

Constructing Category-Specific Models for Monocular Object-SLAM

Parv Parkhiya¹, Rishabh Khawad¹, J. Krishna Murthy¹, Brojeshwar Bhowmick², and K. Madhava Krishna¹

Abstract— We present a new paradigm for real-time object-oriented SLAM with a monocular camera. Contrary to previous approaches, that rely on *object-level* models, we construct *category-level* models from CAD collections which are now widely available. To alleviate the need for huge amounts of labeled data, we develop a rendering pipeline that enables synthesis of large datasets from a limited amount of manually labeled data. Using data thus synthesized, we learn category-level models for object deformations in 3D, as well as discriminative object features in 2D. These category models are instance-independent and aid in the design of object landmark observations that can be incorporated into a generic monocular SLAM framework. Where typical object-SLAM approaches usually solve only for object and camera poses, we also estimate object shape on-the-fly, allowing for a wide range of objects from the category to be present in the scene. Moreover, since our 2D object features are learned discriminatively, the proposed object-SLAM system succeeds in several scenarios where sparse feature-based monocular SLAM fails due to insufficient features or parallax. Also, the proposed category-models help in object instance retrieval, useful for Augmented Reality (AR) applications. We evaluate the proposed framework on multiple challenging real-world scenes and show — to the best of our knowledge — first results of an instance-independent monocular object-SLAM system and the benefits it enjoys over feature-based SLAM methods.

I. INTRODUCTION

Simultaneous Localization and Mapping (SLAM) finds various real-world applications such as autonomous navigation, visual inspection, mapping, and surveillance. Monocular cameras have evolved as popular choices for SLAM, especially on platforms such as hand-held devices and Micro Aerial Vehicles (MAVs). Most state-of-the-art monocular SLAM systems [1] operate on geometric primitives such as points, lines, and planar patches. Others operate directly on images, without the need for expensive feature extraction steps [2]. However, both these sets of approaches lack the ability to provide a rich semantic description of the scene.

Recognizing and keeping track of objects in a scene will enable a robot to build meaningful maps and scene descriptions. Object-SLAM is a relatively new paradigm [3]–[5] towards achieving this goal. Summarized in one sentence, object-SLAM attempts to augment SLAM with object information so that robot localization, object location estimation (in some cases, object pose estimation too), and mapping are achieved in a unified framework.

There are two dominant paradigms in object-SLAM research, depending on the way objects are characterized in

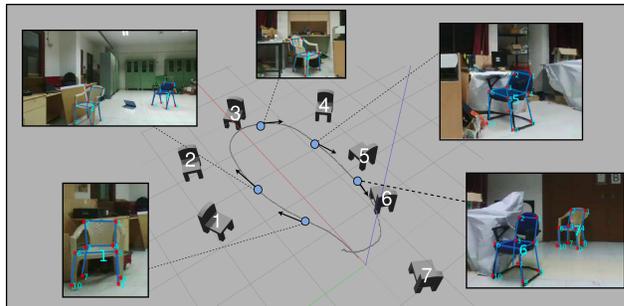


Fig. 1: Sample output of the proposed object-SLAM system. Given RGB images from a monocular camera (here from a quadrotor), we estimate the trajectory of the camera, as well as the poses and shapes of various objects in the scene. The proposed *category-specific models* can be used either in an incremental setting for online SLAM, or in a batch setting for offline factor-graph optimization.

the SLAM framework. In the first paradigm [4], [6], *object-level* (instance-specific) models are assumed to be available beforehand. However, the very nature of monocular SLAM with scale ambiguity coupled with the loss in information due to projection onto the image plane renders this paradigm infeasible for monocular object-SLAM systems. The second paradigm [7], [8], assumes a *generic model*, regardless of the object category. For instance, [8] models all objects as ellipsoids, and [5], [9] model all objects as cuboids. Both these approaches suffer a few disadvantages. Relying on object-level models will result in the need to have precise object models for all instances of an object category. On the other hand, generic models do not give much information about an object beyond the object category label. In many applications, such as manipulation for instance, it is advantageous to know the object pose.

In this paper, we propose a new paradigm for monocular object-SLAM, that combines the best-of-both-worlds. To enjoy the expressive power of instance-specific models yet retain the simplicity of generic models, we construct *category-specific models*, i.e., the object category is modeled as a whole. We employ the widely used linear subspace model [10]–[12] to characterize an object category and define object observations as factors in a SLAM factor graph [13], [14]. In our object-SLAM formulation, we do not assume any knowledge about the instance (interchangeably referred to as *shape*) of the object. Rather, we explicitly solve for the object shape in a joint formulation. The object-SLAM backend estimates robot trajectory and map, as well as poses and shapes of all objects in the scene.

Naturally, one would expect that a lot of data will be

¹ Parv Parkhiya, Rishabh Khawad, J. Krishna Murthy, and K. Madhava Krishna are with the Robotics Research Center, KCIS, IIT Hyderabad, India. albert.author@papercept.net

² Brojeshwar Bhowmick is with Tata Consultancy Services, Kolkata, India.

needed to learn category-specific models that generalize well across object instances, and rightly so. Datasets such as ShapeNet, SceneNet, ObjectNet have made available CAD collections of various object categories. We exploit the ready availability of such CAD collections to construct our category models. These category models capture the deformation modes of objects in 3D. Correspondingly, we leverage recent successes of Convolutional Neural Networks (CNNs) in keypoint localization [10], [15]–[17] to train 2D object feature extractors. To alleviate the need for large amounts of manually annotated training data, we design a *render* pipeline, along the lines of RenderForCNN [18], to synthesize enormous amounts of training data for category model learning. The presented render pipeline takes in a small volume of manually annotated data and synthesizes a large dataset that can be used to efficiently train a 2D object feature extraction network. We show that feature detectors learned from the render pipeline are more precise compared to ones that are learned over real data alone, which corroborates claims in [18].

