

SemanticLoop: loop closure with 3D semantic graph matching

Junfeng Yu, Shaojie Shen

Abstract—Loop closure can effectively correct the accumulated error in robot localization, which plays a critical role in the long-term navigation of the robot. Traditional appearance-based methods rely on local features and are prone to failure in ambiguous environments. On the other hand, object recognition can infer objects’ category, pose, and extent. These objects can serve as stable semantic landmarks for viewpoint-independent and non-ambiguous loop closure. However, there is a critical object-level data association problem due to the lack of efficient and robust algorithms.

We introduce a novel object-level data association algorithm, which incorporates IoU, instance-level embedding, and detection uncertainty, formulated as a linear assignment problem. Then, we model the objects as TSDF volumes and represent the environment as a 3D graph with semantics and topology. Next, we propose a graph matching-based loop detection based on the reconstructed 3D semantic graphs and correct the accumulated error by aligning the matched objects. Finally, we refine the object poses and camera trajectory in an object-level pose graph optimization.

Experimental results show that the proposed object-level data association method significantly outperforms the commonly used nearest neighbor method in accuracy. Our graph matching-based loop closure is more robust to environmental appearance changes than existing appearance-based methods.

I. INTRODUCTION

The long-term autonomous navigation of mobile robots is critical for many applications (e.g., self-driving cars and service robots). However, accumulated errors will inevitably occur in robot localization due to sensor noise. In order to correct the accumulated drift, robots need to perceive the environment in real-time and recognize previously visited places (i.e., loop closure). Although loop closure has been studied extensively in Simultaneous Localization and Mapping (SLAM), it is still considered a well-defined but highly challenging problem to solve in the general sense.

Classical appearance-based methods typically reformulate loop closure as an image retrieval problem. They represent the environment as a database of images. Then the current image is matched with the ones in the database to retrieve the most similar candidate(s) in appearance. These methods generally use visual descriptors to represent images for more efficient retrieval. The Bag-of-Words (BoW [1]) extracted from local features (e.g., ORB [2]) is one of the most effective models. Many existing SLAM systems (e.g., ORB-SLAM2 [3], VINS-Mono [4]) used BoW and demonstrated impressive performance. These approaches are flexible and general. However, they still face many challenges. For example, when the appearance changes due to lighting or

viewpoint differences, the local features may change dramatically, and the classical methods fail. Moreover, appearance-based methods tend to ignore the geometric structure of the environment, which may lead to false positives in repetitive environments.

On the other hand, semantics and geometric structures are usually invariant to appearance changes. For example, a chair remains a chair, whether observed during the day or night or from different viewpoints. Recently, deep learning has made significant progress on perceptual tasks such as object detection and instance segmentation (e.g., Mask R-CNN [5]), motivating the incorporation of semantics into SLAM systems to improve the localization accuracy (e.g., SLAM++ [6], Fusion++ [7]). However, due to the generalization problem, deep learning models often suffer from noise (e.g., false detections, misclassifications) in the working environment. This perceptual noise can easily lead to incorrect object-level data associations, which introduces erroneous semantics. Although this inaccurate semantic information can seriously affect the accuracy and robustness of localization, existing works tend to ignore this critical problem.

Contributions: To address the above challenges, we propose an RGBD-based semantic mapping system with loop closure. Specifically, our main contributions are as follows:

- We introduce a novel object-level data association method that combines IoU, instance-level embedding, and detection uncertainty into a linear assignment formulation, constructing an accurate 3D semantic map insensitive to the noises from deep learning models and odometry drift.
- We propose a 3D semantic graph matching-based loop closure approach that couples semantics and topology of the instances in a quadratic assignment formulation, making the loop closure more robust to appearance changes in the scene.
- To maintain a globally consistent map, we introduce an object-level pose graph optimization that includes the odometry and loop closure constraints to optimize the camera trajectories and object poses jointly.

Moreover, we evaluate the proposed methods on the public TUM RGBD benchmark [8] and SceneNN dataset [9] to verify their effectiveness.

II. RELATED WORK

A. Data Association in Semantic SLAM

SLAM++ [6] is a pioneering work in the direction of semantic SLAM, which first used real-world objects (e.g., tables, chairs) as landmarks. The data association relies on

*The authors are with the Department of Electronic and Computer Engineering, Hong Kong University of Science and Technology, Hong Kong, China. (e-mail: jyubm@connect.ust.hk; eeshaojie@ust.hk).

the Point-Pair feature (PPF)-based 3D object recognition. However, it required a pre-built database of CAD models, making the system less universal. In their later work, Fusion++ [7] utilized a reconstruction-by-segmentation strategy to build a TSDF volume for each object, which solved the problem of relying on an offline CAD database. The data association depended on the IoU between the mask from instance segmentation and the mask projected from the TSDF volume. However, the data associations can become ambiguous when odometry drifts or objects are occluded.

Another line of study is the probabilistic data association. The core idea is to use 'soft' instead of 'hard' data association to put the data association uncertainty into the SLAM backend. The most representative work is Gaussian PDA [10], which proposed using the Expectation-Maximization (EM) algorithm to solve this discrete-continuous optimization problem. The EM algorithm needs to be solved iteratively. However, recalculating the combinatorial number of historical data associations is infeasible for computational reasons. Semantic MM [11] advanced this stream by approximating the data associations with a Max-Marginalization (MM) technique, which solved the computational complexity problem by assuming that future observations will not affect past data associations. However, not optimizing past data association weights may result in a low probability of getting the correct data association, especially when there are many ambiguities due to odometry drift.

