

# Scale Estimation with Dual Quadrics for Monocular Object SLAM

Shuangfu Song<sup>2</sup>, Junqiao Zhao<sup>\*1</sup>, Tiantian Feng<sup>2</sup>, Chen Ye<sup>1</sup>, Lu Xiong<sup>3</sup>

**Abstract**—The scale ambiguity problem is inherently unsolvable to monocular SLAM without the metric baseline between moving cameras. In this paper, we present a novel scale estimation approach based on an object-level SLAM system. To obtain the absolute scale of the reconstructed map, we formulate an optimization problem to make the scaled dimensions of objects conform to the distribution of their sizes in the physical world, without relying on any prior information about gravity direction. The dual quadric is adopted to represent objects for its ability to describe objects compactly and accurately, thus providing reliable dimensions for scale estimation. In the proposed monocular object-level SLAM system, semantic objects are initialized first from fitted 3-D oriented bounding boxes and then further optimized under constraints of 2-D detections and 3-D map points. Experiments on indoor and outdoor public datasets show that our approach outperforms existing methods in terms of accuracy and robustness.

## I. INTRODUCTION

Recently, object-level SLAM has achieved significant progress thanks to breakthroughs in convolutional neural networks (CNNs). A series of monocular object SLAM systems [1], [2], [3] are proposed. These semantically enriched SLAM algorithms can improve robot intelligence for scene understanding and human-robot interaction. However, due to the limitation of monocular sensors, the absolute scale between the reconstructed map and the physical world is unknown, which is known as the scale ambiguity problem. This largely hinders the application of these algorithms in actual scenarios. Hence, it is crucial and necessary to estimate the absolute scale in monocular SLAM systems.

For monocular cameras, the absolute scale is unobservable but can be inferred from prior knowledge such as the sizes of objects in the scene [4]. Based on this insight, some previous studies [5], [6] attempted to use the Bayesian framework to infer the absolute scale by associating estimated heights of observed objects with their priors. However, there are some limitations: First, when measuring the height of an object, it has to resort to the direction of gravity as the reference. Objects that are not upright are treated as anomalous. Second, only specific objects with prominent height can be used in the algorithm, other low objects such as “keyboard”, “book”,

etc., are discarded. As a result, the applicability and accuracy of these methods are limited.

In this paper, we propose a novel absolute scale estimation approach by exploring all reliable dimensions of general objects without resorting to the prior gravity direction. Therefore, more constraints from general objects can reduce the uncertainty of scale estimation and make the scaled dimensions of objects closer to those of the physical world.

To obtain accurate dimension estimation of objects for scale inferencing, we develop a monocular object SLAM system using the dual quadric as the three-dimensional (3-D) object representation. A dual quadric can flexibly fit various 3-D shapes tightly with compact parameters, thus the size of the object can be well measured. Furthermore, quadrics in dual vector space are well defined in projective geometry [7] and can be robustly initialized [8]. We first initialize the objects with a two-stage strategy based on an improved version of our previous work [8]. Then accurate dimensions can be extracted from dual quadrics reconstructed during semantic mapping. The Metric-tree [9], a repository of sizes of more than 900 object categories, is used to provide object size priors. We further introduce uncertainties to Metric-tree priors to ensure that the result will not be skewed by the prior sizes with large deviation. At last, an optimization problem is established to estimate the absolute scale factor.

In summary, our contributions are as follows:

- 1 An accurate and robust scale estimation approach exploiting dual quadrics and object size priors.
- 2 A monocular object SLAM system with a semantic mapping module, which can build object-oriented maps accurately and timely.

To the best of our knowledge, this is the first work that takes advantage of dual quadrics to estimate the scale factor within a monocular object SLAM system.

## II. RELATED WORK

### A. Scale Estimation

To overcome the monocular scale ambiguity problem, an intuitive way is to leverage extra sensors, such as inertial measurement units (IMUs) [10], LiDAR [11]. Without augmented sensors, at least one absolute reference needs to be integrated. Reference information used in existing approaches can be divided into three categories: camera setup parameters, learned models, and prior information of the semantic object.

On ground-based vehicle platforms, cameras are usually mounted at a specific height above the ground plane. [12], [13], [14] proposed to estimate unscaled camera height by

<sup>\*</sup>This work is supported by the National Key Research and Development Program of China (No. 2021YFB2501104), the National Natural Science Foundation of China (No. 41871370)

<sup>2</sup>Shuangfu Song and Tiantian Feng are with the School of Surveying and Geo-Informatics, Tongji University, Shanghai 200092, China (e-mail: 1911204@tongji.edu.cn; fengtiantian@tongji.edu.cn).

<sup>1</sup>Junqiao Zhao and Chen Ye are with the Department of Computer Science and Technology, School of Electronics and Information Engineering, Tongji University, Shanghai 201804, China (e-mail: zhaojunqiao@tongji.edu.cn; yechen@tongji.edu.cn).

<sup>3</sup>Lu Xiong is with the Institute of Intelligent Vehicle, Tongji University, Shanghai, 201804 China (e-mail: xiong.lu@tongji.edu.cn).

extracting the ground plane parameters derived from map points. Then, with the known camera height relative to the ground plane, the absolute scale can be obtained. These approaches are accurate and effective, but cannot cover more general scenarios, for example, where only unanchored hand-held cameras are available.

In [15], a monocular SLAM system was proposed to recover the scale by incorporating a depth prediction module based on deep convolutional neural fields. A similar idea is proposed in DVSO [16], in which deep depth predictions are used as virtual stereo measurements. The difference is that [16] trained the network with sparse depth reconstructions from Stereo DSO [17], instead of using depth measurements captured by a LiDAR. Another learning-based method [18] trained a space-coarse, but depth-accurate CNN for depth prediction, and directly estimated the scale factor for every frame. These approaches have achieved impressive performance on KITTI datasets. However, it is still a challenge for them to generalize to different environments.

