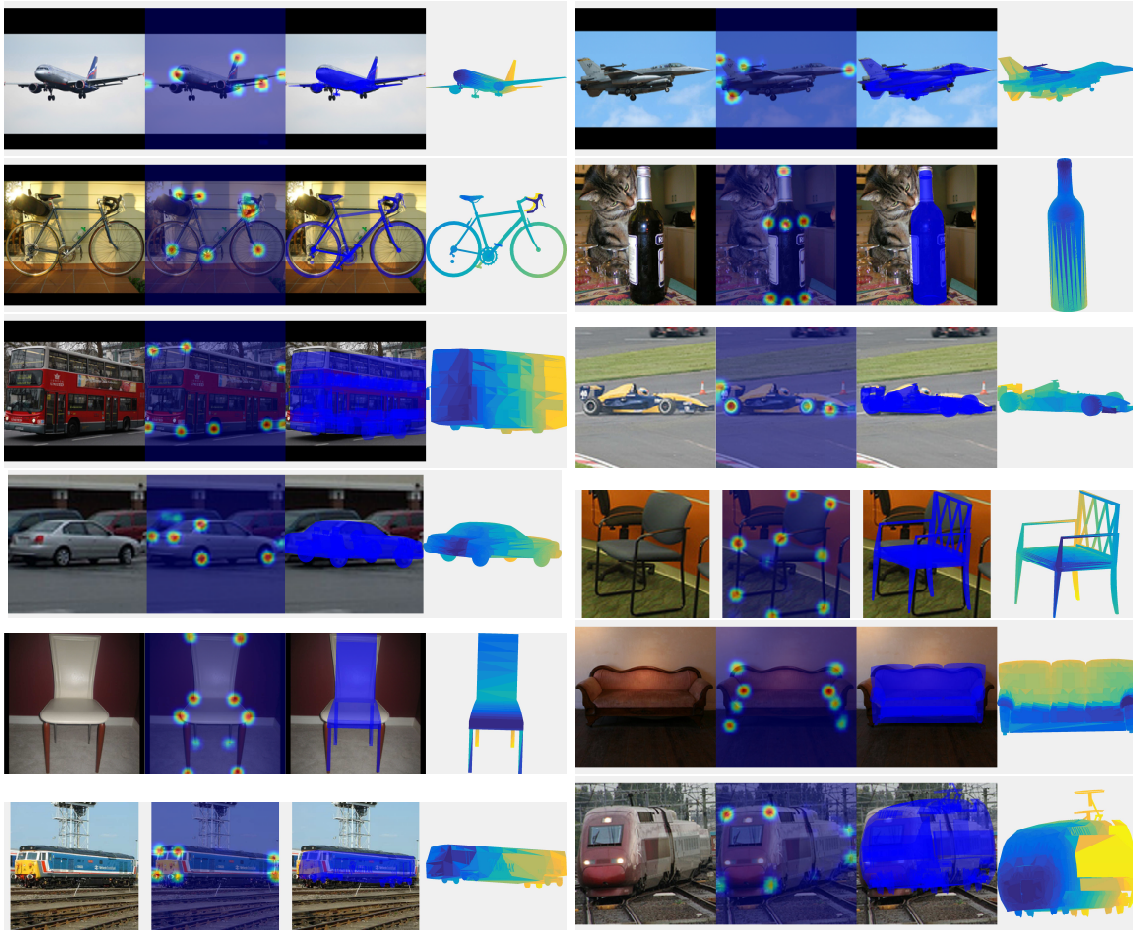


Figure 6. Example results of our approach on PASCAL3D+. For each example from left-to-right: the RGB image of the object, heatmap responses for the keypoints of the specific class, the CAD model projected to 2D after pose estimation, and the CAD model visualized in 3D.

was used to position the mobile manipulator relative to the object. First, the mobile manipulator approached the target object and oriented itself to grasp an object of that object class. Then a second algorithm used a short-range depth sensor to select grasp locations from the generated pointcloud. This approach is applicable, despite variation within an object class, as long as the variation between grasp locations in a class is smaller than the size of the robot’s reachable workspace.

Category-level pose estimation has also been used for semantic SLAM (Bowman et al., 2017), in which application the ability to localize a wide variety of objects within a class is critical to populating the generated map with sufficient density. In the context of semantic SLAM, being able to locate the surface of an object precisely is less important than being able to estimate consistent and accurate poses for every encountered object of the target classes. Thus keypoint-based methods that provide category level pose can perform well while methods that densely reconstruct the surface of an object are not required.

6.2. Instance-based pose recovery: Benchmark for 6-DoF Object Pose Estimation

Our method is equally applicable in estimating object-level pose when the 3D object model and the camera intrinsic are known. This is the classic 6-DoF pose task, and fits well with many robotics

applications, where recognizing 6-DoF pose allows better interactions with a known object. We demonstrate the effectiveness of our method on this task by comparing its performance to state-of-the-art methods on three object pose datasets.

To apply our method, we first set $\mathbf{c}_i = \mathbf{0}$ for all i in Equation 2. We derive \mathbf{B}_0 by choosing recognizable landmarks on the 3D model. For symmetric objects as defined by the BOP benchmark, we place landmarks on geometries or surface texture that break the symmetries. As a result, our pose solution is more accurate than the minimum required. No additional step is taken to handle symmetries. For each training image having a ground truth, 6-DoF pose, we reproject these 3D landmarks to generate 2D keypoints, from which the keypoint detector can be trained. The optimization procedure is the same as in category-level estimation but with \mathbf{c} fixed at zero. Additionally, we fine-tune a Faster-RCNN (Ren et al., 2015) object detector for each dataset to provide cropping bounding boxes for our method following standard practice (Hodaň et al., 2020).

The three object-pose datasets we use are LineMOD-Occluded (Brachmann et al., 2014), YCB-Video (Calli et al., 2015) and TUD-Light (Hodaň et al., 2018). More specifically, we use the train-test split provided by the new BOP benchmark on these three datasets (Hodaň et al., 2020). Using a standard split allows us to compare our method with other state-of-the-art methods while minimizing the performance difference that stems from using different data-generation techniques. The training set comprises synthetic data from physically-based rendering, and real data is added when it is available for that dataset.

Evaluation metrics. On each dataset we evaluate our method with the three metrics used in the BOP19/20 challenge.