We evaluate our object-SLAM system on multiple challenging real-world sequences and present — to the best of our knowledge — first steps towards instance-independence in monocular object-SLAM. Since we use discriminative 2D features on objects, our system is robust to conditions such as strong rotation, where monocular SLAM approaches usually face catastrophic failure. We present both incremental and batch versions of our object-SLAM pipeline and demonstrate its advantages qualitatively and quantitatively over feature-based visual SLAM approaches [1]. Finally, we show that, using our category-level models, one can perform object instance retrieval and this can be used in many Augmented Reality (AR) applications for overlaying object models in a scene. Fig. 1 illustrates an output from our pipeline. Objects are consistently embedded onto the robot’s trajectory and their 3D models are rendered.

II. RELATED WORK

Nearly all state-of-the-art SLAM systems [1], [2], [19] rely on pose-graph (or other factor graph) optimization [20], [21]. In this section we review related work on object-SLAM, and outline certain limitations in them that form motivating factors for the proposed approach.

A. Object-SLAM

With recent advances and subsequent stabilization of SLAM systems, the community has been devoting attention to incorporating objects into the SLAM framework. To this end, some recent approaches for object-oriented SLAM have been proposed [3]–[5], [7]–[9].

Most of these works rely on depth information from RGB-D or stereo sensors [4]–[7]. In [4], [6], assume an instance-level model of the objects are known a priori. In [4], a real-time 3D object detection algorithm is applied on an RGB-D image stream and these objects are fused along with odometry information in a pose graph optimization scheme. Similarly in [6], a framework for multi-robot object-SLAM

is proposed. Again, each robot is equipped with an RGB-D sensor and object models are available a priori.

There is another paradigm in which no instance-level models are available a priori. In [5], association and object poses are solved for jointly, in a factor graph framework, using data from RGB-D cameras. Among monocular object-SLAM/SfM approaches, [8], [9] fall under this paradigm. In such approaches, objects are modeled as bounding boxes [7], [9] or as ellipsoids [8].

Our method hence falls under a *third* paradigm, where we assume category-models, and not instance-level models.

B. Object-Category Models

Over the last few years, object-category models have been applied to several problems in monocular vision. In [10]–[12], category-level models are employed to obtain object reconstructions from single images. These approaches demonstrate that the loss of information in the monocular imaging process can be compensated for by incorporating priors on shapes of objects belonging to a particular category.

We use these category models and exploit them to design object observation factors that can be easily incorporated into monocular SLAM, and also generalize across several instances from the category, without the need for modeling all possible instances from the category.

C. Keypoint Localization Using CNNs

Convolutional Neural Networks (CNNs) have been the driving reason behind recent advances in object detection [22], [23] and object keypoint localization [15]–[17], [24]. When run on a GPU, these CNNs are capable of processing image frames with latencies of about 100-300 milliseconds and form important components of our pipeline.

D. Render Pipelines for Data Synthesis

With the advent of CAD model collections such as [25], 3D data is now in abundance. In [18], image synthesis using a rendering engine is proposed as an alternative to training on real images annotated manually. Models trained for the task of object viewpoint prediction on rendered data (and subsequently finetuned on a smaller dataset comprising real data) are shown to outperform models that have been trained on (larger) real datasets alone. Our experiments corroborate this fact for the task of object keypoint prediction as well.

We build on several of the components described here, however we craft the outputs to create object factors that can be augmented to factor graphs [13] constructed using monocular SLAM approaches. The entire pipeline is summarized in Fig. 2 and is explained in the subsequent sections.

III. CONSTRUCTING CATEGORY-SPECIFIC MODELS

In this section, we describe the category-level models that we employ in our object-SLAM system. Our key insight is that shapes of objects from the same category are not *arbitrary*. Instead they all follow a set pattern that can be learned by examining several instances. We adopt the widely used linear subspace model [10]–[12], [15] and represent objects as 3D wireframes.