To exploit the semantic object priors, methods in [5] and [6] adopt the Bayesian framework incorporating the height priors of objects to infer the global scale. Recent work [9] improved this approach by introducing not only heights but all available dimensions of objects into the scale estimation. However, these approaches all require the gravity direction known and objects to be placed upright to measure the size of objects.

### B. Object SLAM

Classic visual SLAM systems have been well explored in representing geometric primitives like points, lines, and planes. SLAM++ [19] is considered to be the first object-level SLAM system that uses RGB-D detection and represents objects by matching prior models. Similarly, [20] leveraged a large object database and uses bags of words to identify objects. Later, more general object representations were adopted to break through the restriction of the pre-prepared object database. [21] represented objects as spheres and incorporated them into bundle adjustment as extended points. QuadricSLAM [1] was proposed to represent objects flexibly using dual quadrics. CubeSLAM [2] described objects using cuboids. Inspired by them, EAO-SLAM [3] adopted both quadrics and cuboids to represent objects based on their different prior shapes. However, cuboids require additional constraints, such as uprightness, for the reconstruction, as its projection is not well defined in the projective geometry.

## III. ABSOLUTE SCALE ESTIMATION

### A. Problem Formulation

We use the ellipsoid, a closed surface form of the dual quadric, to describe a general 3-D object. It has nine parameters: 6 DoF pose and three dimensions i.e. length, width and height. We denote  $d$  as a unscaled dimension of reconstructed objects,  $\tilde{d}$  as its real length. The absolute scale factor  $s$  can

be expressed as

$$s = \frac{\tilde{d}}{d}. \quad (1)$$

Then the reconstructed map can be corrected by this scale. However, each dimension in each object can calculate a different local scale factor through its prior length, which is globally inconsistent; This requires us to use all dimensions of all objects as conditions to find a global optimal scale.

We assume that each dimension of general objects is conforming to a Gaussian distribution  $N(\mu, \sigma)$ , and dimensions are independent of each other. With the set of all dimensions  $D = \{d_i\}$  of objects, the conditional probability of scale  $s$  is presented as

$$P(s|D) \propto \prod_i P(d_i|s). \quad (2)$$

Our goal is to find the scale which maximizes the probability according to the distribution of dimensions in the physical world. The likelihood can then be presented as

$$P(d_i|s) = N(\mu_i, \sigma_i) \quad (3)$$

where  $\mu_i$  and  $\sigma_i$  are prior parameters to describe the distribution of  $d_i$  in the physical world. The error between scaled  $d_i$  and its expectation is defined as

$$e_i = \mu_i - sd_i. \quad (4)$$

The maximum likelihood estimation (MLE) can be formulated as a least-squares optimization problem

$$s^* = \arg \min_s \sum_i \left\| \frac{e_i}{\sigma_i} \right\|^2. \quad (5)$$

### B. Object Uncertainty Model

The uncertainty of a dimension, as shown in Equation 3, can be modeled by size statistics from the Metric-tree [9]. However, the accuracy of the scale estimate is also affected by the quality of the object reconstruction. We define a confidence  $c$  value to evaluate the reliability of an object based on observations used for dual quadrics construction:

$$c = \frac{w_1 c_{det} + w_2 c_{pt} + w_3 c_{vis}}{w_1 + w_2 + w_3} \quad (6)$$

where  $w_1, w_2, w_3$  are weight parameters.

$c_{det}$  is the confidence of 2-D detection and is derived by:

$$c_{det} = \frac{1}{n} \sum_k p_k$$

where  $p_k$  is the probability of detection obtained from the 2-D object detector in view  $k$ .

$c_{pt}$  is the confidence derived from  $N_p$ , the total number of map points associated with the object:

$$c_{pt} = \min(1, \max(0, \log_a N_p)).$$

$c_{vis}$  is the confidence of visibility and is derived by:

$$c_{vis} = \min(1, \max(0, \log_b N_o))$$

where  $N_o$  is the total number of 2D detections.  $a, b$  are super parameters set to be 10 and 15 respectively in our implementation.

Lastly, Equation 5 can be reformulated by integrating the confidence weight as:

$$s^* = \arg \min_s \sum_i \left\| \frac{c_i e_i}{\sigma_i} \right\|^2. \quad (7)$$

### C. Object Dimensions Selection

Besides the uncertainty, the reliability of the dimensions of an object should also be taken into account. For example, dimensions like the height of “bottle”, length of “spoon” are more stable than the thickness of “book” and “keyboard”. This is because small dimensions are difficult to estimate accurately during the mapping process. Therefore, objects are classified into three categories according to their dimensions. We denote three dimensions of a object as  $d_1, d_2, d_3$  following  $d_1 \geq d_2 \geq d_3 > 0$ . Then the shape feature of an object can be defined as:

$$L_d = \frac{d_1 - d_2}{d_1} \quad P_d = \frac{d_2 - d_3}{d_1} \quad S_d = \frac{d_3}{d_1} \quad (8)$$

where  $L_d, P_d, S_d$  are linearity, planarity, and scattering, respectively [22]. For an object with  $S_d < 0.3$ , it belongs to the pole-like when  $L_d > 0.5$ ; it belongs to the disk-like object when  $P_d > 0.5$ . For a pole-like object, only the longest dimension of it is stable enough to be used. For a disk-like object, the shortest dimension of it should be discarded. For others, all dimensions can provide useful hints for scale inferencing.

### D. Outliers Elimination

False 2-D object detections lead to the erroneous association between 3D objects and their corresponding size priors, therefore should be eliminated.