- Visible Surface Discrepancy (VSD) calculates the average 3D distance between the estimated and the ground-truth objects on the visible part of the object
- Maximum Symmetry-Aware Surface Distance (MSSD) computes the maximum 3D distance between the estimated and the ground-truth objects. MSSD is similar to the classical 6-DoF pose metric used to compute average distance, but Hodan et al. (Hodaň et al., 2018) argue that MSSD is less sensitive to object geometry
- Maximum Symmetry-Aware Projection Distance (MSPD) computes the maximum 2D distance between projections of the estimated and the ground-truth objects, and is more useful for application in virtual reality, where reprojection accuracy is more important

Additionally, for symmetric objects, we computed MSSD and MSPD over all global symmetric transforms of the object and retaining the lowest error. For each of the three metrics, average recall (AR) is calculated over different thresholds, above which an estimation is considered correct.

We show extensive comparisons in Tables 4, 5, 6. On the overall benchmark, CosyPose (Labbé et al., 2020) achieves state-of-the-art results and outperforms the other methods in most of the metrics. Our method, despite its broader scope, is competitive with the other methods in the 2020 BOP benchmark. In fact, for the YCB-V dataset, our method outperforms all other methods except CosyPose.

Table 4. Quantitative evaluation on LineMOD-Occluded (Brachmann et al., 2014). Comparison to RGB-based 6-DoF pose estimation methods. Statistics are taken from the BOP benchmark (Hodaň et al., 2020).

Method	test input	Avg.	AR_{mspd}	AR_{mssd}	AR_{vsd}
SSD-6D w/o ref (Kehl et al., 2017)	rgb	0.139	0.285	0.083	0.047
Pix2Pose (Park et al., 2019)	rgb	0.363	0.550	0.307	0.233
Leaping from 2D to 6D (Liu et al., 2020)	rgb	0.525	0.781	0.444	0.350
CDPN (Li et al., 2019)	rgb	0.624	0.815	0.612	0.445
CosyPose (Labbé et al., 2020)	rgb	0.633	0.812	0.606	0.480
Ours	rgb	0.612	0.795	0.581	0.459

Table 5. Quantitative evaluation on YCB-Video (Calli et al., 2015). Comparison to RGB-based 6-DoF pose estimation methods. Statistics are taken from the BOP benchmark (Hodaň et al., 2020).

Method	test input	Avg.	AR _{mspd}	AR _{mssd}	AR _{vsd}
Pix2Pose (Park et al., 2019)	rgb	0.457	0.571	0.429	0.372
CDPN (Li et al., 2019)	rgb	0.532	0.631	0.570	0.396
Leaping from 2D to 6D (Liu et al., 2020)	rgb	0.543	0.687	0.499	0.443
CosyPose (Labbé et al., 2020)	rgb	0.821	0.850	0.842	0.772
Ours	rgb	0.618	0.700	0.605	0.549

Table 6. Quantitative evaluation on TUD-Light (Hodaň et al., 2018). Comparison to RGB-based 6-DoF pose estimation methods. Statistics are taken from the BOP benchmark (Hodaň et al., 2020).

Method	test input	Avg.	AR _{mspd}	AR _{mssd}	AR _{vsd}
Pix2Pose (Park et al., 2019)	rgb	0.420	0.641	0.364	0.255
Leaping from 2D to 6D (Liu et al., 2020)	rgb	0.751	0.922	0.716	0.614
CDPN (Li et al., 2019)	rgb	0.772	0.925	0.793	0.597
CosyPose (Labbé et al., 2020)	rgb	0.823	0.973	0.807	0.689
Ours	rgb	0.784	0.950	0.777	0.622

Notably, our method has the second-highest VSD score across the three datasets. VSD calculates discrepancy on only the visible part of the objects. Because our optimization procedure weights the predicted keypoints by their confidence, and occluded keypoints often have lower confidence, our method favors poses that better match visible object-parts.

6.3. Semi-Automated Data Labeling

To apply keypoint-based pose estimation to new object classes, we must be able to rapidly collect images and annotate them with keypoints. We demonstrate a significant reduction to the required labeling effort with our projection based labeling tool.

To show the utility of our labeling tool, we answer several questions:

1. How well do the keypoints generated by the labeling tool match the keypoints labeled by human annotators?
2. Do keypoint localization models trained with the annotations from the tool achieve good performance?
3. Does the tool reduce the effort to label a set of images?
4. How does the keypoint-based refinement method compare with the model-based refinement method?

To answer these questions, a subset of the collected data was labeled by human annotators using a traditional annotation tool, whereby the user is prompted to select all object keypoints in randomly selected images. Histograms of the pixel-wise distances between manual and projected keypoint annotations for each class are given in Figure 10, showing approximately zero-mean distributions with different deviations for different objects. These differences are a function of both the accuracy of the reconstruction and camera pose estimation, as well as the accuracy of the manual annotations, which are not always reliable as ground truth surrogates, particularly in the case of occluded keypoints as illustrated in Figure 7.

Table 7 compares the performance of the keypoint-level refinement and the object-level refinement against manually-annotated keypoints. For the visible case, where manual annotation is a better ground truth proxy, the object-level refinement outperforms the keypoint-level refinement. Therefore, we use the object-level refinement as the only refinement step for subsequent analysis, so following

Table 7. Evaluation of our refinement techniques, keypoint-level (feature matching with jump edge refinement) and object-level refinement (ICP). Error is measured as the ℓ_2 distances in pixels for the refinement methods with respect to manual annotations. The object-level context used by ICP to refine all keypoints performs better than the per-keypoint adjustment of feature matching with jump edge refinement, so we adopt the object-level method as our default refinement technique.

	All		Visible		Occluded	
	Keypoint	Object	Keypoint	Object	Keypoint	Object
Barrel	11.26 ± 0.3	11.05 ± 0.3	11.04 ± 0.3	10.77 ± 0.3	12.78 ± 1.0	12.96 ± 1.0
Barrier	11.90 ± 0.2	11.92 ± 0.2	9.33 ± 0.2	9.47 ± 0.2	14.14 ± 0.4	14.07 ± 0.4
Crate	11.16 ± 0.2	11.14 ± 0.2	9.27 ± 0.2	9.07 ± 0.1	14.66 ± 0.4	14.98 ± 0.4
Gas Can	12.58 ± 0.2	12.19 ± 0.2	11.06 ± 0.3	10.42 ± 0.3	14.73 ± 0.4	14.67 ± 0.4
Hedgehog	11.64 ± 0.2	11.75 ± 0.2	11.64 ± 0.2	11.71 ± 0.2	11.79 ± 1.1	13.20 ± 1.2
Robot	10.43 ± 0.2	10.02 ± 0.2	8.48 ± 0.2	7.96 ± 0.2	14.67 ± 0.6	14.52 ± 0.6

Table 8. Visible versus occluded disparities for manual annotations. ℓ_2 distances in pixels for both the semi-automated keypoint projection method as well as object-level refinement are given with respect to manual annotations. We show the difference for all the keypoints, for only the keypoints that are visible to the camera, and for only the keypoints that are occluded because they are on the far side of the object. There is a clear difference between the offsets for visible versus occluded keypoints, demonstrating the difficulty of accurately annotating occluded keypoints by hand.