Recently, QuadricSLAM [12] and CubeSLAM [13] explored the use of ellipsoids and cuboids as object representations, which were extracted from multi-view and single-view 2D bounding boxes, respectively. In QuadricSLAM, they overlooked the data association problem and focused on ellipsoid initialization. In the following work, [14] used the BoW model as object representation and formulated the data association as a linear assignment problem. In CubeSLAM, the data association relied on feature point matching, and the bounding box that shared the most feature points was selected. However, data association in these works depended on traditional feature points or descriptors without semantics, which may be subject to failure on textureless objects (e.g., TV).

Similar to Fusion++ [7], we use a reconstruction-by-segmentation strategy. The dense object model can represent objects' pose and shape while providing sufficient semantics. In contrast to the commonly used nearest neighbor method, our data association method combines IoU, instance-level embedding, and detection uncertainty into a linear assignment formulation. Therefore, it is more robust to textureless objects, deep learning model noises, and odometry drift.

B. Loop Closure in Semantic SLAM

SLAM++ [6] performed loop closure by matching the local object graph with the long-term object graph. They treated objects as vertices and their x-axes as normal directions to extract the PPFs and then reused the same 3D object recognition algorithm as in data association. In Fusion++ [7],

they extracted 3D BRISK for object models and applied the 3D-3D RANSAC algorithm between them to perform loop detection, which was extremely slow (more than 780ms) even on modern GPU platforms.

Recent approaches attempted to incorporate more semantics to address the loop closure problem in cases with extreme appearance changes. X-View [15] proposed a novel loop detection idea based on semantic graphs. The system constructed 2D semantic graphs using image sequences with instance segmentation, in which vertices were semantic blobs and edges represented proximity relations. Loop closure depended on matching the random walk descriptors between vertices. The random walk descriptor contained topological information of the semantic graph, making it highly robust to seasonal and significant viewpoint changes. A series of follow-up works had extended the idea in X-View to 3D [16], to edit distance minimization-based matching algorithm [17], and to semantic histogram-based descriptor [18] (to be faster). However, when there are many objects in the graph, random walks tend to lose information, and the performance may degrade severely. Graph matching in [17] often suffers false alarms when duplicated objects with similar topologies are in the graph. Moreover, matching descriptors between the query and target graphs is inefficient when the number of random walks is large.

In our system, we perform loop closure by constructing 3D semantic maps online. To improve the efficiency and robustness of loop closure detection, we perform a geometric graph matching between semantic graphs rather than descriptors matching between vertices. In addition, our system is a complete pipeline, including a pose graph optimization to maintain the camera and object poses.

III. METHOD

A. Overview

Figure 1 visualizes the pipeline proposed in our work. There are mainly three modules: semantic mapping, graph-based loop closure, and pose graph optimization.

B. Semantic Mapping

From RGBD input, we utilize off-the-shelf RGBD odometry to obtain relative poses. The underlying assumption is that we can build a 3D object map based on the local consistency of the odometry and our data association algorithm. An instance segmentation network processes the RGB frame in a separate thread to detect bounding boxes, masks, and semantic labels. Then the bounding boxes are fed into an instance-feature learning network to extract the instance-level embeddings. Based on the outputs of the neural network thread, the multi-frame instance tracking algorithm filter out perceptual noises and integrate a local map at the current position. Then, the object-level data association algorithm matches detections in the local map with the objects in the global map based on semantics and camera poses. When no match occurs, we create a new TSDF volume and add it to the global map. When an object is associated, we utilize an approach similar to Fusion++ [7] to fuse the new

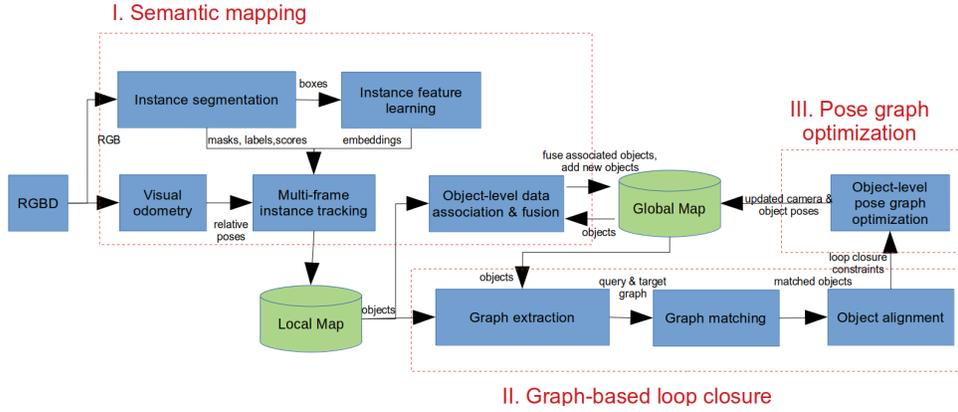


Fig. 1. Overview of the proposed semantic mapping system with loop closure.

measurements into the TSDF volume and use an averaging scheme to update a probability distribution over the semantic label.