We found that most local scale factors estimated are consistent with each other. Therefore, outliers caused by false detection can be detected and eliminated by statistical methods such as the boxplot. First, all local scales are sorted in ascending order. Then the interquartile range (IQR) is defined as  $IQR = Q_3 - Q_1$ , where  $Q_1$  is the first quartile and  $Q_3$  is the third quartile. A dimension will be discarded if its local scale factor is less than  $Q_1 - 1.5IQR$  or greater than  $Q_3 + 1.5IQR$ .

### E. Implementation

Our scale estimation pipeline is shown in Figure 1. We first select all stable dimensions from map objects and calculate their local scales with prior dimension distribution for outliers elimination. Then we combine the prior variance of the dimension and the confidence of the reconstructed object to weight each error term of the scale optimization process. Scale estimation is embedded in the back-end optimization process of our SLAM system as described in Section IV. Every time the object map is updated, a new scale factor will be calculated automatically by g2o [23], and then be used to scale the whole map.

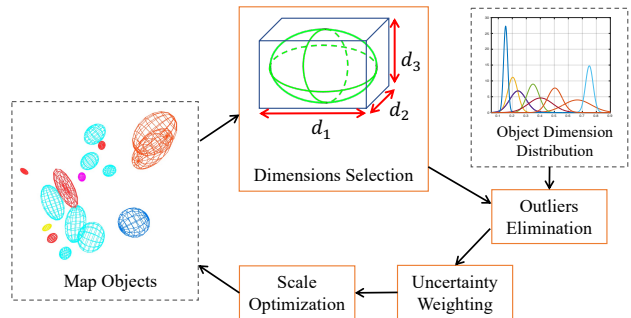


Fig. 1. The pipeline of scale estimation.

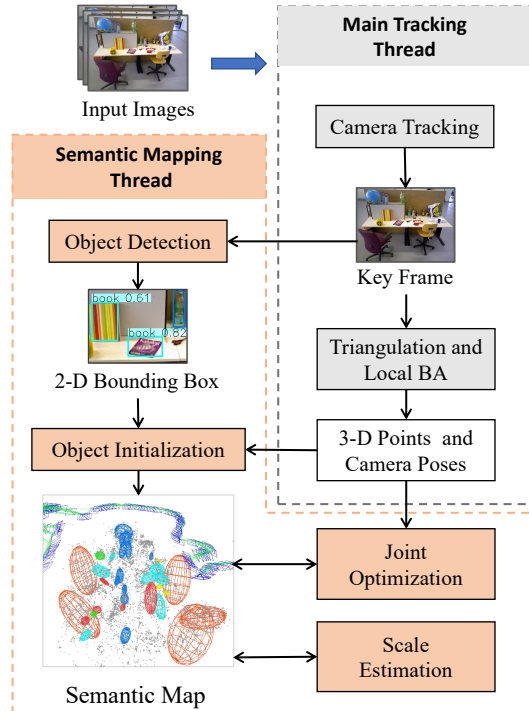


Fig. 2. The proposed object SLAM system framework.

## IV. SEMANTIC MAPPING

### A. Overview

To acquire accurate object size for scale estimation, we build an object SLAM system based on ORB-SLAM2 [24]. Only keyframes created by camera tracking are used for object detection. Then dual quadrics are initialized from 2-D bounding boxes and 3-D points in the semantic mapping thread. At last, all objects, points and camera poses are jointly optimized to update the semantic map. The framework of the system is shown in Figure 2.

### B. Dual Quadric Initialization

A 3-D dual quadric can be initialized from a set of planar constraints constructed by multi-view detections for this object [1]. However, this method only works when enough detections with diverse viewpoints are accumulated. Our previous work [8] constructed a 3-D convex hull from map points to provide valid planar constraints for dual

quadric initialization. In this paper, we propose an efficient 3-D oriented bounding box (OBB)-based approach for dual quadric initialization.

We first associate 3-D map points to their corresponding object if points are projected into the 2-D bounding box of this object. These map points are then clustered to filter out background points similar to [8]. Afterwards, OBB is extracted from associated map points to directly provide initial parameters (orientation, position, size) for dual quadrics. Here, the adopted covariance-based OBB fitting algorithm is proposed in [25]. With the accumulation of measurements, the dual quadric will be re-initialized from 2-D detections as done in [8].

### C. Joint Optimization

After initialization, objects are further optimized jointly with other map components. Denote the set of camera poses, dual quadric object and points as  $X = \{x_i \in SE(3)\}$ ,  $Q = \{q_j \in \mathbb{R}^{4 \times 4}\}$  and  $P = \{p_k \in \mathbb{R}^3\}$ , respectively. Three types of measurement errors are introduced into the bundle adjustment.

1) *Camera-Object Measurement Error*: A dual quadric, represented by a  $4 \times 4$  symmetric matrix  $Q^*$ , can be projected onto an image plane to obtain a dual conic represented by  $3 \times 3$  symmetric matrix  $C^*$ , following this rule:

$$C^* = H Q^* H^T \quad (9)$$

where  $H = K[R|\mathbf{t}]$  is the camera projection matrix composed of camera intrinsic and extrinsic parameters. As done in [26], we can obtain the predicted 2-D bounding box  $\hat{\mathbf{b}} = [\hat{u}_{max}, \hat{v}_{max}, \hat{u}_{min}, \hat{v}_{min}]$  of the dual quadric by

$$\hat{u}_{max}, \hat{u}_{min} = \frac{1}{C_{3,3}^*} \left( C_{1,3}^* \pm \sqrt{C_{1,3}^{*2} - C_{1,1}^* C_{3,3}^*} \right) \quad (10)$$