	All		Visible		Occluded	
	Projected	Refined	Projected	Refined	Projected	Refined
Barrel	11.29 ± 0.3	11.05 ± 0.3	11.07 ± 0.3	10.77 ± 0.3	12.78 ± 1.0	12.96 ± 1.0
Barrier	12.07 ± 0.2	11.92 ± 0.2	9.71 ± 0.2	9.47 ± 0.2	14.14 ± 0.4	14.07 ± 0.4
Crate	10.87 ± 0.2	11.14 ± 0.2	8.82 ± 0.1	9.07 ± 0.1	14.66 ± 0.4	14.98 ± 0.4
Gas Can	12.44 ± 0.2	12.19 ± 0.2	10.82 ± 0.3	10.42 ± 0.3	14.73 ± 0.4	14.67 ± 0.4
Hedgehog	11.60 ± 0.2	11.75 ± 0.2	11.59 ± 0.2	11.71 ± 0.2	11.79 ± 1.1	13.20 ± 1.2
Robot	9.90 ± 0.2	10.02 ± 0.2	7.72 ± 0.2	7.96 ± 0.2	14.67 ± 0.6	14.52 ± 0.6

mentions of the process refer only to object-level refinement. Table 8 compares the performance of naive projection to keypoint refinement, again using the projected pixel distance to manual annotations as an error metric. The results indicate that keypoint refinement may not be necessary for sequences with very good reconstructions, although we have identified some sequences that clearly benefit from refinement, such as shown in Figure 8. Table 7 and Table 8 also demonstrate the impact of visible versus occluded keypoints on manual annotation accuracy, with occluded keypoints always having higher distances between projected and manual annotations. While manual annotations for clearly visible keypoints, user error notwithstanding, can be considered as approximate ground truth, manual keypoints for occluded keypoints cannot.

For asymmetric objects, we ensured that each data collection contained a full loop around the object of interest, made up of a continuous sequence of images capturing the object from all angles; however for symmetric objects (hedgehog and barrel), pose estimation is ambiguous about the axis of symmetry. Therefore, for such classes, each data sequence contains data for only one side of the object. As a consequence, these sequences have very few occluded keypoints, which results in similar distance results for both visible and occluded keypoints. Details on the distribution of occluded versus visible keypoints in the full dataset are shown in Figure 9a, and the distribution of data used to compare manual and semi-automated annotations are shown in Figure 9b.

Next, we analyzed the performance of the keypoint localization model trained with keypoints from different sources. We trained three keypoint localization models for each class: (i) with the manual annotations, (ii) with the semi-automated annotations, and (iii) with the refined annotations. We then evaluate the performance of these models on a held-out set of images with manually labeled keypoints, as shown in Table 9. Keypoint localization models trained with our refined keypoints

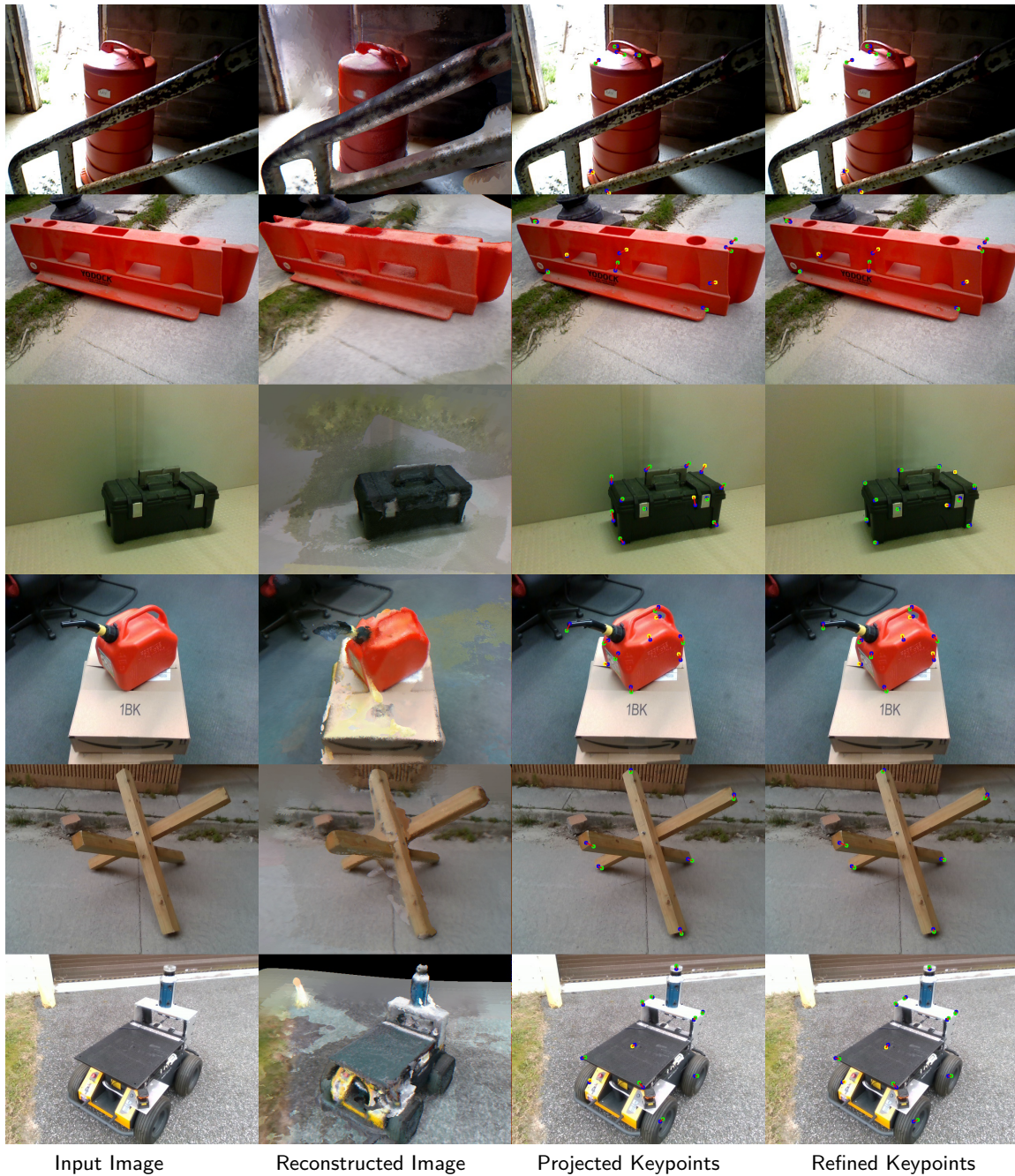


Figure 7. Reconstructed images and semi-automated keypoint projection. Keypoints are visualized as manual (blue dot), projected and visible (green dot), or projected and occluded (yellow dot). Red lines show the distances between manual keypoints and corresponding projected/refined keypoints.