Object-Level Data Association: The core of semantic mapping is a critical object-level data association problem. Suppose at frame i , there are M detections from the instance segmentation network, denoted as $\mathcal{S} \triangleq \{s_k\}_{k=1}^M$. Each detection is represented as $s_k = (m_k, l_k, c_k, e_k)$, where m_k is the binary mask, l_k is the semantic label, c_k is the confidence score, and e_k is the embedding from the instance feature learning network. Meanwhile, we have N object landmarks in the object map, denoted as $\mathcal{O} \triangleq \{o_j\}_{j=1}^N$. Each object is represented as $o_j = (V_j, T_{wo_j}, l_j, m_j, E_j)$, where V_j is the TSDF volume, T_{wo_j} is the pose, l_j is the semantic label, m_j is the predicted binary mask, and E_j is a set storing all the matched embeddings from past matches.

The object-level data association needs to find as many matches as possible by assigning at most one object to each detection and at most one detection to each object, such that the total matching cost is minimized. Since we usually have more landmarks than detections, i.e., $N \geq M$, we can reformulate this problem as a 2D rectangular assignment problem as follows:

$$\begin{aligned}
 & \min_{\mathbf{A}} \sum_{j=1}^N \sum_{k=1}^M \text{Tr}(\mathbf{A}^T \mathbf{L}) \\
 & \text{subject to } \mathbf{A}(j, k) \in \{0, 1\}, \forall j, k \\
 & \sum_{j=1}^N \mathbf{A}(j, k) = 1, \forall k \\
 & \sum_{k=1}^M \mathbf{A}(j, k) \leq 1, \forall j
 \end{aligned} \quad (1)$$

where \mathbf{A} is a $N \times M$ assignment matrix, and \mathbf{L} is a $N \times M$ cost matrix. The equality constraint means that every column (detection) is assigned to a row (landmark). The inequality constraint means that not every row (landmark) is assigned to a column (detection). Note that due to unobserved new objects or false detections, the matching cost between a

detection and a landmark may surpass a certain threshold, then this association should be discarded.

The assignment matrix \mathbf{A} can be defined as follows:

$$\mathbf{A}(j, k) = \begin{cases} 1, & \text{if } o_j \text{ is matched with } s_k \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

A matching cost between object o_j and detection s_k can be calculated based on the binary masks, embeddings, and semantic labels as follows:

$$\begin{aligned}
 \mathbf{L}(j, k) &= 1.0 - \mathbf{W}(j, k) p(l_k | l_j) \\
 \mathbf{W}(j, k) &= \lambda \mathbf{IoU}(j, k) + (1.0 - \lambda) \mathbf{d}_f(j, k) \\
 \mathbf{IoU}(m_j, m_k) &= \frac{\sum m_k \cap m_j}{\sum m_k + \sum m_j - \sum m_k \cap m_j} \\
 \mathbf{d}_f(e_j, e_k) &= e_j e_k
 \end{aligned} \quad (3)$$

where \mathbf{IoU} is calculated between the mask m_k and the predicted mask m_j . The metric distance (\mathbf{d}_f) of embeddings is computed between the instance embedding e_k and every embedding $e_j \in E_j$. We use the cosine distance and choose the maximum among all $\mathbf{d}_f(e_j, e_k)$. A hyperparameter λ is used to balance the \mathbf{IoU} and \mathbf{d}_f terms. The probability distribution $p(l_k | l_j)$ corresponds to the confusion matrix of the instance segmentation network and is learned offline. The problem defined in equation (1) can be solved using the shortest augmenting path algorithm, described in [19].

Multi-Frame Instance Tracking: Since noises of the deep learning model are likely to lead to erroneous data associations, we perform a multi-frame instance tracking to filter out the perceptual noises. We reuse the same formulation as in object-level data association, and the main difference is to use the mask in the previous frame instead of the predicted mask. We only keep the instances that are tracked over a certain number of times.

C. Graph-based Loop Closure

Directly matching between vertices through random walk descriptions is inefficient and does not fully exploit the topology in the graph. First, we extract the query and target graphs from the local and global maps. Then, we perform a geometric graph matching to find the correspondences

between the query and target graphs by considering both edge and vertex similarity. Next, we estimate the drift errors by aligning the matched objects. See an example in Figure 2.

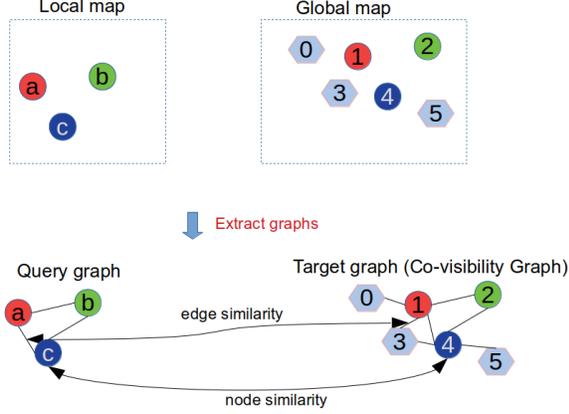


Fig. 2. An example shows the graph extraction and graph matching steps. The object in the map forms a graph, where the vertex is composed of the object’s center and semantics. The edge represents the metric distance and co-visibility relationship between the objects. Graph matching aims to find the correspondences between the query and target graphs by considering both edge and vertex similarity.

Graph Extraction: All objects in the map form a graph $G = (V, E)$, where vertices $v_j \in V$ contains the center, semantic label l_j and embeddings of object o_j . Each edge $e_{j_1, j_2} \in E$ is given by the Euclidean distance between the centers of object o_{j_1} and object o_{j_2} . In order to maintain the topology between the objects, we adopt a co-visibility strategy. That is, we add an edge between the nodes only when the corresponding objects are observed in the same local map. In this way, we extract the query G_q and target G_t graphs from the local and global maps.