$$\hat{v}_{max}, \hat{v}_{min} = \frac{1}{C_{3,3}^*} \left( C_{2,3}^* \pm \sqrt{C_{2,3}^{*2} - C_{2,2}^* C_{3,3}^*} \right) \quad (11)$$

where  $[u_{max}, v_{max}], [u_{min}, v_{min}]$  represent the top left and bottom right corners of the 2-D box, respectively. Denote the 2-D bounding box measurement observed by the object detector as  $\mathbf{b}$ . Given a camera pose  $x$ , and a 3-D object  $q$ , the 4-D re-projected error can be defined as

$$\mathbf{e}(x, q) = \hat{\mathbf{b}} - \mathbf{b}. \quad (12)$$

2) *Camera-Point Measurement Error*: We use the standard 3-D point reprojection error, the same as in ORB-SLAM2 [24]:

$$\mathbf{e}(x, p) = \pi(T_c^{-1}p) - z \quad (13)$$

where  $\pi(\cdot)$  is the camera projective function,  $T_c$  is the camera pose, and  $z$  is the pixel coordinate measurement.

3) *Object-Point Measurement Error*: To exploit constraints between map objects and map points, we introduce a novel measurement error:

$$\mathbf{e}(q, p) = \max(0, \sqrt{p^T Q p + 1} - 1) \quad (14)$$

where  $Q$  is the primal quadric matrix with  $Q_{4,4} = -1$ . As shown in Figure 3, when point  $p$  is outside the ellipsoid,  $p^*$

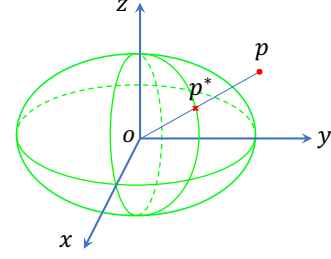


Fig. 3. The measurement error between quadrics and points.

is the intersection point between the line segment  $\overline{op}$  and the quadric surface. the error represents the ratio  $\overline{op^*}/\overline{op}$ . We use the max operator to encourage the quadric to wrap the associated map points. When a point lies inside the quadric, its measurement error will be zero.

Lastly, combining all measurement error items, the optimization problem can be formulated as

$$X^*, Q^*, P^* = \arg \min_{\{X, Q, P\}} \left\{ \sum_{i,j} \|e(x_i, q_j)\|_{\Omega_{ij}}^2 + \sum_{i,k} \|e(x_i, p_k)\|_{\Omega_{ik}}^2 + \sum_{j,k} \|e(q_j, p_k)\|_{\Omega_{jk}}^2 \right\} \quad (15)$$

where  $\Omega$  is the covariance matrix of different error measurements for Mahalanobis norm. This nonlinear least-squares problem can be solved efficiently using Levenberg-Marquardt algorithm.

## V. EXPERIMENTS

We evaluate the performance of our scale estimation approach on the TUM RGB-D data [27] as an indoor dataset and the KITTI tracking sequence [28] as an outdoor dataset. For monocular object SLAM, only RGB images are input and YOLOv3 [29] are adopted for 2-D object detection. We adopt the relative scale error (RSE) as the evaluation metric of scale estimation:

$$\text{RSE} = \frac{|\hat{s} - s|}{s} \quad (16)$$

where  $\hat{s}$  is our estimated global scale factor of the scene, and  $s$  is the ground truth scale factor which can be obtained from the similarity transformation  $S = \{s, R, \mathbf{t}\}$  between the trajectory estimated by our SLAM algorithm and the ground truth trajectory as following:

$$S^* = \arg \min_{\{s, R, \mathbf{t}\}} \sum_i \|p_i - sR\hat{p}_i - \mathbf{t}\|^2 \quad (17)$$

where  $\hat{p}_i, p_i$  represent the position of estimated trajectory, groundtruth trajectory respectively.

For the indoor dataset with common objects, we use the size information provided by the Metric-tree [9], which is collected from the internet, and construct a size distribution for each class of objects. For the outdoor dataset, only the ‘‘car’’ object is used as reference, and its prior size distribution is based on the statistics of KITTI 3-D object annotations. The remaining few classes of objects are either rarely seen in the scene or difficult to initialize successfully and are therefore discarded.

### A. Scale Estimation Results

There are two works [5], [9] similar to ours. The former first measures the heights of observed objects and then integrates the corresponding height priors into a Bayesian framework to infer the absolute scale. The latter uses a structure-from-motion (SFM) system running offline to reconstruct 3-D scenes and a pixel-wise instance segmentation method to extract objects from the reconstructed point cloud, which is difficult to compare fairly with our approach. Therefore, the former is chosen as our baseline. Both methods use identical object size priors and a scale consistency-based outlier rejection algorithm. We selected three sequences from the TUM dataset and six sequences from the KITTI dataset for evaluation. The selection criteria were two-fold: first, there should be representative semantic landmarks in the sequence; second, our SLAM system could run well on the sequence. We ran each algorithm ten times on all sequences to calculate the average RSE and its standard deviation.

As shown in Table I and Table II, our method achieves an average RSE of around 5% and 10% on TUM and KITTI sequences respectively, which is better than the baseline [5]. In the KITTI dataset, our method outperforms the baseline in most cases, except for the sequences 0000 and 0014, which are very short and have a limited number of objects available for scale estimation (from 2 to 4).

In the fr1\_desk sequence, the baseline method performed not well because there are few objects with stable height to use as reference, as shown in the far left of Figure 5(a). Our method effectively utilizes the size information of “book” and “keyboard” in this scene to obtain more accurate and robust results.

In the other two sequences of TUM, our method performs better despite the fact that the baseline method has more object height cues. This is because the baseline method measures object height based on 2D detection results, which means it is susceptible to object occlusion. Meanwhile, our method exploits multi-view observations of the object to reduce the uncertainty of dimension measurement, which improves the performance of scale estimation.