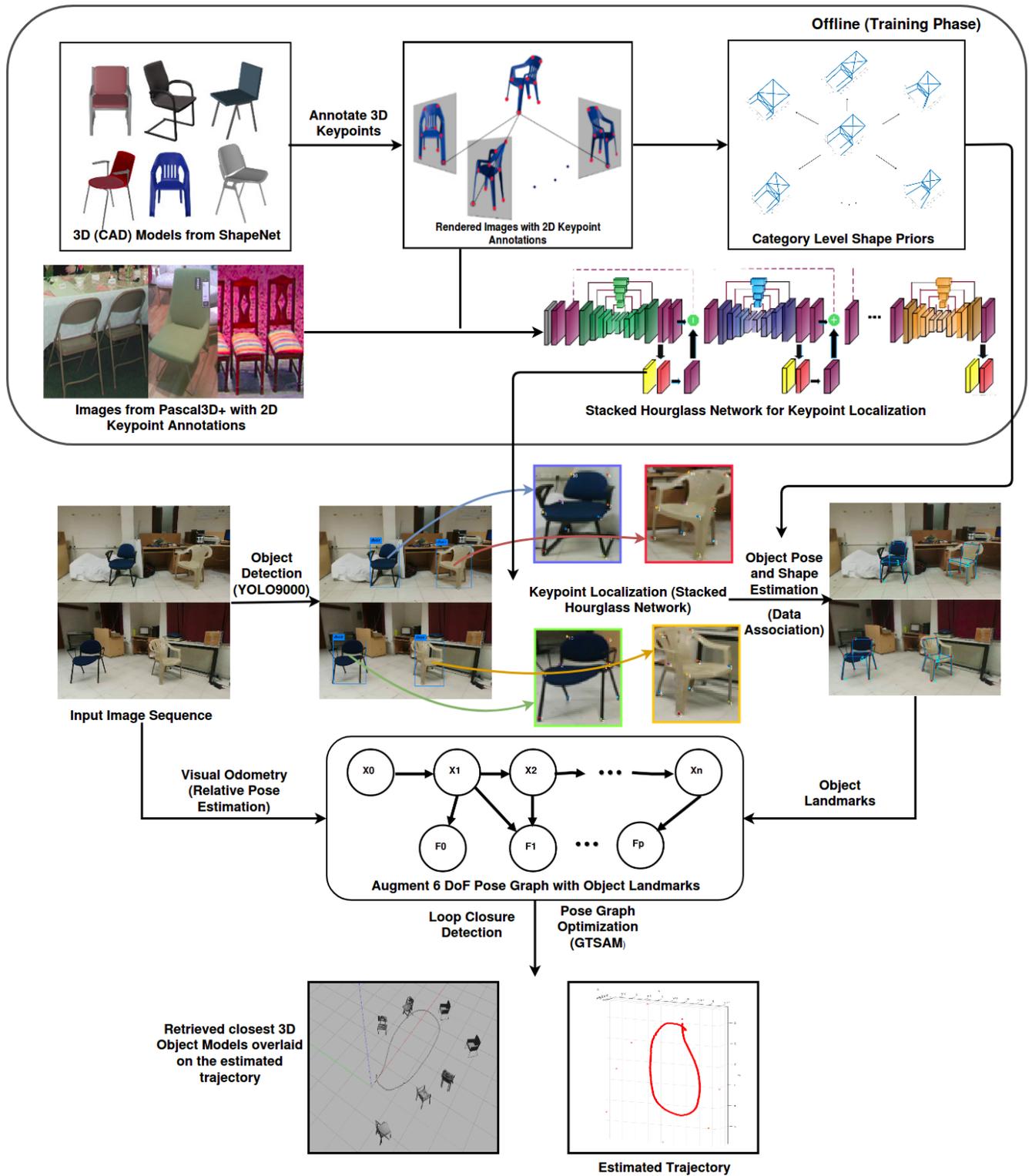


Fig. 2: Complete pipeline for learning and incorporating category-level models for object-SLAM

A. Category-Level Model

In the proposed approach, we lay an emphasis on the use of *category-level* models as opposed to *instance-level* models for objects. To construct a category level model, each object is first characterized as a set of 3D locations of salient points

that are common across all instances of the category. For example, such salient points for the *chair* category could be legs of the chair, corners of the chair backrest, and so on. These points are common across all instances of the category. We denote by \mathbf{S} the $3K$ -vector comprising of an ordered

collection of 3D locations of these K salient points, which we refer to as *keypoints* for that object category. Keypoints for the *chair* category which we use throughout this work, are shown in Fig. 2.

According to the linear subspace model [10], the space of actual shapes of a category are confined to a much lower dimensional subspace of \mathbb{R}^{3K} . This is easy to see, since there are several dependencies that exist among the K keypoints, such as subsets of the keypoints being planar, symmetric, etc. Any object can then be expressed as a *mean shape* that can be deformed along linearly independent directions (*basis vectors*) using *shape parameters* (coefficients of the linear combination).

Mathematically, if $\bar{\mathbf{S}}$ is the mean shape of the category, and \mathbf{V}_b s are a deformation basis obtained from PCA over a collection of aligned ordered 3D CAD models (such as [25]),

$$\mathbf{S} = \bar{\mathbf{S}} + \sum_{b=1}^B \lambda_b \mathbf{V}_b = \bar{\mathbf{S}} + \mathbf{V}\mathbf{\Lambda} \quad (1)$$

where B is the number of basis vectors (the top- B eigenvectors after PCA). Here, \mathbf{V} represents the learnt $3K \times B$ deformation modes. $\mathbf{\Lambda}$ is a K -vector containing the deformation coefficients. Varying the deformation coefficients produces various shapes from the learnt shape subspace, as shown in Fig. 2 (panel: *Category Level Shape Priors*).

B. Discriminative Feature Extraction

Since our category-level models rely on keypoints for that category, we need reliable discriminative feature extractors that localize the corresponding keypoints in 2D images. For the same, we adopt the *stacked hourglass* model introduced in [17] and subsequently modified in [15].

C. Render Pipeline

Inspired by RenderForCNN [18], we implement our customized render pipeline for generating huge amounts of synthetic keypoint annotated chair images using a small set of 3D annotated keypoints. We briefly summarize the steps in our render pipeline, and how we exploit its advantages for learning 3D category-models as well as discriminative 2D feature extractors.

Render Pipeline for 3D Category-level Models: First, we choose a collection of CAD models of chairs, comprising of about 250 chairs sampled from the ShapeNet [25] repository. For each chair, we synthesize a few (typically 8) 2D images with predetermined viewpoints (azimuth, elevation, and camera-tilt angles). Keypoints in these images are then annotated (in 2D) manually, and then triangulated to 3D to obtain 3D keypoint locations on the CAD model.

Since the models are already assumed to be aligned, performing a Principal Component Analysis (PCA) over the (mean-subtracted) 3D keypoint locations results in the deformation basis (eigenvector obtained from PCA). This constitutes the category-level model learning phase.