Graph Matching: With a query graph $G_q = (V_q, E_q)$ of size M and a target graph $G_t = (V_t, E_t)$ of size N , the graph matching step aims to find a correspondence between graphs, which fits both vertex’s attributes (e.g., semantic label, embeddings) and graph topology (e.g., Euclidean distance between the co-visible objects). The problem can be reformulated as a quadratic assignment problem as follows:

$$\begin{aligned}
 & \max_{\mathbf{A}} \sum_{e_{j_1, j_2} \in E_q} \sum_{e_{k_1, k_2} \in E_t} (\mathbf{A}(j_1, k_1) \mathbf{A}(j_2, k_2) \\
 & \quad \mathbf{L}(j_1, j_2, k_1, k_2)) \\
 & = \max_{\mathbf{A}} \mathbf{vec}(\mathbf{A})^\top \mathbf{S} \mathbf{vec}(\mathbf{A}) \\
 & \text{subject to } \mathbf{A}(j, k) \in \{0, 1\}, \forall j, k \\
 & \quad \sum_{j=1}^N \mathbf{A}(j, k) \leq 1, \forall k \\
 & \quad \sum_{k=1}^M \mathbf{A}(j, k) \leq 1, \forall j
 \end{aligned} \tag{4}$$

where \mathbf{A} is a $N \times M$ assignment matrix, and \mathbf{L} is a $N \times M \times N \times M$ reward tensor, The \mathbf{vec} operator vectorizes a

matrix into a column vector. The reward matrix \mathbf{S} is a square matrix of size NM , which is constructed by unfolding the reward tensor L . The diagonal elements of the reward matrix are the matching reward for the nodes, and the off-diagonal elements are the matching reward for the edges.

The assignment matrix \mathbf{A} can be defined as follows:

$$\mathbf{A}(j, k) = \begin{cases} 1, & \text{if } v_j \in V_t \text{ is matched with } v_k \in V_q \\ 0, & \text{otherwise} \end{cases} \tag{5}$$

A matching reward can be defined as follows:

$$\begin{aligned}
 \mathbf{L}(j_1, j_2, k_1, k_2) &= \mathbf{d}_v(j_1, k_1) \mathbf{d}_v(j_2, k_2) \mathbf{d}_e(e_{j_1, j_2}, e_{k_1, k_2}) \\
 \mathbf{d}_v(j, k) &= \begin{cases} 1, & \text{if } l_j = l_k \\ 0, & \text{otherwise} \end{cases} \\
 \mathbf{d}_e(e_{j_1, j_2}, e_{k_1, k_2}) &= \exp(-\mu \|e_{j_1, j_2} - e_{k_1, k_2}\|_2)
 \end{aligned} \tag{6}$$

where functions \mathbf{d}_v and \mathbf{d}_e calculate the similarity of vertices and edges respectively, and μ is a hyperparameter.

The problem defined in equation (4) is NP-hard. However, we can solve it approximately using the spectral methods described in [20]. We first relax the integral constraints on A , such that the elements of A can take real values between $[0, 1]$. Since only the relative values between the elements of A matter, we can fix the norm of $\mathbf{vec}(A)$ to 1. According to the Raleigh’s ratio theorem, the $\mathbf{vec}(A^*)$ that maximizes $\mathbf{vec}(A)^T \mathbf{S} \mathbf{vec}(A)$ is the principal eigenvector of \mathbf{S} . Moreover, a key constraint of \mathbf{S} is that it is element-wise non-negative. Therefore, by the Perron-Frobenius theorem, the elements of $\mathbf{vec}(A^*)$ are non-negative, i.e., between $[0, 1]$. Next, in order to obtain an assignment matrix from A^* , i.e., a matrix with elements in $\{0, 1\}$ and proper row/column sums, [20] proposed to use a greedy algorithm to discretize the $\mathbf{vec}(A^*)$. Our main difference from [20] is that we reuse the linear assignment formulation in equation (1) with a cost matrix A^* to obtain an assignment matrix.

Object Alignment: We estimate the relative transformations by registering the point clouds extracted from the TSDFs of the matched objects. In order to get an accurate result for this wide baseline alignment, we first use the PFPF [21]-based 3D-3D RANSAC to perform an initial coarse alignment, then use ICP for refinement.

D. Pose Graph Optimization

The pose graph contains both object and camera nodes. Each node contains an $\mathbf{SE}(3)$ transformation. For frame i with instance segmentation, we create a new camera node T_{wi} . We fix the first camera node as the world coordinate system. When a new object o_j is added, we create a corresponding object node T_{wo_j} . The object’s coordinate system is attached to the object’s center, and the coordinate axes are aligned with the world coordinate axes. Each $\mathbf{SE}(3)$ measurement is a relative transformation constraint between the corresponding nodes. The measurement $Z_{i, i+1}$ between camera nodes represents the relative pose estimate from frame i to $i+1$. The measurement $Z_{o_j, i}$ between the object and camera nodes denotes the pose of frame i expressed