consistently achieved performance comparable to the model trained with the manually labeled keypoints, and had a lower error on four out of six object classes. The model trained with our refined keypoints perform worse only on the barrel and the hedgehog, the only two objects with rotational symmetry.

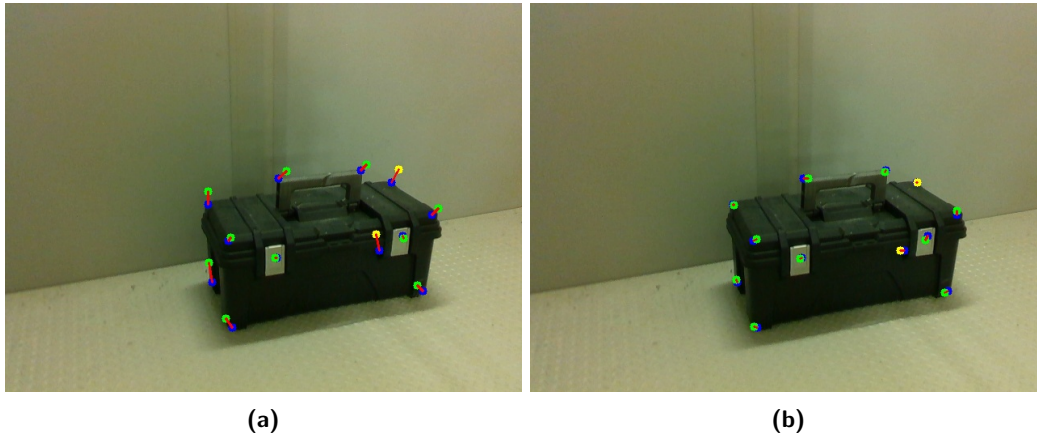


Figure 8. Demonstration of object-level refinement. Yellow circles indicate that a keypoint is occluded, green circles indicate visible keypoints, and blue circles indicate the manual annotations. For sequences with errors in sensor-trajectory estimation, in this case due to rapid camera motion with few available background features, naive projection from model to image frame, Figure 8a, exhibits poor alignment to manual annotations. After object-level refinement, Figure 8b, the keypoints are substantially improved.

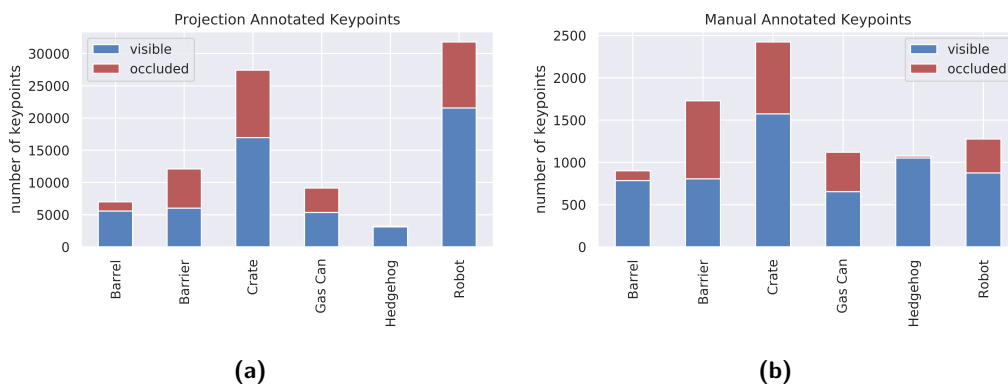


Figure 9. Total number of visible versus occluded keypoints in the our dataset for all projected frames 9a and manually annotated frames 9b (note difference in scale between the plots). As a result of not performing full loops around symmetric objects such as the Barrel and Hedgehog, the number of occluded keypoints for symmetric objects is much smaller than for other objects.

We conclude with an analysis of annotation effort required. Measurements of the required time are shown in Figure 11. As expected, using semi-automated projection allows the user to label an order of magnitude more frames, within the same quantity of time, as they would be able to label with a standard annotation tool. The annotations for each object class, with approximately 100 total frames per class labeled using standard image annotation, take about the same time as it takes to annotate a 3D reconstructed model, which we project to label thousands of images per class. For this analysis, we subsampled the number of projected frames by a factor of 10; therefore the true number of potentially annotated frames using our method is even higher than shown. We chose to subsample the output frames to reduce the total projection time as well as to avoid training with essentially redundant camera angles.

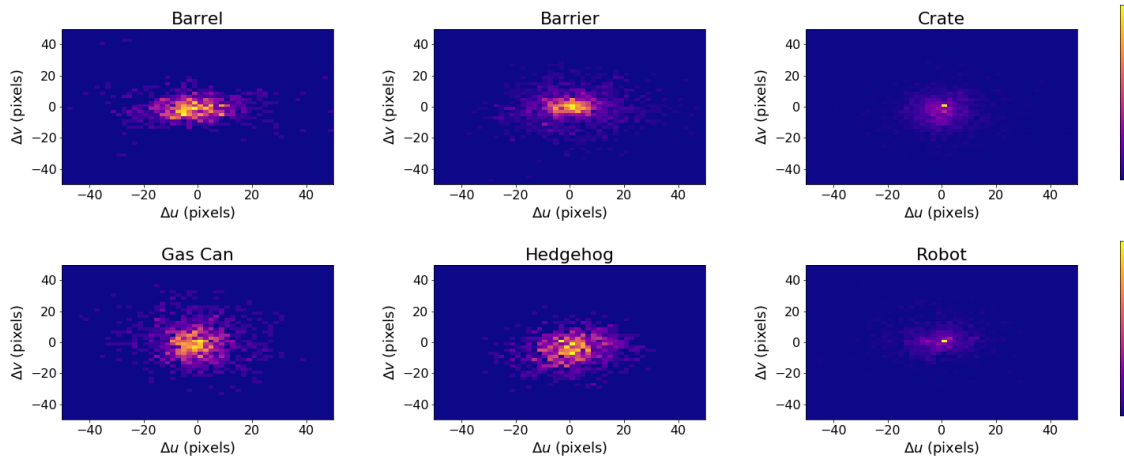


Figure 10. Pixel-wise distance between manual annotations and projected keypoints. The distance is a function of both the scene reconstruction pipeline accuracy as well as the manual annotation accuracy.