In addition, the scale estimation in the fr3-office sequence is slightly worse than in the other two scenes. The main reason is that the dimension of some bottles in this scene exceeds the regular size in the Metric-tree [9].

In KITTI long sequences (except 0000 and 0014 with a duration of less than 30 seconds), with relatively sufficient semantic objects as reference and multi-view observations for optimization, our algorithm achieves satisfactory results. However, it can be noted that the advantage of our method is not as obvious on outdoor datasets as on indoor datasets. The reasons are two-fold: One is that the camera moves along a straight path most of the time in outdoor scenes, in which case the quality of dual quadric mapping is not as good as in indoor scenes with the orbit trajectory. Second, only the “car” object can be used as reference for scale estimation in KITTI sequences, and there is no cross-validation with objects of other classes, which makes the overall result prone to bias.

TABLE I  
COMPARISON OF SCALE ESTIMATION ON TUM DATASET

| Seq        | Sucar’s [5] |      | Ours        |      |
|------------|-------------|------|-------------|------|
|            | RSE(%)      | std  | RSE(%)      | std  |
| fr1-desk   | 13.89       | 5.78 | <b>4.35</b> | 3.51 |
| fr2-desk   | 5.50        | 2.58 | <b>3.63</b> | 1.22 |
| fr3-office | 9.49        | 3.02 | <b>6.27</b> | 2.60 |

TABLE II  
COMPARISON OF SCALE ESTIMATION ON KITTI DATASET

| Seq  | Sucar’s [5]  |      | Ours        |      |
|------|--------------|------|-------------|------|
|      | RSE(%)       | std  | RSE(%)      | std  |
| 0000 | <b>11.09</b> | 2.02 | 17.58       | 2.76 |
| 0001 | 11.43        | 1.94 | <b>9.68</b> | 2.12 |
| 0007 | 8.50         | 1.76 | <b>7.77</b> | 2.19 |
| 0009 | 6.84         | 1.55 | <b>4.33</b> | 2.49 |
| 0011 | 9.64         | 1.88 | <b>9.05</b> | 2.24 |
| 0014 | <b>9.14</b>  | 2.02 | 13.44       | 2.12 |

### B. Odometry Evaluation Results

In order to further demonstrate the effect of scale estimation, we evaluated the scale-recovered camera pose error using absolute trajectory error (ATE) proposed in [27]. As a comparison, we also tested ORB-SLAM2 in monocular, stereo, and RGB-D modes. Note that no additional scale alignment was performed when calculating the ATE. The results are shown in Table III and Table IV. It can be seen that our method significantly reduces the error compared to the monocular mode of ORB-SLAM2. A centimeter-level accuracy is achieved in the TUM indoor scenes, and the average ATE is about 7 m in the KITTI outdoor scenes, which makes sense in some practical applications such as Augmented Reality and Virtual Reality. It can also be seen that the object size information is a weaker metric reference compared to the stereo baseline or RGB-D depth measurements due to its relatively large deviation.

### C. Qualitative Results

The qualitative evaluation results of scale estimation and semantic mapping are shown in Figure 4 and Figure 5 respectively. Figure 4 shows the comparison between the ground truth and the scale-recovered trajectory estimated by our system without further scale alignment. EVO [30] was used to visualize the trajectories. Figure 5(a) shows the projection of dual quadrics onto the images. It can be seen that the reconstructed dual quadrics can capture objects in the scene accurately. Figure 5(b) shows that the object-oriented

TABLE III  
COMPARISON OF POSE ERROR ON TUM DATASET

| Seq        | Ours         |       | ORB-SLAM2 (mono) |       | ORB-SLAM2 (RGB-D) |       |
|------------|--------------|-------|------------------|-------|-------------------|-------|
|            | ATE(m)       | std   | ATE(m)           | std   | ATE(m)            | std   |
| fr1-desk   | <u>0.045</u> | 0.031 | 0.062            | 0.061 | <b>0.016</b>      | 0.004 |
| fr2-desk   | <u>0.065</u> | 0.021 | 0.914            | 0.085 | <b>0.022</b>      | 0.006 |
| fr3-office | <u>0.139</u> | 0.054 | 1.210            | 0.019 | <b>0.010</b>      | 0.003 |

The **best** and second-best results are highlighted.

TABLE IV  
COMPARISON OF POSE ERROR ON KITTI DATASET

| Seq  | Ours         |      | ORB-SLAM2 (mono) |      | ORB-SLAM2 (stereo) |      |
|------|--------------|------|------------------|------|--------------------|------|
|      | ATE(m)       | std  | ATE(m)           | std  | ATE(m)             | std  |
| 0000 | 2.94         | 0.71 | 16.39            | 0.35 | <b>1.44</b>        | 0.01 |
| 0001 | <u>8.17</u>  | 3.03 | 79.51            | 0.75 | <b>3.03</b>        | 0.04 |
| 0007 | <u>7.54</u>  | 1.22 | 73.29            | 0.27 | <b>3.39</b>        | 0.04 |
| 0009 | <u>10.00</u> | 3.29 | 163.96           | 0.40 | <b>5.97</b>        | 0.16 |
| 0011 | <u>10.79</u> | 1.85 | 64.23            | 0.22 | <b>1.63</b>        | 0.05 |
| 0014 | <u>1.67</u>  | 0.27 | 11.89            | 0.28 | <b>0.39</b>        | 0.01 |

The **best** and second-best results are highlighted.

maps built by our system can express the environment well and are understandable with semantic information.