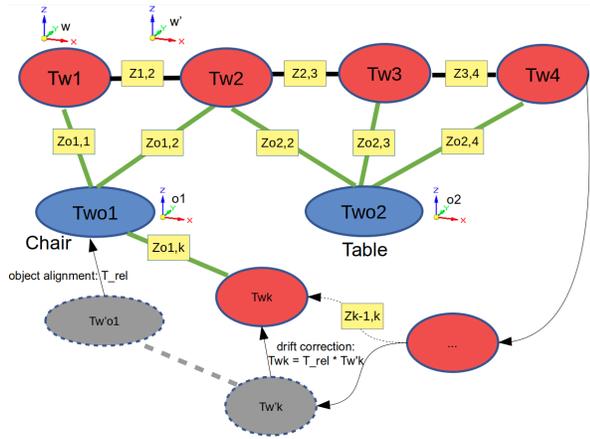


Fig. 3. An example pose graph shows six camera pose nodes T_{w_i} (red circles) and two object pose nodes $T_{w_{o_j}}$ (blue circles). The object-camera constraints $Z_{o_j,i}$ (green edges) represent the pose of frame i in the coordinate system of the object o_j . The camera-camera constraints $Z_{i,i+1}$ (black edges) denote the odometry measurements between frame i and $i+1$. When a loop is detected, the drift error T_{rel} is obtained by aligning the object in the local map (node $T_{w'_{o_1}}$) with the object in the global map (node $T_{w_{o_1}}$). Then, the camera pose T_{w_k} and object-camera constraint $Z_{o_1,k}$ are corrected using the estimated drift error.

in the object o_j 's coordinate system. See an example in Figure 3.

Object-Level Pose Graph Optimization: After getting the drift errors through object alignment, we can calculate the corrected camera pose and add the new object-camera constraints to the pose graph. Then we can further refine the entire camera trajectory $\mathcal{X} = \{T_{w,i}\}_{i=1}^T$ and object poses $\mathcal{O} = \{T_{w,o_j}\}_{j=1}^M$ in a object-level pose graph optimization. We minimize the error terms for all measurement constraints as follows:

$$\begin{aligned} \mathcal{X}, \mathcal{O} = \arg \min_{\mathcal{X}, \mathcal{O}} & \sum_{Z_{i,i+1}} \|e_{cc}(T_{w,i}, T_{w,i+1})\|_{\Sigma_{t,t+1}} \\ & + \sum_{Z_{o_j,i}} \|e_{oc}(T_{w,i}, T_{w,o_j})\|_{\Sigma_{o_j,i}} \quad (7) \\ e_{cc}(T_{w,i}, T_{w,i+1}) & = \log(Z_{i,i+1}^{-1} T_{w,i}^{-1} T_{w,i+1}) \\ e_{oc}(T_{w,i}, T_{w,o_j}) & = \log(Z_{o_j,i}^{-1} T_{w,o_j}^{-1} T_{w,i}) \end{aligned}$$

where e_{cc} and e_{oc} are the measurement error terms for the camera-camera $Z_{i,i+1}$ and object-camera measurement constraint $Z_{o_j,i}$ respectively. $\|e\|_{\Sigma} = e^T \Sigma^{-1} e$ is the Mahalanobis distance and \log is the logarithmic map of $\mathbf{SE}(3)$. We solve this nonlinear least squares problem using the Levenberg-Marquart algorithm in the Ceres solver. After optimization, we update the object and camera poses before initializing new objects.

IV. EXPERIMENTS

We use the open-source Mask R-CNN implementation of Matterport as the instance segmentation network, and the weights are pre-trained on the Microsoft COCO dataset [22]. We use [23] as the instance-level feature learning network, and the weights are fine-tuned on the SceneNet RGBD dataset [24].

We conduct experiments on the public TUM RGBD and SceneNN dataset to evaluate the performance of the proposed methods. The TUM RGBD dataset [8] consists of real-time RGB, depth images, and ground-truth trajectories. In addition to RGB, depth, and camera pose ground-truth, the SceneNN dataset [9] provides instance segmentation ground-truth.

A. Data Association Performance

Metric: The idea is to find the correspondences between ground-truth object IDs and object IDs in the data association algorithm. Similar to [14], we reformulate this problem as a linear assignment problem. The first partite set S consists of the object IDs in the data association algorithm. The second partite set F consists of the ground-truth object IDs. The reward function $w(j, k)$ on edge $e(j, k) \in F \times S$ can be defined as the number of identical bounding boxes shared by the object ID $j \in F$ and object ID $k \in S$. By solving this problem, the sum of the reward on all matched edges is the number of correct associations. We take the ratio of the number of correct associations to the number of all ground-truth bounding boxes as the accuracy of the data association algorithm.

We compare the proposed method with the commonly used nearest neighbor method in semantic SLAM systems. In our experiments, we add random noises to the ground-truth camera poses and observe how the accuracy changes with different noise levels. We conduct experiments on SceneNN 021, 025, and 231 sequences, as shown in Figure 4. Due to the inaccuracy of the camera poses, it is easy to cause the prediction to deviate from the measurement, resulting in ambiguities in the data association. Results show that the proposed method is more robust to localization noise than the baseline method due to incorporating instance-level embeddings and performing global minimum cost matching.