Table 9. Keypoint prediction performance on test data. ℓ_2 Distance of predicted keypoints locations in pixel space from manually labeled test set. The models trained with data from the refined labeling pipeline outperform the models trained with the manually labeled data at predicting the locations of manually labeled keypoints on four out of six object classes, while requiring less labeling effort.

	Manual	Projected	Refined
Barrel	6.7 ± 0.7	7.0 ± 0.8	8.9 ± 0.9
Barrier	14.7 ± 1.0	15.1 ± 1.1	14.1 ± 1.1
Crate	6.5 ± 0.6	22.9 ± 0.96	5.6 ± 0.6
Gascan	13.4 ± 0.8	5.7 ± 0.56	8.4 ± 0.8
Hedgehog	1.6 ± 0.1	2.0 ± 0.1	1.9 ± 0.1
Robot	7.2 ± 0.7	5.0 ± 0.6	3.0 ± 0.3

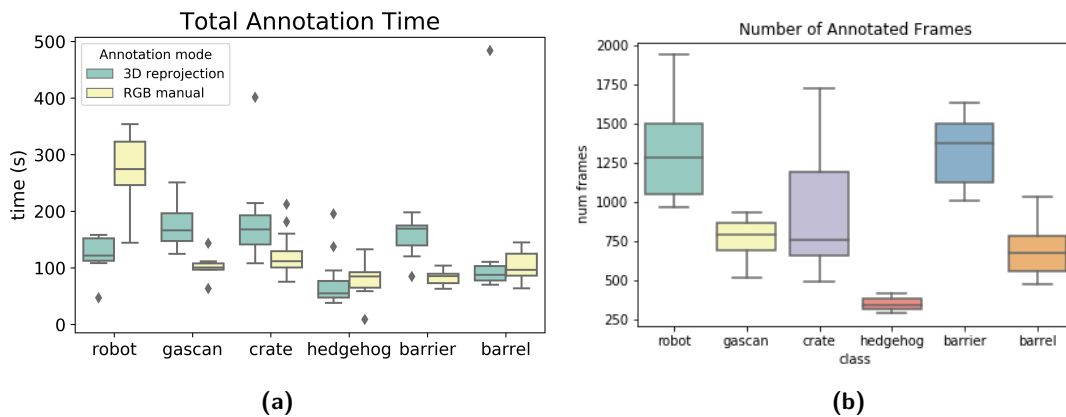


Figure 11. Analysis between manual versus semi-automated projection labeling. For each collected sequence for a given class, human annotation times (a) and sequence length (b) are shown. A human annotator is able to annotate a 3D model in approximately the same time it takes to annotate 100 individual image frames. The average number of annotated frames for each class is between one and two orders of magnitude greater for the same effort when using our annotation tool.










7. Conclusion

In this paper, we presented an efficient method to estimate the continuous 6-DoF pose of an object from a single RGB image. Capitalizing on the robust, semantic keypoint predictions provided by a state-of-the-art convnet, we proposed a pose optimization scheme that fits a deformable shape model to the 2D keypoints and recovers the object's 6-DoF pose. To ameliorate the effect of false detections, our pose optimization scheme integrates heatmap values, which reflect predictive confidence, to model each detection's certainty. Both the weak perspective and the full perspective cases were investigated. Moreover, we incorporated a technique for generating data from unlabeled videos, enabling semi-automatic labeling, such that our pipeline can be trained with minimal effort during the annotation stage. The experimental validation includes evaluation on instance-based pose estimation on the LineMOD-Occluded, YCB-Video, and TUD-Light datasets as well as evaluation on category-level pose estimation on the PASCAL3D+ dataset. Additionally, our method is accompanied by an efficient implementation with a running time under 0.3 seconds, making it a good fit for near-realtime robotics applications.

Acknowledgments

This research was sponsored by the U.S. Army Research Laboratory (ARL) under Cooperative Agreement W911NF-10-2-0016. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of ARL or the U.S. Government. We would like to thank Shiyani Patel for her help with data collection and analysis, as well as RCTA program leaders, in particular Stuart Young, Dilip Patel, Dave Baran, and Geoff Slipher.

ORCID

Karl Schmeckpeper¹  <https://orcid.org/0000-0003-4989-2022>
 Philip R. Osteen²  <https://orcid.org/0000-0001-8266-9848>
 Yufu Wang¹  <https://orcid.org/0000-0001-9907-8382>
 Georgios Pavlakos^{3,*}  <https://orcid.org/0000-0001-5821-1909>
 Kenneth Chaney¹  <https://orcid.org/0000-0003-1768-6136>
 Wyatt Jordan²  <https://orcid.org/0000-0002-5099-8388>
 Xiaowei Zhou^{4,*}  <https://orcid.org/0000-0003-1926-5597>
 Konstantinos G. Derpanis⁵  <https://orcid.org/0000-0001-5097-3802>
 Kostas Daniilidis¹  <https://orcid.org/0000-0003-0498-0758>

References

- Besl, P. J. and McKay, N. D. (1992). Method for registration of 3-d shapes. In *Sensor fusion IV: control paradigms and data structures*, volume 1611, pages 586–606. International Society for Optics and Photonics.
- Boumal, N., Mishra, B., Absil, P.-A., and Sepulchre, R. (2014). Manopt, a Matlab toolbox for optimization on manifolds. *Journal of Machine Learning Research*, 15:1455–1459.
- Bowman, S. L., Atanasov, N., Daniilidis, K., and Pappas, G. J. (2017). Probabilistic data association for semantic slam. In *ICRA*, pages 1722–1729.
- Brachmann, E., Krull, A., Michel, F., Gumhold, S., Shotton, J., and Rother, C. (2014). Learning 6D object pose estimation using 3D object coordinates. In *ECCV*, pages 536–551.
- Brachmann, E., Michel, F., Krull, A., Ying Yang, M., Gumhold, S., et al. (2016). Uncertainty-driven 6D pose estimation of objects and scenes from a single rgb image. In *CVPR*, pages 3364–3372.
- Calli, B., Singh, A., Walsman, A., Srinivasa, S., Abbeel, P., and Dollar, A. M. (2015). The ycb object and model set: Towards common benchmarks for manipulation research. In *2015 international conference on advanced robotics (ICAR)*, pages 510–517. IEEE.