#### D. Ablation Study

As explained in Section III, three schemes are designed to enhance our method: outliers elimination, dimensions selection and uncertainty model. We conduct an ablation study to evaluate their effects on scale estimation. The results are shown in Table V. It can be seen that the relative scale errors are significantly diminished after the process of outliers rejection based on scale consistency. This is because we get rid of the objects with wrong semantic labels from the false 2-D object detections. Then, by introducing more reliable dimensions in the scale optimization, the accuracy and robustness of our method are further improved, especially in the KITTI sequences. Lastly, we adopt the uncertainty model to assign appropriate weights to each error term, and the best performance is achieved.

Furthermore, we believe that the number of dimensions of objects involved in scale estimation is closely related to the reliability of results. We run the proposed SLAM system on 7 sequences, recording the number of dimensions and the corresponding RSE. As shown in Figure 6, with the increasing number of dimensions, the results of scale estimation tend to be more stable and accurate.

#### E. Runtime Analysis

The experiments are carried out on Intel i7-9700K CPU with 3.6GHz and an Nvidia GeForce RTX 2080Ti with 11GB of memory. On all sequences, the maximum time per frame is around 22ms in the camera tracking thread. When processing keyframes, object detection and semantic mapping together take about 80ms per image. Scale estimation module takes less than 1ms. In summary, our system can run in real time.

## VI. CONCLUSION

In this paper, we present an accurate and robust scale estimation approach that takes the object size priors as the absolute reference. We develop a monocular object SLAM system to reconstruct objects as dual quadrics to provide reliable dimensions for scale estimation with no resort to assumptions on the gravity direction. In return, the estimated scale factor can be used to recover the absolute scale of the whole map built by our SLAM system. Quantitative and qualitative experiments demonstrate the outstanding

performance of our method. In the future, we intend to improve our algorithm to cope with the scale drift problem which often occurs in the large-scale outdoor environment. Moreover, methods of collecting and modeling prior object size information are also worth exploring.

## REFERENCES

- [1] L. Nicholson, M. Milford, and N. Sünderhauf, "QuadricSLAM: Dual Quadrics From Object Detections as Landmarks in Object-Oriented SLAM," *IEEE Robotics and Automation Letters*, vol. 4, no. 1, pp. 1–8, 2019.
- [2] S. Yang and S. Scherer, "CubeSLAM: Monocular 3-D Object SLAM," *IEEE Transactions on Robotics*, vol. 35, no. 4, pp. 925–938, 2019.
- [3] Y. Wu, Y. Zhang, D. Zhu, Y. Feng, S. Coleman, and D. Kerr, "EAO-SLAM: Monocular Semi-Dense Object SLAM Based on Ensemble Data Association," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2020, pp. 4966–4973.
- [4] T. Konkle and A. Oliva, "A Real-World Size Organization of Object Responses in Occipitotemporal Cortex," *Neuron*, vol. 74, no. 6, pp. 1114–1124, 2012.
- [5] E. Sucar and J.-B. Hayet, "Probabilistic Global Scale Estimation for MonoSLAM Based on Generic Object Detection," in *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*. IEEE, 2017, pp. 988–996.
- [6] E. Sucar and J.-B. Hayet, "Bayesian Scale Estimation for Monocular SLAM Based on Generic Object Detection for Correcting Scale Drift," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 5152–5158.
- [7] R. Hartley and A. Zisserman, *Multiple View Geometry in Computer Vision*. Cambridge, U.K.: Cambridge Univ. Press, 2004.
- [8] S. Chen, S. Song, J. Zhao, T. Feng, C. Ye, L. Xiong, and D. Li, "Robust Dual Quadric Initialization for Forward-Translating Camera Movements," *IEEE Robotics and Automation Letters*, vol. 6, no. 3, pp. 4712–4719, 2021.
- [9] S. Zhang, X. Li, Y. Liu, and H. Fu, "Scale-aware Insertion of Virtual Objects in Monocular Videos," in *2020 IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*. IEEE, 2020, pp. 36–44.
- [10] G. Nützi, S. Weiss, D. Scaramuzza, and R. Siegwart, "Fusion of IMU and Vision for Absolute Scale Estimation in Monocular SLAM," *Journal of Intelligent & Robotic Systems*, vol. 61, no. 1, pp. 287–299, 2011.
- [11] Z. Zhang, R. Zhao, E. Liu, K. Yan, and Y. Ma, "Scale Estimation and Correction of the Monocular Simultaneous Localization and Mapping (SLAM) Based on Fusion of 1D Laser Range Finder and Vision Data," *Sensors*, vol. 18, no. 6, p. 1948, 2018.
- [12] S. Song and M. Chandraker, "Robust Scale Estimation in Real-Time Monocular SFM for Autonomous Driving," in *2014 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2014, pp. 1566–1573.
- [13] Dingfu Zhou, Y. Dai, and Hongdong Li, "Reliable scale estimation and correction for monocular Visual Odometry," in *2016 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2016, pp. 490–495.
- [14] X. Wang, H. Zhang, X. Yin, M. Du, and Q. Chen, "Monocular Visual Odometry Scale Recovery Using Geometrical Constraint," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 988–995.
- [15] X. Yin, X. Wang, X. Du, and Q. Chen, "Scale Recovery for Monocular Visual Odometry Using Depth Estimated with Deep Convolutional Neural Fields," in *2017 IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2017, pp. 5871–5879.
- [16] N. Yang, R. Wang, J. Stückler, and D. Cremers, "Deep virtual stereo odometry: Leveraging deep depth prediction for monocular direct sparse odometry," in *Proceedings of the European Conference on Computer Vision (ECCV)*. Springer International Publishing, 2018, pp. 835–852.
- [17] R. Wang, M. Schworer, and D. Cremers, "Stereo DSO: Large-Scale Direct Sparse Visual Odometry With Stereo Cameras," in *Proceedings of the IEEE International Conference on Computer Vision*. IEEE, 2017, pp. 3903–3911.
- [18] W. N. Greene and N. Roy, "Metrically-Scaled Monocular SLAM using Learned Scale Factors," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 43–50.