B. Loop Detection Performance

We label two frames as a loop closure if they observe more than two objects in common. To prevent adjacent frames from being labeled as loop closures, the difference in frame indices between them needs to be greater than 500. If 50% of the edges in the local graph can be matched, then we consider the loop detection successful. We compare the proposed loop detect method with ORB-SLAM2 and the random walk descriptor-based graph matching in [17]. The authors of [17] have not released its source code before the submission. To have a fair comparison, we faithfully reimplemented [17] on our own. In our experiment, we set the number of random walks to 200 and the walk depth to 4. The other parameters are consistent with the paper. The element of the random walk descriptor consists of the semantic label and embedding of the object. Table I shows the results of loop closure detection on TUM RGBD sequences. The results show that, compared with ORB-SLAM2, our method can achieve 100% accuracy and yield more true positives on these three sequences, although it detects fewer loop closure candidates due to its stricter graph matching. Compared

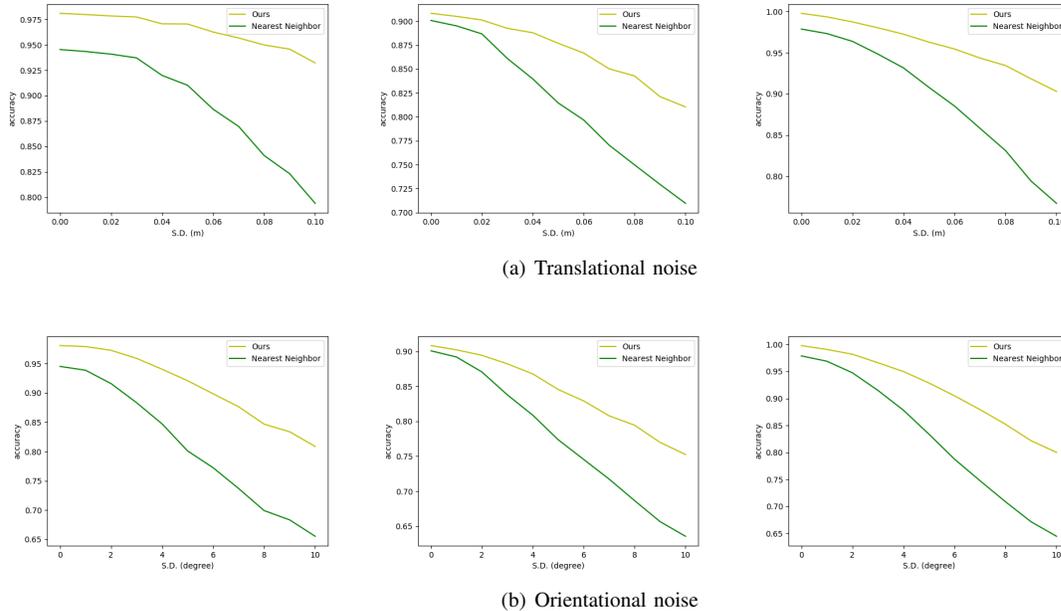


Fig. 4. Results of object-level data association accuracy according to the different noise levels on the public SceneNN dataset. The three columns are 021, 025, and 231 sequences. (a) accuracy at different translational noise levels. (b) accuracy at different orientational noise levels.

to the random walk descriptor-based method, our method achieves a higher recall and accuracy.

TABLE I. Loop detection results on TUM RGBD dataset.

Sequence	Metric	Ours	ORB-SLAM2	[17]
fr1 room	Detections	21	176	16
	True Positives	21	26	16
	False Positives	0	150	0
	After Verification	21	0	16
fr2 desk	Detections	44	188	35
	True Positives	44	38	35
	False Positives	0	150	0
	After Verification	44	1	35
fr3 office	Detections	76	250	59
	True Positives	76	34	51
	False Positives	0	216	8
	After Verification	76	1	51

C. Loop Closure Results

Figure 5 shows two challenging scenes on the TUM RGBD dataset. The first column is the result of Superglue [25] feature matching. The second column is the result of semantic graph matching. Results show that the state-of-the-art learning-based matching method can not effectively perform loop closures in these challenging cases. However, our semantic graph matching-based method can associate measurements from different viewpoints to the object landmarks in the map and thus is highly robust to significant viewpoint differences.

Figure 6 shows four false loop detections on the TUM RGBD and SceneNN datasets. The first row is the failure cases of random walk descriptor-based graph matching in [17], but successful based on our method. The reasons are two folded: firstly, random walks tend to lose information

when there are many objects in the global map, and secondly, random walk descriptor-based graph matching often suffers false alarms when duplicated objects with similar topologies are in the global map. The above two reasons explain why Table I shows that our method can achieve better recall and accuracy than the random walk descriptor-based method. The second row is the failure cases of our method on the SceneNN dataset. Results show that partially reconstructed objects with inaccurate centers may cause false loop detections. Moreover, our method cannot handle scenes with the same object layout, e.g., a computer lab with many repeated object layouts. Combining feature points and semantics may solve this problem. However, if there are multiple associated objects, we can eliminate these false loop detections by checking for topology and the spatial distance consistency between the matching objects in the two maps.

D. Localization Performance

We evaluate the localization performance on the TUM RGBD and SceneNN datasets. We compare our approach with ORB-SLAM2. Table II shows the Root Mean Square translational Error (RMSE) of the trajectories. Note that we obtained the results of ORB-SLAM2 with the loop closing thread turned on. On fr2_desk and fr3_office sequences, our method is on par with ORB-SLAM2 since both methods detect enough loop closures. However, our method exhibits better localization performance on the other sequences as it can detect more challenging loop closures.