Render Pipeline for 2D Feature Extraction

For training 2D feature extractors, we use the 3D keypoint annotated CAD models to synthesize 2D images with

Blender [26] as the rendering engine. Apart from using the RenderForCNN [18] framework to randomly sample view parameters and generate 2D images with overlaid backgrounds, we also perform a projection of the 3D keypoints to 2D using the same view parameters to obtain 2D keypoint annotations for training the stacked hourglass architecture [15]–[17] for keypoint localization.

Advantages of the Proposed Render Pipeline

- Only 2D annotation is required. Annotation on 3D CAD models requires expensive labeling effort. We, however, perform annotations on a small set of rendered 2D images and triangulate them to 3D.
- Keypoint annotations will be available even for occluded parts. This is one major drawback of keypoint localization architectures trained over real data. Since there are no reliable estimates of keypoint locations for occluded parts, these labels are not present in the ground-truth in datasets such as PASCAL3D+ [27]. We demonstrate that our models perform much better due to the availability of occluded keypoint locations as well.
- Using only small volumes of labeled data, we can generate millions of training examples for the keypoint localization CNN.

An overview of the render pipeline developed is illustrated in Fig. 2.

IV. OBJECT MEASUREMENT FACTORS USING CATEGORY-MODELS

In this section, we describe how the category-models thus learned can be incorporated into a monocular SLAM backend. We first introduce the notation we use throughout the section. $\mathbf{T}_{ij} \in SE(3)$ denotes a rigid-body transform that takes a 3D point expressed in the camera frame at time i and expresses it with respect to the camera frame at time j . Hence, $\mathbf{T}_{ij} = \mathbf{T}_j \mathbf{T}_i^{-1}$ (\mathbf{T}_* denotes a transformation matrix that transforms points from frame $*$ to the world frame W). Each such matrix \mathbf{T}_{ij} is a 4×4 matrix and takes the following form.

$$\mathbf{T}_{ij} = \begin{bmatrix} \mathbf{R}_{ij} & \mathbf{t}_{ij} \\ \mathbf{0} & 1 \end{bmatrix} \text{ where } \mathbf{R}_{ij} \in SO(3), \mathbf{t}_{ij} \in \mathbb{R}^3 \quad (2)$$

Each camera frame i is related to camera frame j in the following manner. If ${}^i\mathbf{X}$ denotes the 3D coordinates of a world point ${}^w\mathbf{X}$ in the i th camera frame, ${}^j\mathbf{X}$ is given by ${}^j\mathbf{X} = \mathbf{T}_{ij} {}^i\mathbf{X}$. We use $\pi : \mathbb{R}^4 \mapsto \mathbb{R}^2$ that projects 3D homogeneous coordinates to 2D Euclidean coordinates. We subsume the camera intrinsics f_x, f_y, c_x, c_y into the function π .

$$\pi \begin{pmatrix} X \\ Y \\ Z \\ 1 \end{pmatrix} = \begin{pmatrix} \frac{f_x X}{Z} + c_x \\ \frac{f_y Y}{Z} + c_y \end{pmatrix} \quad (3)$$

We use $\text{Log}(\cdot)$ to denote the logarithm map that maps an element of the group $SE(3)$ to a corresponding element in the exponential coordinates of its tangent space $se(3)$.

A. 3D Pose-SLAM

We formally define the 3D pose-SLAM problem as follows. Given a set \mathcal{Z} of *relative* pose measurements $\hat{\mathbf{T}}_{ij}$ among robot poses i, j , where $i, j \in \{1..N\}$, estimate \mathbf{T}_i for all $i \in \{1..N\}$ such that the log-likelihood of the observations (relative pose measurements) is maximized. This reduces to the problem of minimizing the observation errors (minimizing the negative log likelihood). We assume that each relative-pose measurement $\hat{\mathbf{T}}_{ij}$ has an associated uncertainty Σ_{ij} .

$$\min_{\mathbf{T}_i, i \in \{1..N\}} \mathcal{E}_{pose} = \sum_{\hat{\mathbf{T}}_{ij} \in \mathcal{Z}} \|\text{Log}(\hat{\mathbf{T}}_{ij}^{-1} \mathbf{T}_j \mathbf{T}_i^{-1})\|_{\Sigma_{ij}} \quad (4)$$

One popular approach of minimizing the error in (4) is the use of factor graphs [13], [21]. We adopt this approach since it naturally extends to an incremental optimization framework that enables online processing.

B. Object Observation Factors

Using the linear subspace model illustrated in (1), we now design object observation factors that are suitable for instance-independent monocular object-SLAM. Given 2D locations of object keypoints \mathbf{s} , we use the pose and shape adjustment scheme proposed in [10] to *lift* the 2D keypoints to 3D, thereby estimating the shape and pose of the object from just a single image.

We assume that, from the i^{th} pose, the robot observes M objects (M could also be zero). The set of object observations in frame i is denoted \mathcal{O}_i , and each observation in the set is indexed as $\hat{\mathbf{T}}_i^{\mathcal{O}_m}$. We denote the number of keypoints of the object category by K . The k^{th} keypoint of the m^{th} object observed are denoted by \mathbf{s}_k^m . The pose and shape of each object can then be computed by minimizing the following keypoint reprojection error (see [10] for details).

$$\min_{\hat{\mathbf{T}}_i^{\mathcal{O}_m}, \Lambda_m} \left(\sum_{m=1}^M \sum_{k=1}^K \|\pi(\mathbf{K} \hat{\mathbf{T}}_i^{\mathcal{O}_m} (\bar{\mathbf{S}} + \mathbf{V} \Lambda_m)) - \mathbf{s}_k^m\|_2^2 \right) + \rho(\Lambda_m) \quad (5)$$

In the above equation, $\rho(\Lambda_m)$ denotes an appropriate regularizers (such as an L_2 norm, for instance) to prevent shape parameter (Λ) estimates from deviating too much from the category-model.