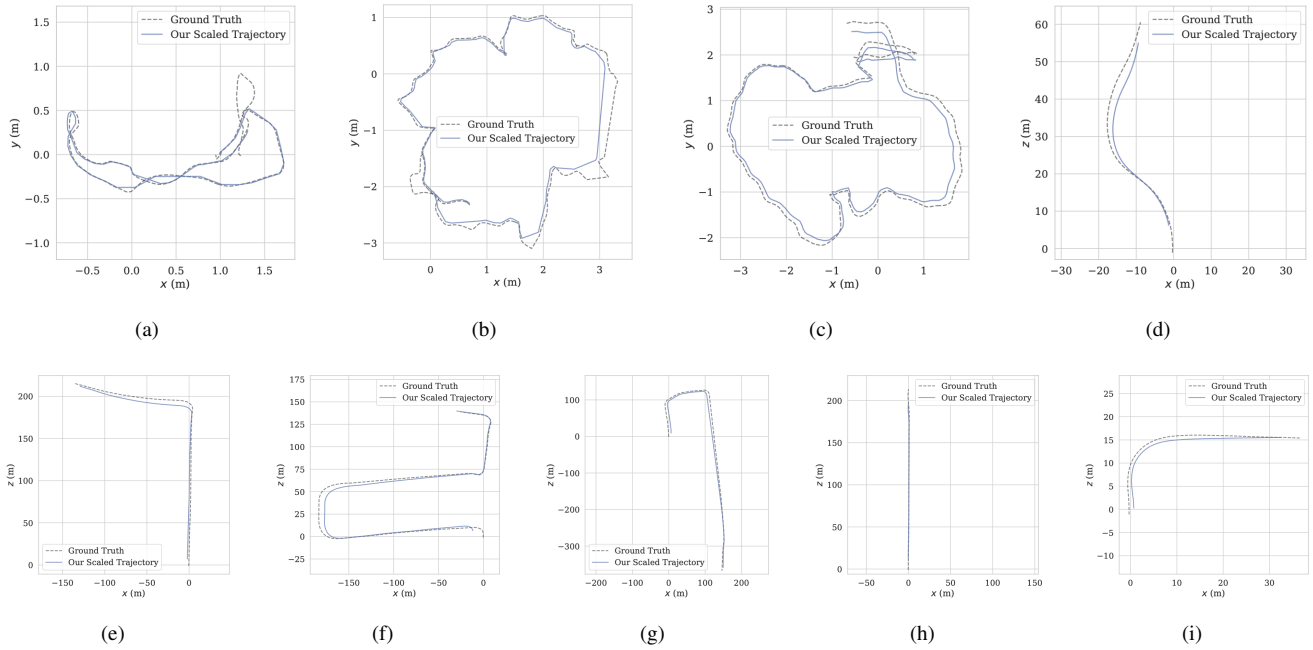


Fig. 4. Comparisons between the ground truth (gray) and our scaled trajectories (blue). (a), (b) and (c) are results on TUM fr1-desk, fr2-desk and fr3-desk, respectively. (d)-(i) are results on KITTI tracking sequences. (d) 0000, (e) 0001, (f) 0007, (g) 0009, (h) 0011 and (i) 0014.

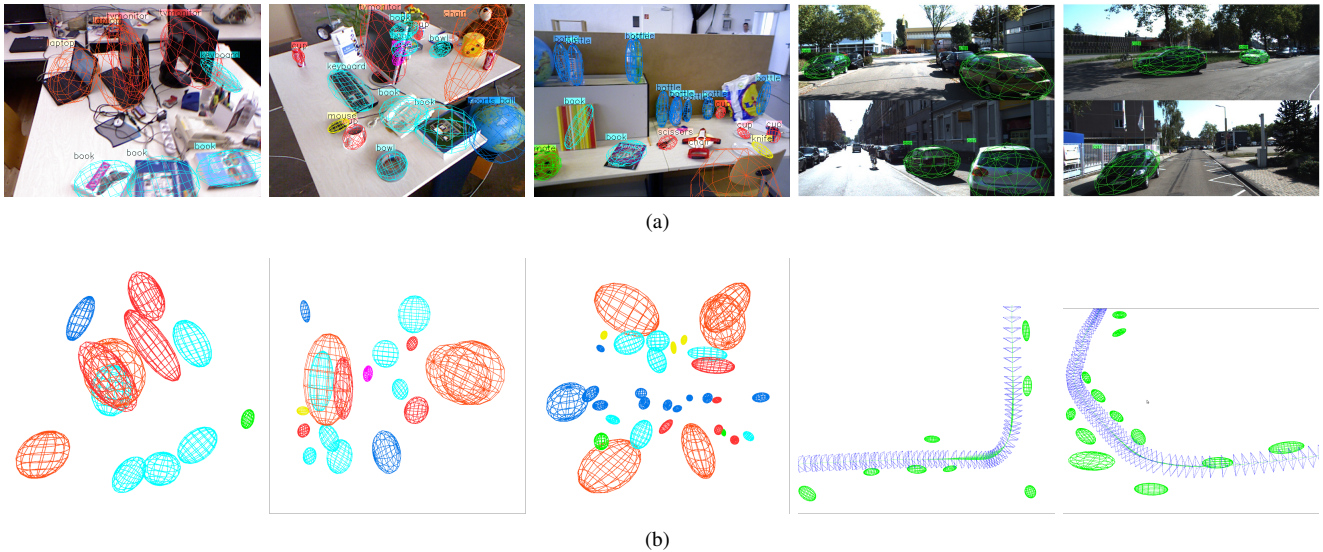


Fig. 5. The qualitative results of scale estimation and semantic mapping. From left to right, columns are results on sequence TUM fr1-desk, fr2-desk and fr3-desk, KITTI tracking 0001 and 0009 respectively. (a) The 2-D projection of reconstructed quadrics on the image. (b) The object map built by our object SLAM.