- Cao, C., Weng, Y., Lin, S., and Zhou, K. (2013). 3D shape regression for real-time facial animation. *TOG*, 32(4):41.
- Cao, Z., Sheikh, Y., and Banerjee, N. (2016). Real-time scalable 6DOF pose estimation for textureless objects. In *ICRA*.
- Carreira, J., Agrawal, P., Fragkiadaki, K., and Malik, J. (2016). Human pose estimation with iterative error feedback. In *CVPR*.
- Castrejon, L., Kundu, K., Urtasun, R., and Fidler, S. (2017). Annotating object instances with a polygon-RNN.
- Chen, Y. and Medioni, G. (1992). Object modeling by registration of multiple range images. *Image Vision Comput.*, 10:145–155.
- Choi, S., Zhou, Q.-Y., and Koltun, V. (2015). Robust reconstruction of indoor scenes. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5556–5565.
- Choi, S., Zhou, Q.-Y., Miller, S., and Koltun, V. (2016). A large dataset of object scans. *arXiv:1602.02481*.
- Collet, A., Berenson, D., Srinivasa, S. S., and Ferguson, D. (2009). Object recognition and full pose registration from a single image for robotic manipulation. In *ICRA*, pages 48–55.
- Collet, A., Martinez, M., and Srinivasa, S. S. (2011). The MOPED framework: Object recognition and pose estimation for manipulation. *IJRR*, 30(10):1284–1306.
- Cootes, T. F., Taylor, C. J., Cooper, D. H., and Graham, J. (1995). Active shape models- Their training and application. *CVIU*, 61(1):38–59.
- Dai, J., He, K., and Sun, J. (2016). Instance-aware semantic segmentation via multi-task network cascades. In *CVPR*.
- Do, T.-T., Cai, M., Pham, T., and Reid, I. (2018). Deep-6dpose: Recovering 6d object pose from a single rgb image. *arXiv preprint arXiv:1802.10367*.
- Doumanoglou, A., Kouskouridas, R., Malassiotis, S., and Kim, T.-K. (2016). Recovering 6D object pose and predicting next-best-view in the crowd. In *CVPR*, pages 3583–3592.
- Fidler, S., Dickinson, S. J., and Urtasun, R. (2012). 3D object detection and viewpoint estimation with a deformable 3D cuboid model. In *NIPS*, pages 620–628.
- Fischler, M. A. and Bolles, R. C. (1981). Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography. *Commun. ACM*, 24(6):381–395.
- Gu, C. and Ren, X. (2010). Discriminative mixture-of-templates for viewpoint classification. In *ECCV*, pages 408–421.
- He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *CVPR*.
- Hinterstoisser, S., Cagniart, C., Ilic, S., Sturm, P. F., Navab, N., Fua, P., and Lepetit, V. (2012). Gradient response maps for real-time detection of textureless objects. *TPAMI*, 34(5):876–888.
- Hodan, T., Barath, D., and Matas, J. (2020). Epos: estimating 6d pose of objects with symmetries. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11703–11712.
- Hodaň, T., Michel, F., Brachmann, E., Kehl, W., Glent Buch, A., Kraft, D., Drost, B., Vidal, J., Ihrke, S., Zabulis, X., Sahin, C., Manhardt, F., Tombari, F., Kim, T.-K., Matas, J., and Rother, C. (2018). BOP: Benchmark for 6D object pose estimation. *European Conference on Computer Vision (ECCV)*.
- Hodaň, T., Sundermeyer, M., Drost, B., Labbé, Y., Brachmann, E., Michel, F., Rother, C., and Matas, J. (2020). BOP challenge 2020 on 6D object localization. *European Conference on Computer Vision Workshops (ECCVW)*.
- Hou, T., Ahmadyan, A., Zhang, L., Wei, J., and Grundmann, M. (2020). Mobilepose: Real-time pose estimation for unseen objects with weak shape supervision. *arXiv preprint arXiv:2003.03522*.
- Hu, Y., Hugonot, J., Fua, P., and Salzmann, M. (2019). Segmentation-driven 6d object pose estimation. In *CVPR*.
- Hua, B.-S., Pham, Q.-H., Nguyen, D. T., Tran, M.-K., Yu, L.-F., and Yeung, S.-K. (2016). Scenenn: A scene meshes dataset with annotations. In *3DV*.
- Izadi, S., Kim, D., Hilliges, O., Molyneaux, D., Newcombe, R., Kohli, P., Shotton, J., Hodges, S., Freeman, D., Davison, A., and Fitzgibbon, A. (2011). Kinectfusion: Real-time 3d reconstruction and interaction using a moving depth camera. In *UIST '11 Proceedings of the 24th annual ACM symposium on User interface software and technology*, pages 559–568. ACM.
- Kazhdan, M. and Hoppe, H. (2013). Screened poisson surface reconstruction. *ACM Trans. Graph.*, 32(3).
- Ke, L., Li, S., Sun, Y., Tai, Y.-W., and Tang, C.-K. (2020). Gsnet: Joint vehicle pose and shape reconstruction with geometrical and scene-aware supervision. In *European Conference on Computer Vision*, pages 515–532. Springer.