We estimate each object's shape (Λ) and pose ($\hat{\mathbf{T}}_i^{\mathcal{O}_m}$) in each frame i by alternating minimization of the error term in (5), with respect to the pose and shape parameters. If the same object has been associated successfully across multiple frames, we also exploit temporal consistency [15] for more precise estimates.

C. Object-SLAM

The object pose observations arising from (5) now form additional factors in the SLAM factor graph. If $\hat{\mathbf{T}}_i^{\mathcal{O}_m}$ denotes the pose of object m with respect to the i^{th} camera frame, for each object node in the factor graph, the following pose error is to be minimized.

TABLE I: Keypoint localization accuracy (2D) for the stacked hourglass network trained using the proposed render pipeline

Test Accuracy	Train Accuracy		
	Synthetic only	Real only	Synthetic + Real
Synthetic data	90.67%	85.90%	94.51%
Pascal3D+	75.99%	87.46%	93.4%
Sequence 1 (Ours)	21.11%	66.18%	95.93%

$$\min_{\substack{\mathbf{T}_i, i \in \{1..N\} \\ \hat{\mathbf{T}}_i^{\mathcal{O}_m}, m \in \{1..M\}}} \mathcal{E}_{obj} = \sum_{i=1}^N \sum_{m=1}^M \|\text{Log}(\left(\hat{\mathbf{T}}_{\mathcal{O}_m}^i\right)^{-1} \mathbf{T}_i^{-1} \mathbf{T}^{\mathcal{O}_{\phi(m)}})\| \quad (6)$$

Here, $\mathbf{T}^{\mathcal{O}_{\phi(m)}}$ denotes the pose of object $\mathcal{O}_{\phi(m)}$ with respect to the global coordinate frame, where $\phi(m)$ denotes the data association function which assigns a globally unique identifier to each distinct object observed so far. The data association pipeline is briefly outlined in Section V.

Finally, the object-SLAM error that jointly estimates robot poses as well as object poses from relative pose observations is given by

$$\mathcal{E} = \mathcal{E}_{pose} + \mathcal{E}_{obj} \quad (7)$$

V. IMPLEMENTATION DETAILS

We use the publicly available GTSAM [21] framework for constructing and optimizing the proposed factor graph model. Robot pose observations are obtained from a visual odometry/SLAM approach. We use odometry information published by ORB-SLAM [1] for this purpose. However, monocular SLAM inherently suffers from scale ambiguity, i.e., the absolute scale of reconstruction cannot *usually* be recovered. To recover actual scale, we use information from a range sensor to estimate the height of the drone (camera) above the ground. Thereby, using objects that lie on the ground plane (chairs, for instance), we can compute the absolute scale factor by backprojection via the ground plane [10], [28]. Moreover, ShapeNet [25] CAD models for the categories we use are available in metric scale. Hence, a rough prior on object dimensions is available and can be used to initialize the solution for faster convergence [10].

A. Object Detection and Data Association

We use the YOLO9000 [22] object detector to detect objects of interest and perform non-maximum suppression. Each bounding box detection is then passed to the keypoint localization network.

Between successive frames, objects are tracked using a greedy tracker based on the Hungarian algorithm. For long-term tracking and for Object Loop Closure detection (OLC), we use the estimated shape and pose parameters as costs in the Hungarian algorithm. This has shown to be an extremely simple, yet effective data association method in [29].

The object shape and pose are optimized using Ceres solver [30] and are cast as observations into the factor graph. The factor graph can be optimized in incremental or batch mode, and with or without object loop closures.

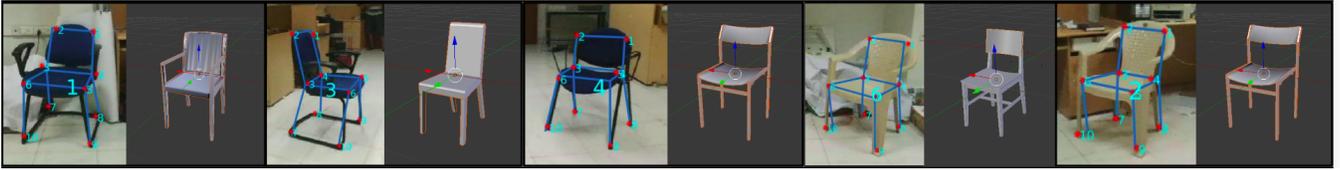


Fig. 3: Qualitative results for instance retrieval. Query images and their retrieved CAD models are shown.

VI. RESULTS

A. Object-SLAM Evaluation

In this section, we present experimental results upon evaluating various modes of operation of the proposed object-SLAM approach and the render pipeline on multiple real-world sequences. By explicitly modeling objects in a SLAM framework, we achieve substantially lower object localization errors over feature-based SLAM [1]. Also, we highlight the use of discriminatively learned features that enable us to obtain trajectory and object pose estimates even in scenarios where monocular SLAM approaches fail, such as strong rotational motion.

B. Dataset

To demonstrate object SLAM in a real-world setting, we collected a dataset comprising of 4 monocular video sequences from robots operating in office and laboratory environments. These sequences were collected from a micro aerial vehicle (MAV) flying at a constant height above the floor. For one of the sequences, fiducial markers were placed in the environment and were used to estimate ground truth camera pose and object locations. For two other sequences, no markers were placed in the environment, to deny any undue advantage enjoyed by monocular SLAM. In these sequences, all object locations were measured a priori and this is used in evaluating the object localization accuracy. There was a *rotation-only* run, where the MAV rotates more-or-less *in-place*. In all these runs, we evaluate object localization precision for ORB-SLAM [1] and for various flavors of the proposed object-SLAM approach.