TABLE V  
ABLATION STUDY OF DIFFERENT OPTIONAL SELECTING RESULTS

| out <sup>1</sup> dim <sup>2</sup> uncer <sup>3</sup> | fr1-desk    |             | fr2-desk    |             | fr3-office  |             | 0000         |             | 0001        |             | 0007        |             | 0009        |             | 0011        |             | 0014         |             |
|--|-------------|-------------|-------------|-------------|-------------|-------------|--------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|--------------|-------------|
|  | RSE(%)      | std         | RSE(%)      | std         | RSE(%)      | std         | RSE(%)       | std         | RSE(%)      | std         | RSE(%)      | std         | RSE(%)      | std         | RSE(%)      | std         | RSE(%)       | std         |
|  | 10.81       | 6.56        | 9.10        | 3.40        | 14.51       | 4.59        | 29.42        | 4.78        | 15.25       | 2.83        | 22.56       | <b>1.89</b> | 10.87       | 4.56        | 19.89       | 3.14        | 22.36        | 4.53        |
| ✓  | 4.84        | 3.37        | 5.98        | 1.94        | 7.88        | 3.70        | —            | —           | 13.48       | 3.85        | 19.42       | 3.90        | 4.96        | 3.56        | 18.78       | <b>2.00</b> | —            | —           |
| ✓ ✓  | 4.48        | <b>3.37</b> | 3.79        | 1.28        | 6.44        | 2.91        | 17.65        | 2.83        | 10.54       | 2.71        | 8.38        | 2.87        | 4.39        | 2.58        | 9.67        | 2.97        | <b>13.42</b> | <b>2.11</b> |
| ✓ ✓ ✓  | <b>4.35</b> | 3.51        | <b>3.63</b> | <b>1.22</b> | <b>6.27</b> | <b>2.60</b> | <b>17.28</b> | <b>2.76</b> | <b>9.68</b> | <b>2.12</b> | <b>7.77</b> | 2.19        | <b>4.33</b> | <b>2.49</b> | <b>9.05</b> | 2.24        | 13.44        | 2.12        |

<sup>1</sup> Outliers elimination. <sup>2</sup> Dimensions selection. <sup>3</sup> Uncertainty model. “—” means that outliers elimination does not work in sequence 0000 and 0014 because too few dimensions are available before dimensions selection.

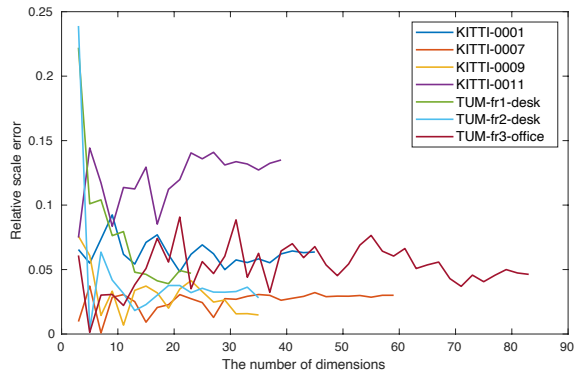


Fig. 6. The relation between scale accuracy and the number of available dimensions of objects used in scale estimation.

- [19] 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 *2013 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2013, pp. 1352–1359.
- [20] D. Gálvez-López, M. Salas, J. D. Tardós, and J. M. M. Montiel, "Real-time monocular object SLAM," *Robotics and Autonomous Systems*, vol. 75, pp. 435–449, 2016.
- [21] D. Frost, V. Prisacariu, and D. Murray, "Recovering Stable Scale in Monocular SLAM Using Object-Supplemented Bundle Adjustment," *IEEE Transactions on Robotics*, vol. 34, no. 3, pp. 736–747, 2018.
- [22] M. Weinmann, B. Jutzi, and C. Mallet, "Semantic 3D scene interpretation: A framework combining optimal neighborhood size selection with relevant features," *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, vol. II-3, pp. 181–188, 2014.
- [23] R. Kümmerle, G. Grisetti, H. Strasdat, K. Konolige, and W. Burgard, "G2o: A general framework for graph optimization," in *2011 IEEE International Conference on Robotics and Automation*. IEEE, 2011, pp. 3607–3613.
- [24] R. Mur-Artal and J. D. Tardós, "ORB-SLAM2: An Open-Source SLAM System for Monocular, Stereo, and RGB-D Cameras," *IEEE Transactions on Robotics*, vol. 33, no. 5, pp. 1255–1262, 2017.
- [25] S. A. Gottschalk, "Collision queries using oriented bounding boxes," Ph.D., The University of North Carolina at Chapel Hill, 2000.
- [26] K. Ok, K. Liu, K. Frey, J. P. How, and N. Roy, "Robust Object-based SLAM for High-speed Autonomous Navigation," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 669–675.
- [27] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers, "A benchmark for the evaluation of RGB-D SLAM systems," in *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2012, pp. 573–580.
- [28] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? The KITTI vision benchmark suite," in *2012 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2012, pp. 3354–3361.
- [29] J. Redmon and A. Farhadi, "YOLOv3: An Incremental Improvement," *arXiv e-prints*, 2018.
- [30] M. Grupp, "evo: Python package for the evaluation of odometry and slam." <https://github.com/MichaelGrupp/evo>, 2017.