- Kehl, W., Manhardt, F., Tombari, F., Ilic, S., and Navab, N. (2017). Ssd-6d: Making rgb-based 3d detection and 6d pose estimation great again. In *Proceedings of the IEEE international conference on computer vision*, pages 1521–1529.
- Kessens, C. C., Fink, J., Hurwitz, A., Kaplan, M., Osteen, P. R., Rocks, T., Rogers, J., Stump, E., Quang, L., DiBlasi, M., Gonzalez, M., Patel, D., Patel, J., Patel, S., Weiker, M., Bowkett, J., Detry, R., Karumanchi, S., Burdick, J., Matthies, L., Oza, Y., Agarwal, A., Dornbush, A., Likhachev, M., Schmeckpeper, K., Daniilidis, K., Kamat, A., Choudhury, S., Mandalika, A., and Srinivasa, S. (2020). Toward fieldable human-scale mobile manipulation using RoMan. In Pham, T., Solomon, L., and Rainey, K., editors, *Artificial Intelligence and Machine Learning for Multi-Domain Operations Applications II*, volume 11413, pages 418 – 437. International Society for Optics and Photonics, SPIE.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). ImageNet classification with deep convolutional neural networks. In *NIPS*, pages 1106–1114.
- Labbé, Y., Carpentier, J., Aubry, M., and Sivic, J. (2020). Cosypose: Consistent multi-view multi-object 6d pose estimation. In *European Conference on Computer Vision*, pages 574–591. Springer.
- LeCun, Y., Boser, B. E., Denker, J. S., Henderson, D., Howard, R. E., Hubbard, W. E., and Jackel, L. D. (1989). Backpropagation applied to handwritten zip code recognition. *Neural Computation*, 1(4):541–551.
- Lee, C.-Y., Xie, S., Gallagher, P., Zhang, Z., and Tu, Z. (2015). Deeply-supervised nets. In *AISTATS*, volume 2, page 6.
- Lepetit, V., Moreno-Noguer, F., and Fua, P. (2009). EPnP: An accurate $O(n)$ solution to the PnP problem. *IJCV*, 81(2):155–166.
- Li, Y., Wang, G., Ji, X., Xiang, Y., and Fox, D. (2018). Deepim: Deep iterative matching for 6d pose estimation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 683–698.
- Li, Z., Wang, G., and Ji, X. (2019). 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.
- Liu, J., Zou, Z., Ye, X., Tan, X., Ding, E., Xu, F., and Yu, X. (2020). Leaping from 2d detection to efficient 6dof object pose estimation. In *Computer Vision – ECCV 2020 Workshops*, pages 707–714.
- Long, J., Zhang, N., and Darrell, T. (2014). Do convnets learn correspondence? In *NIPS*, pages 1601–1609.
- Lowe, D. G. (2004). Distinctive image features from scale-invariant keypoints. *IJCV*, 60(2):91–110.
- Manuelli, L., Gao, W., Florence, P., and Tedrake, R. (2019). Kpam: Keypoint affordances for category-level robotic manipulation. *International Symposium on Robotics Research (ISRR)*.
- Marion, P., Florence, P. R., Manuelli, L., and Tedrake, R. (2018). Label fusion: A pipeline for generating ground truth labels for real rgbd data of cluttered scenes. In *ICRA*.
- Massa, F., Aubry, M., and Marlet, R. (2014). Convolutional neural networks for joint object detection and pose estimation: A comparative study. *CoRR*, abs/1412.7190.
- Michel, F., Kirillov, A., Brachmann, E., Krull, A., Gumhold, S., Savchynskyy, B., and Rother, C. (2017). Global hypothesis generation for 6D object pose estimation. In *CVPR*.
- Montserrat, D. M., Chen, J., Lin, Q., Allebach, J. P., and Delp, E. J. (2019). Multi-view matching network for 6d pose estimation. *arXiv preprint arXiv:1911.12330*.
- Mousavian, A., Anguelov, D., Flynn, J., and Kosecka, J. (2017). 3d bounding box estimation using deep learning and geometry. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7074–7082.
- Muja, M., Rusu, R. B., Bradski, G. R., and Lowe, D. G. (2011). REIN - A fast, robust, scalable recognition infrastructure. In *ICRA*, pages 2939–2946.
- Murthy, J. K., Krishna, G., Chhaya, F., and Krishna, K. M. (2017). Reconstructing vehicles from a single image: Shape priors for road scene understanding. *ICRA*.
- Narayanan, P., Yeh, B., Holmes, E., Martucci, S., Schmeckpeper, K., Mertz, C., Osteen, P., and Wigness, M. (2020). An integrated perception pipeline for robot mission execution in unstructured environments. In *Artificial Intelligence and Machine Learning for Multi-Domain Operations Applications II*, volume 11413, page 1141318. International Society for Optics and Photonics.
- Newell, A., Yang, K., and Deng, J. (2016). Stacked hourglass networks for human pose estimation. In *ECCV*.
- Osteen, P. R., Owens, J. L., and Kaukeinen, B. (2019). Reducing the cost of visual DL datasets. In Pham, T., editor, *Artificial Intelligence and Machine Learning for Multi-Domain Operations Applications*, volume 11006, pages 121 – 139. International Society for Optics and Photonics, SPIE.

- Park, K., Patten, T., and Vincze, M. (2019). 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.
- Pavlakos, G., Zhou, X., Chan, A., Derpanis, K. G., and Daniilidis, K. (2017). 6-DoF object pose from semantic keypoints. *ICRA*, pages 2011–2018.
- Peng, S., Liu, Y., Huang, Q., Zhou, X., and Bao, H. (2019). Pvnnet: Pixel-wise voting network for 6dof pose estimation. In *CVPR*, pages 4561–4570.
- Pepik, B., Stark, M., Gehler, P. V., and Schiele, B. (2012). Teaching 3D geometry to deformable part models. In *CVPR*, pages 3362–3369.
- Qin, Z., Fang, K., Zhu, Y., Fei-Fei, L., and Savarese, S. (2019). Keto: Learning keypoint representations for tool manipulation. *International Conference on Robotics and Automation (ICRA)*.
- Rad, M. and Lepetit, V. (2017). Bb8: A scalable, accurate, robust to partial occlusion method for predicting the 3d poses of challenging objects without using depth. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3828–3836.
- Ramakrishna, V., Kanade, T., and Sheikh, Y. (2012). Reconstructing 3D human pose from 2D image landmarks. In *ECCV*, pages 573–586.
- Ren, S., He, K., Girshick, R., and Sun, J. (2015). Faster R-CNN: Towards real-time object detection with region proposal networks. In *NIPS*, pages 91–99.
- Rios-Cabrera, R. and Tuytelaars, T. (2013). Discriminatively trained templates for 3D object detection: A real time scalable approach. In *ICCV*, pages 2048–2055.
- Rusu, R. B., Blodow, N., and Beetz, M. (2009). Fast point feature histograms (FPFH) for 3D registration. In *ICRA*, pages 3212–3217.
- Salti, S., Tombari, F., and di Stefano, L. (2014). Shot: Unique signatures of histograms for surface and texture description. *CVIU*, 125:251–264.
- Sohn, K., Berthelot, D., Li, C.-L., Zhang, Z., Carlini, N., Cubuk, E. D., Kurakin, A., Zhang, H., and Raffel, C. (2020). Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *Advances in Neural Information Processing Systems (NeurIPS)*.
- Song, C., Song, J., and Huang, Q. (2020). Hybridpose: 6d object pose estimation under hybrid representations. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 431–440.
- Su, H., Qi, C. R., Li, Y., and Guibas, L. J. (2015a). Render for CNN: viewpoint estimation in images using CNNs trained with rendered 3D model views. In *ICCV*, pages 2686–2694.
- Su, H., Qi, C. R., Li, Y., and Guibas, L. J. (2015b). Render for cnn: Viewpoint estimation in images using cns trained with rendered 3d model views. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2686–2694.
- Sundermeyer, M., Marton, Z.-C., Durner, M., and Triebel, R. (2020). Augmented autoencoders: Implicit 3d orientation learning for 6d object detection. *International Journal of Computer Vision*, 128(3):714–729.
- Tekin, B., Sinha, S. N., and Fua, P. (2018). Real-time seamless single shot 6d object pose prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 292–301.
- Tieleman, T. and Hinton, G. (2012). Lecture 6.5-rmsprop, coursera: Neural networks for machine learning. *University of Toronto, Technical Report*.
- Toshev, A. and Szedgy, C. (2014). DeepPose: Human pose estimation via deep neural networks. In *CVPR*, pages 1653–1660.
- Tremblay, J., To, T., Sundaralingam, B., Xiang, Y., Fox, D., and Birchfield, S. (2018). Deep object pose estimation for semantic robotic grasping of household objects. *Conference on Robot Learning (CoRL)*.
- Tulsiani, S. and Malik, J. (2015). Viewpoints and keypoints. In *CVPR*, pages 1510–1519.
- Vasilopoulos, V. and Koditschek, D. E. (2018). Reactive navigation in partially known non-convex environments. In *International Workshop on the Algorithmic Foundations of Robotics*, pages 406–421. Springer.
- Vasilopoulos, V., Pavlakos, G., Bowman, S. L., Caporale, J. D., Daniilidis, K., Pappas, G. J., and Koditschek, D. E. (2020a). Reactive Semantic Planning in Unexplored Semantic Environments Using Deep Perceptual Feedback. *RAL*, 5(3):4455–4462.
- Vasilopoulos, V., Pavlakos, G., Schmeckpeper, K., Daniilidis, K., and Koditschek, D. E. (2020b). Reactive navigation in partially familiar planar environments using semantic perceptual feedback. *arXiv preprint arXiv:2002.08946*.