A brief description of the dataset is as follows.

- Sequence 1 : An elongated loop with two parallel sides exhibiting dominant straight motion (7.6 meters) (fiducial markers were used for acquiring ground-truth).
- Sequence 2 : Smaller run with smoother turns (no fiducial markers).
- Sequence 3 : 360° Rotation in place without any translation from origin.
- Sequence 4 : 2 meters forward, sharp 90° rotation towards left, 4.2 meters forward

In all these sequences we consider *chairs* as the objects of interest. For a fair evaluation, and to demonstrate the efficacy of the proposed object-SLAM system, we ensure that most of the frames in the sequence contain various chair instances.

C. Render Pipeline

For the task of training CNNs to detect keypoints for the chair category, we use our render pipeline to synthesize a little over a million keypoint annotated images of various models of chairs. We train the stacked hourglass [15], [17]

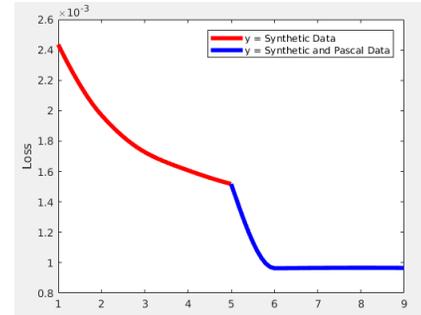


Fig. 4: Progress of training the hourglass CNN using the proposed render pipeline. Loss per epoch when training on synthetic data followed by finetuning on a mix of real and synthetic images is illustrated here.

architecture on this synthesized dataset. We then finetune the CNN on a smaller dataset that comprises of real chair images from the PASCAL3D+ [27] dataset.

Keypoint localization accuracies for several configurations of the hourglass network are shown in TABLE I. We see that CNNs trained on synthetic data alone fail to generalize to real data from our sequences. Similarly, training on real data (from the PASCAL3D+ dataset) alone performs fairly well on test samples from the same dataset, but fails to generalize to other kinds of data, such as our sequence. However, finetuning the CNN on a combination of both real and synthetic data leads to better generalization. Fig. 4 shows the decrease of loss with time when the best performing approach from TABLE I is trained using our render pipeline.

D. Discussion

On the collected dataset, we evaluate the several variants of the proposed object-SLAM approach, viz. batch mode (the factor graph is constructed for the entire scene and is optimized in one go), incremental mode (observations are processed as they arrive). We also evaluate object-SLAM with and without object-based loop closures. TABLE II summarizes the results of these experiments.

For each of the approaches, we evaluate object localization error (in meters), and also the accumulated drift in the trajectory in the X and Z directions. The accumulated drift is computed only in cases where the robot returns to the starting position at the end of the run. Significant drifts in the Z-direction occur in sequence 1 for ORB-SLAM [1], but these are corrected by the object-SLAM system, regardless of the presence or absence of object loop closures. In sequence 4, there is a strong rotation (a perpendicular turn), where ORB-SLAM performs poorly.

Our discriminative features enjoy an advantage over tracked features [1] in that they need to be visible in just a

TABLE II: Quantitative analysis of several variants of the proposed object-SLAM approach. Here, the best and worst object localization errors refer to the errors in the most accurately localized and the most erroneously localized object respectively. We also show the average object localization error and accumulated endpoint drifts, wherever applicable.

Sequence ID	Approach	Mode	# Objects	Object localization Error (metres)			Drift of trajectory	
				Best	Worst	Avg	Z direction (in meters)	X direction (in meters)
1	ORB-SLAM [1]	-	11	0.3162	4.9535	2.8539	4.3	0.9051
	Ours (without object loop closure)	batch		0.2526	2.7017	1.0318	2.051	0.18
	Ours (with object loop closure)	batch		0.2193	1.6584	0.779	0.804	0.326
		incremental		0.2188	1.6756	0.7998	0.79	0.209
2	ORB-SLAM [1]	-	7	0.0607	0.8734	0.3734	0.22	0.6181
	Ours (without object loop closure)	batch		0.0791	0.7306	0.3125	0.19	0.4309
	Ours (with object loop closure)	batch		0.0738	0.6005	0.2419	0.04	0.13
		incremental		0.072	0.5936	0.2364	0.04	0.13
3	ORB-SLAM [1]	-	9	ORB-SLAM Fails to initialize				
	Ours (with object loop closure)	batch		0.0915	0.9998	0.521	0.0194	0.0556
		incremental		0.0825	1.0174	0.5219	0.0045	0.06
4	ORB-SLAM [1]	-	7	0.1113	0.5993	0.3657	0.101	N/A
	Ours (without object loop closure)	batch		0.0661	0.4688	0.3086	0.055	N/A
		incremental		0.0614	0.4678	0.3051	0.059	N/A