E. Runtime and Scalability

We evaluate the average running time of graph match on a Linux system with an Intel Core i7-7700K CPU at 4.20GHz. Table III shows that our method is more than 2.5 times

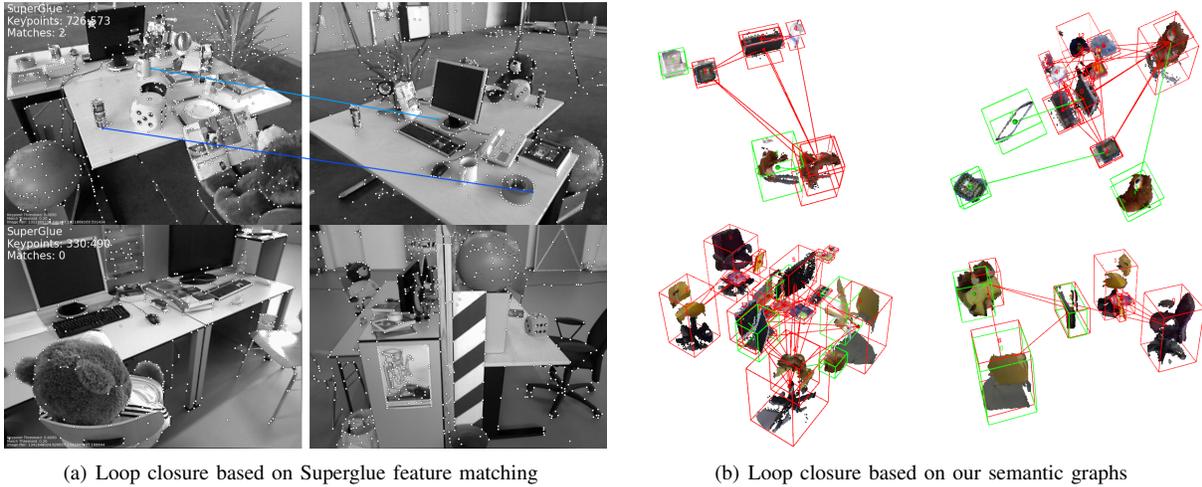


Fig. 5. Two examples of challenging cases (significant viewpoint differences) in loop closure on the public TUM RGBD dataset. (a) Failure matching results based on Superglue. (b) Successful matching results based on our semantic graph matching. Green bounding boxes are objects in the local map. Red bounding boxes are objects in the global map. Green lines represent the correspondences of the objects. Red lines represent the co-visibility relationships.



(a) Failure cases of random walk descriptor-based graph matching in [17]. (Left) Random walks tend to lose information when there are many objects in the global map. As a result, two books (highlighted in the red circle) in the local map failed to match. (Right) Random walk descriptor-based graph matching often suffers false alarms when duplicated objects with similar topologies are in the global map. The monitor in the local map is incorrectly associated (highlighted with the green arrow) with the other monitor with id=5 in the global map.

(b) Failure cases of our method. (Left) The partially reconstructed desk is incorrectly associated (highlighted with the red arrow) because its center is closer to the desk with id=304 than the desk with id=275. (Right) The objects with id=498044 (keyboard), 509336(monitor), 454702(keyboard) have the same layout as objects with id=469041 (keyboard), 483912(monitor), 454702(keyboard). Thus, two objects are incorrectly associated (highlighted with red arrows).

Fig. 6. Four examples of failure cases. (a) Failure cases of random walk descriptor-based graph matching. (b) Failure cases of our method.

TABLE II. Trajectory estimation mean error.

DataSet	Sequence	Ours	ORB-SLAM2
TUM RGBD	fr1 room	0.040	0.044
	fr2 desk	0.008	0.010
	fr3 office	0.009	0.008
SceneNN	021	0.066	0.106
	025	0.086	0.116
	231	0.048	0.061

faster than the random walk descriptor-based graph matching method in [17].

Since the S matrix in equation (4) is a highly sparse matrix (for the sequences in Table III, the sparsity is larger than

0.95). The complexity of computing its principal eigenvectors is usually less than $O(n^{3/2})$, where $n = N \times M$. In our implementation, we call the Spectra library [26], which implements the Arnoldi/Lanczos method to find the principal eigenvectors of large symmetric sparse matrices efficiently. Figure 7 shows how the running time for the graph match steps varies with the number of objects (we set $N = M =$ number of objects in our experiments) and the sparsity of matrix S . The results show that our method scales well as the number of objects increases. For huge maps (containing thousands of objects), we can divide the huge map into several smaller submaps and then perform graph matching between the submaps.

TABLE III. Average running time of graph match.

Sequence	Ours (ms)	[17] (ms)	Sparsity of S
fr1 room	0.39	0.99	0.98
fr2 desk	0.035	0.15	0.96
fr3 office	0.14	0.37	0.95

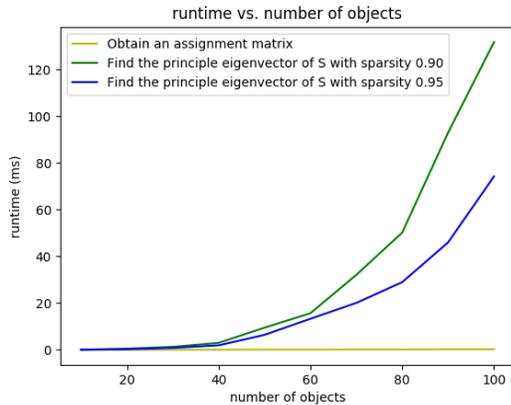


Fig. 7. Running time for the graph matching steps varies with the number of objects and the sparsity of S.