- Wang, C., Martín-Martín, R., Xu, D., Lv, J., Lu, C., Fei-Fei, L., Savarese, S., and Zhu, Y. (2019). 6-pack: Category-level 6D pose tracker with anchor-based keypoints. *International Conference on Robotics and Automation (ICRA)*.
- Wang, G., Manhardt, F., Shao, J., Ji, X., Navab, N., and Tombari, F. (2020). Self6d: Self-supervised monocular 6d object pose estimation. In *European Conference on Computer Vision*, pages 108–125. Springer.
- Wei, S.-E., Ramakrishna, V., Kanade, T., and Sheikh, Y. (2016). Convolutional pose machines. In *CVPR*.
- Whelan, T., Kaess, M., Johannsson, H., Fallon, M., Leonard, J. J., and McDonald, J. (2015). Real-time large-scale dense rgb-d slam with volumetric fusion. *The International Journal of Robotics Research*, 34(4-5):598–626.
- Whelan, T., Salas-Moreno, R. F., Glocker, B., Davison, A. J., and Leutenegger, S. (2016). Elasticfusion: Real-time dense slam and light source estimation. *The International Journal of Robotics Research*, 35(14):1697–1716.
- Xiang, Y., Mottaghi, R., and Savarese, S. (2014). Beyond PASCAL: A benchmark for 3D object detection in the wild. In *WACV*, pages 75–82.
- Xiang, Y., Schmidt, T., Narayanan, V., and Fox, D. (2017). Posecnn: A convolutional neural network for 6d object pose estimation in cluttered scenes. *Robotics: Science and Systems (RSS)*.
- Xie, Q., Dai, Z., Hovy, E., Luong, M.-T., and Le, Q. V. (2019). Unsupervised data augmentation for consistency training. *Advances in Neural Information Processing Systems (NeurIPS)*.
- Xie, Z., Singh, A., Uang, J., Narayan, K. S., and Abbeel, P. (2013). Multimodal blending for high-accuracy instance recognition. In *IROS*, pages 2214–2221.
- Zeng, A., Yu, K.-T., Song, S., Suo, D., Walker, E., Rodriguez, A., and Xiao, J. (2017). Multi-view self-supervised deep learning for 6d pose estimation in the amazon picking challenge. In *2017 IEEE international conference on robotics and automation (ICRA)*, pages 1386–1383. IEEE.
- Zhou, X., Karpur, A., Luo, L., and Huang, Q. (2018). Starmap for category-agnostic keypoint and viewpoint estimation. In *ECCV*, pages 318–334.
- Zhou, X., Leonardos, S., Hu, X., and Daniilidis, K. (2015a). 3D shape estimation from 2D landmarks: A convex relaxation approach. In *CVPR*, pages 4447–4455.
- Zhou, X., Zhu, M., Leonardos, S., Derpanis, K., and Daniilidis, K. (2015b). Sparseness meets deepness: 3D human pose estimation from monocular video. In *CVPR*.
- Zhu, M., Derpanis, K. G., Yang, Y., Brahmabhatt, S., Zhang, M., Phillips, C., Lecce, M., and Daniilidis, K. (2014). Single image 3D object detection and pose estimation for grasping. In *ICRA*, pages 3936–3943.
- Zhu, M., Zhou, X., and Daniilidis, K. (2015). Single image pop-up from discriminatively learned parts. In *ICCV*, pages 927–935.
- Zia, M. Z., Stark, M., Schiele, B., and Schindler, K. (2013). Detailed 3D representations for object recognition and modeling. *TPAMI*, 35(11):2608–2623.
- Zuo, Y., Qiu, W., Xie, L., Zhong, F., Wang, Y., and Yuille, A. L. (2019). Craves: Controlling robotic arm with a vision-based economic system. In *CVPR*, pages 4214–4223.

How to cite this article: Schmeckpeper, K., Osteen, P. R., Wang, Y., Pavlakos, G., Chaney, K., Jordan, W., Zhou, X., Derpanis, K. G., & Daniilidis, K. (2022). Semantic keypoint-based pose estimation from single RGB frames. *Field Robotics*, 2, 147–171.

Publisher’s Note: Field Robotics does not accept any legal responsibility for errors, omissions or claims and does not provide any warranty, express or implied, with respect to information published in this article.