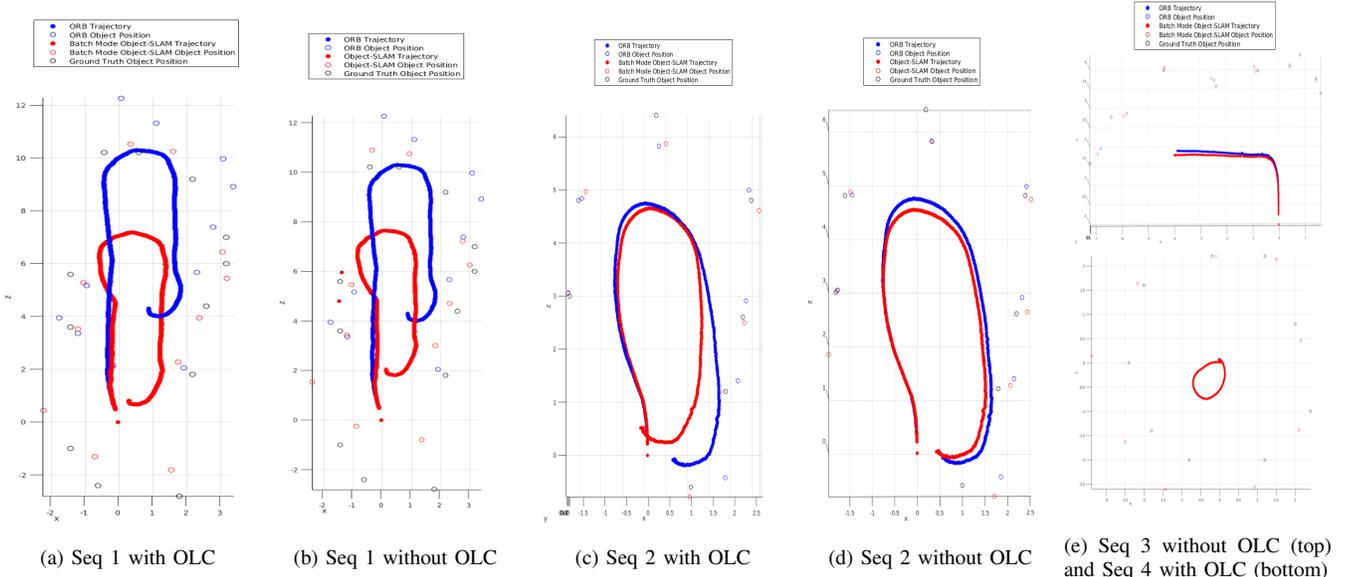


Fig. 5: Estimated Trajectories and Object Locations by ORB-SLAM and Object-SLAM with/without **Object Loop Closure detection (OLC)**

single image to enable object measurements. We demonstrate this in sequence 3 (cf. Fig. 6). In this sequence the robot purely rotates about its origin. While ORB-SLAM [1] fails to initialize for this sequence, the object-SLAM proposed here is able to accurately localize all chairs in the scene. The circle traced out by our object-SLAM systems has a radius roughly equal to that of the drone used in the run.

The incremental version performs competitively, and in cases, better, than the batch-mode object-SLAM. Loop closures based on object observations further boost the accuracy of the proposed approach. For batch mode, we perform a chordal initialization before optimizing the factor graph.

Fig. 5 shows plots of the trajectories estimated by ORB-SLAM [1], as well as from object-SLAM. We also show object location estimates overlaid on ground-truth object locations for the sequence. Sequences 1 and 2 are plotted with and without incorporating object loop closure constraints, to show that object loop closure further reduces endpoint drift.

E. Instance Retrieval

The proposed category-level model certainly achieves instance independence, but instance retrieval itself may be of use in several scenarios, such as Augmented Reality (AR). The shape parameters estimated by our approach can be used to guide an instance search over the keypoint annotated CAD collection. Given a query image (2D), we run a K-Nearest Neighbors search to retrieve the closest instance possible from the CAD collection. In Fig. 3, we present results from running a 5-nearest neighbor search and manually choosing the closest instance.

These instances can then be used to render the objects in the scene as well as the robot trajectory. One such rendering for Sequence 3 is shown in Fig 6.

VII. CONCLUSIONS

In this paper, we have presented an approach for monocular object-SLAM that relies on category-specific models, rather than relying on instance-specific models or generic

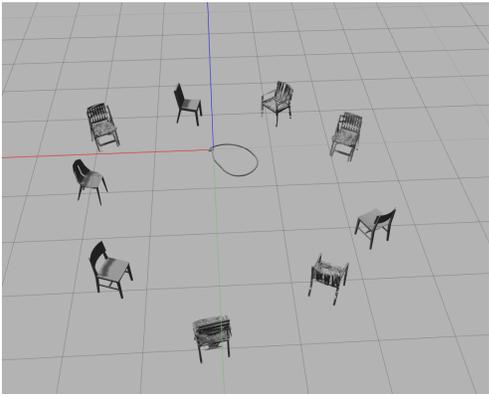


Fig. 6: Rendering objects estimated in a run where the robot rotates in place. ORB-SLAM [1] fails to initialize due to insufficient parallax.



Fig. 7: Qualitative results on another object category (laptops).

models. This is a new paradigm in the nascent field of object-SLAM research. The proposed category-specific models can be adopted to many other such rigid objects, as long as data for the corresponding category is available in the form of aligned CAD models (See Fig. for an example). Future work could focus on further reducing the extent of supervision required to learn such category-specific models.

VIII. ACKNOWLEDGEMENTS

The authors acknowledge the support and funding from Kohli Center for Intelligent Systems (KCIS) IIT Hyderabad, and Tata Consultancy Services (TCS) Innovation Labs India for this work. We also acknowledge the help of Akanksha Baranwal, Roopal Nahar, Danish Sodhi, Karnik Ram, and Sarthak Sharma for their timely assistance.