V. CONCLUSION

We have proposed a novel object-level data association to reconstruct the environment as 3D semantic maps. Then we perform loop closure based on semantic graph matching and object alignment. Finally, we jointly optimize camera trajectories and object poses in an object-level pose graph formulation. We have evaluated our methods on public TUM RGBD and SceneNN datasets. Experimental results demonstrate the effectiveness of the algorithms.

We believe that our method can further address the long-term localization challenges of robots, allowing robots to perceive the world in a more human-like manner. However, we need to address several limitations in future work. The partially reconstructed objects may affect the accuracy and robustness of semantic graph matching and object alignment. We plan to introduce learned representations to provide shape priors for better object reconstruction. Currently, our work only exploits semantic labels, spatial distances, and co-visibility. We plan to expand our algorithm with 6-DoF poses and 3D scene graphs.

REFERENCES

- [1] J. Sivic and A. Zisserman, "Video google: A text retrieval approach to object matching in videos," in *Computer Vision, IEEE International Conference on*, vol. 3. IEEE Computer Society, 2003, pp. 1470–1470.
- [2] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, "Orb: An efficient alternative to sift or surf," in *2011 International conference on computer vision*. Ieee, 2011, pp. 2564–2571.
- [3] 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.
- [4] T. Qin, P. Li, and S. Shen, "Vins-mono: A robust and versatile monocular visual-inertial state estimator," *IEEE Transactions on Robotics*, vol. 34, no. 4, pp. 1004–1020, 2018.

- [5] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2961–2969.
- [6] 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.
- [7] J. McCormac, R. Clark, M. Bloesch, A. Davison, and S. Leutenegger, "Fusion++: Volumetric object-level slam," in *2018 international conference on 3D vision (3DV)*. IEEE, 2018, pp. 32–41.
- [8] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers, "A benchmark for the evaluation of rgb-d slam systems," in *Proc. of the International Conference on Intelligent Robot Systems (IROS)*, Oct. 2012.
- [9] B.-S. Hua, Q.-H. Pham, D. T. Nguyen, M.-K. Tran, L.-F. Yu, and S.-K. Yeung, "Scenenn: A scene meshes dataset with annotations," in *2016 Fourth International Conference on 3D Vision (3DV)*. IEEE, 2016, pp. 92–101.
- [10] S. L. Bowman, N. Atanasov, K. Daniilidis, and G. J. Pappas, "Probabilistic data association for semantic slam," in *2017 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2017, pp. 1722–1729.
- [11] K. J. Doherty, D. P. Baxter, E. Schneeweiss, and J. J. Leonard, "Probabilistic data association via mixture models for robust semantic slam," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 1098–1104.
- [12] 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, 2018.
- [13] S. Yang and S. Scherer, "Cubeslam: Monocular 3-d object slam," *IEEE Transactions on Robotics*, vol. 35, no. 4, pp. 925–938, 2019.
- [14] Z. Qian, K. Patath, J. Fu, and J. Xiao, "Semantic slam with autonomous object-level data association," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 11 203–11 209.
- [15] A. Gawel, C. Del Don, R. Siegwart, J. Nieto, and C. Cadena, "X-view: Graph-based semantic multi-view localization," *IEEE Robotics and Automation Letters*, vol. 3, no. 3, pp. 1687–1694, 2018.
- [16] Y. Liu, Y. Petillot, D. Lane, and S. Wang, "Global localization with object-level semantics and topology," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 4909–4915.
- [17] S. Lin, J. Wang, M. Xu, H. Zhao, and Z. Chen, "Topology aware object-level semantic mapping towards more robust loop closure," *IEEE Robotics and Automation Letters*, vol. 6, no. 4, pp. 7041–7048, 2021.
- [18] X. Guo, J. Hu, J. Chen, F. Deng, and T. L. Lam, "Semantic histogram based graph matching for real-time multi-robot global localization in large scale environment," *IEEE Robotics and Automation Letters*, vol. 6, no. 4, pp. 8349–8356, 2021.
- [19] D. F. Crouse, "On implementing 2d rectangular assignment algorithms," *IEEE Transactions on Aerospace and Electronic Systems*, vol. 52, no. 4, pp. 1679–1696, 2016.
- [20] M. Leordeanu and M. Hebert, "A spectral technique for correspondence problems using pairwise constraints," 2005.
- [21] R. B. Rusu, N. Blodow, and M. Beetz, "Fast point feature histograms (fpfh) for 3d registration," in *2009 IEEE international conference on robotics and automation*. IEEE, 2009, pp. 3212–3217.
- [22] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European conference on computer vision*. Springer, 2014, pp. 740–755.
- [23] N. Wojke and A. Bewley, "Deep cosine metric learning for person re-identification," in *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2018, pp. 748–756.
- [24] J. McCormac, A. Handa, S. Leutenegger, and A. J. Davison, "Scenenet rgb-d: Can 5m synthetic images beat generic imagenet pre-training on indoor segmentation?" in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2678–2687.
- [25] P.-E. Sarlin, D. DeTone, T. Malisiewicz, and A. Rabinovich, "Superglue: Learning feature matching with graph neural networks," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 4938–4947.
- [26] "Spectra library," 3 2022. [Online]. Available: <https://github.com/yixuan/spectra>