REFERENCES

- [1] M. J. M. M. Mur-Artal, Raúl and J. D. Tardós, “ORB-SLAM: a versatile and accurate monocular SLAM system,” *IEEE Transactions on Robotics*, vol. 31, no. 5, pp. 1147–1163, 2015.
- [2] J. Engel, T. Schöps, and D. Cremers, “Lsd-slam: Large-scale direct monocular slam,” in *European Conference on Computer Vision*. Springer, 2014, pp. 834–849.
- [3] N. Sünderhauf, T. Pham, Y. Latif, M. Milford, and I. D. Reid, “Meaningful maps - object-oriented semantic mapping,” *CoRR*, vol. abs/1609.07849, 2016. [Online]. Available: <http://arxiv.org/abs/1609.07849>
- [4] R. F. Salas-Moreno, R. A. Newcombe, H. Strasdat, P. H. Kelly, and A. J. Davison, “Slam++: Simultaneous localisation and mapping at the level of objects,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2013, pp. 1352–1359.
- [5] B. Mu, S.-Y. Liu, L. Paull, J. Leonard, and J. P. How, “Slam with objects using a nonparametric pose graph,” in *Intelligent Robots and Systems (IROS), 2016 IEEE/RSJ International Conference on*. IEEE, 2016, pp. 4602–4609.
- [6] S. Choudhary, L. Carlone, C. Nieto, J. Rogers, Z. Liu, H. I. Christensen, and F. Dellaert, “Multi robot object-based slam,” in *International Symposium on Experimental Robotics*. Springer, 2016, pp. 729–741.

- [7] N. Sünderhauf, F. Dayoub, S. McMahon, M. Eich, B. Upcroft, and M. Milford, “Slam-quo vadis? in support of object oriented and semantic slam,” 2015.
- [8] M. Crocco, C. Rubino, and A. Del Bue, “Structure from motion with objects,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 4141–4149.
- [9] D. Gálvez-López, M. Salas, J. D. Tardós, and J. Montiel, “Real-time monocular object slam,” *Robotics and Autonomous Systems*, vol. 75, pp. 435–449, 2016.
- [10] J. K. Murthy, G. S. Krishna, F. Chhaya, and K. M. Krishna, “Reconstructing vehicles from a single image: Shape priors for road scene understanding,” in *Proceedings of the IEEE Conference on Robotics and Automation*, 2017.
- [11] M. Z. Zia, M. Stark, and K. Schindler, “Towards scene understanding with detailed 3d object representations,” *International Journal of Computer Vision*, vol. 112, no. 2, pp. 188–203, 2015.
- [12] S. Tulsiani, A. Kar, J. Carreira, and J. Malik, “Learning category-specific deformable 3d models for object reconstruction,” *IEEE transactions on pattern analysis and machine intelligence*, 2016.
- [13] M. Kaess, H. Johannsson, R. Roberts, V. Ila, J. J. Leonard, and F. Dellaert, “isam2: Incremental smoothing and mapping using the bayes tree,” *The International Journal of Robotics Research*, vol. 31, no. 2, pp. 216–235, 2012.
- [14] D. M. Rosen, M. Kaess, and J. J. Leonard, “Rise: An incremental trust-region method for robust online sparse least-squares estimation,” *IEEE Transactions on Robotics*, vol. 30, no. 5, pp. 1091–1108, 2014.
- [15] J. K. Murthy, S. Sharma, and M. Krishna, “Shape priors for real-time monocular object localization in dynamic environments,” in *Proceedings of the IEEE Conference on Intelligent Robots and Systems*, 2017.
- [16] G. Pavlakos, X. Zhou, A. Chan, K. Derpanis, and K. Daniilidis, “6-dof object pose from semantic keypoints,” in *Proceedings of the IEEE Conference on Robotics and Automation (In Press)*, 2017.
- [17] A. Newell, K. Yang, and J. Deng, “Stacked hourglass networks for human pose estimation,” in *European Conference on Computer Vision*. Springer, 2016, pp. 483–499.
- [18] H. Su, C. R. Qi, Y. Li, and L. J. Guibas, “Render for cnn: Viewpoint estimation in images using cnns trained with rendered 3d model views,” in *The IEEE International Conference on Computer Vision (ICCV)*, December 2015.
- [19] C. Forster, Z. Zhang, M. Gassner, M. Werlberger, and D. Scaramuzza, “Svo: Semidirect visual odometry for monocular and multicamera systems,” *IEEE Transactions on Robotics*, vol. 33, no. 2, pp. 249–265, 2017.
- [20] R. Kümmerle, G. Grisetti, H. Strasdat, K. Konolige, and W. Burgard, “g 2 o: A general framework for graph optimization,” in *Robotics and Automation (ICRA), 2011 IEEE International Conference on*. IEEE, 2011, pp. 3607–3613.
- [21] F. Dellaert *et al.*, “Gtsam,” URL: <https://borg.cc.gatech.edu>, 2012.
- [22] J. Redmon and A. Farhadi, “Yolo9000: Better, faster, stronger,” *arXiv preprint arXiv:1612.08242*, 2016.
- [23] R. Girshick, “Fast r-cnn,” in *International Conference on Computer Vision (ICCV)*, 2015.
- [24] S. Tulsiani and J. Malik, “Viewpoints and keypoints,” in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2015, pp. 1510–1519.
- [25] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, *et al.*, “Shapenet: An information-rich 3d model repository,” *arXiv preprint arXiv:1512.03012*, 2015.
- [26] Blender Online Community, *Blender - a 3D modelling and rendering package*, Blender Foundation, Blender Institute, Amsterdam, 2017. [Online]. Available: <http://www.blender.org>
- [27] Y. Xiang, R. Mottaghi, and S. Savarese, “Beyond pascal: A benchmark for 3d object detection in the wild,” in *IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2014.
- [28] S. Song and M. Chandraker, “Joint sfm and detection cues for monocular 3d localization in road scenes,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3734–3742.
- [29] J. K. M. K. M. K. Sarthak Sharma, Junaid Ahmed Ansari, “Beyond pixels: Leveraging geometry and shape cues for online multi-object tracking,” in *Proceedings of the IEEE Conference on Robotics and Automation (In Press)*, 2018.
- [30] S. Agarwal, K. Mierle, and Others, “Ceres solver,” <http://ceres-solver.org>.