# SGAligner: 3D Scene Alignment with Scene Graphs

Sayan Deb Sarkar [1], Ondrej Miksik [2], Marc Pollefeys [1,2], Daniel Barath [1], Iro Armeni [1]

[1]Department of Computer Science, ETH Zurich, Switzerland
[2]Microsoft Mixed Reality & AI Lab, Zurich, Switzerland

sgaligner.github.io

## Abstract

*Building 3D scene graphs has recently emerged as a topic in scene representation for several embodied AI applications to represent the world in a structured and rich manner. With their increased use in solving downstream tasks (e.g., navigation and room rearrangement), can we leverage and recycle them for creating 3D maps of environments, a pivotal step in agent operation? We focus on the fundamental problem of aligning pairs of 3D scene graphs whose overlap can range from zero to partial and can contain arbitrary changes. We propose SGAligner, the first method for aligning pairs of 3D scene graphs that is robust to in-the-wild scenarios (i.e., unknown overlap – if any – and changes in the environment). We get inspired by multi-modality knowledge graphs and use contrastive learning to learn a joint, multi-modal embedding space. We evaluate on the 3RScan dataset and further showcase that our method can be used for estimating the transformation between pairs of 3D scenes. Since benchmarks for these tasks are missing, we create them on this dataset. The code, benchmark, and trained models are available on the project website.*

## 1. Introduction

Generating accurate 3D maps of environments is a key focus in computer vision and robotics, being a fundamental component for agents and machines to operate within the scene, make decisions, and perform tasks. As such, these maps should be actionable, *i.e.*, containing information (such as objects, instances, their position, and relationship to other elements) that allows agents to perform an action and represented such that it is easily scalable, updateable, and shareable. Recently, 3D scene graphs [2, 44, 36, 20] have emerged as a topic in scene representation, providing a structured and rich way to represent the world. Not only do they fit the above requirements, but they can also be a lighter-weight [7] and more privacy-aware representation of 3D scenes than the predominantly used 3D point clouds or
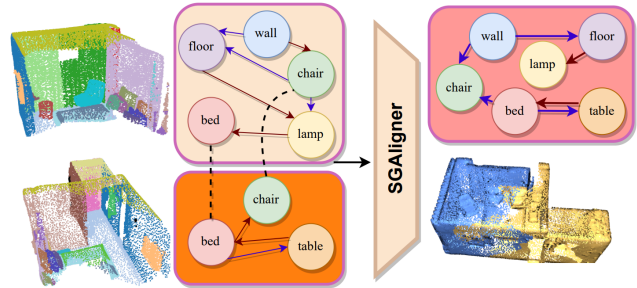


Figure 1. **SGAligner.** We address the problem of aligning 3D scene graphs of environments using multi-modal learning and leverage the output for the downstream task of 3D point cloud registration. Our approach operates directly on 3D scene graph level, is fast and robust to real-world scenarios.

voxel grids – hence being easier and safer to share across agents and humans operating in the same scene [24, 51].

Given their potential, 3D scene graphs are increasingly used in embodied agents as a representation – commonly built on the fly – to perform robotic navigation [39, 37, 35, 24, 7] and task completion [13, 1, 34, 19, 22]. Since more and more agents are already building 3D scene graphs for downstream tasks, we investigate how to leverage and recycle them for creating 3D maps of the environment – a pivotal step in the agent operation – directly on the scene graph level. Specifically, we examine the fundamental problem of aligning partial 3D scene graphs of an environment that originate from different observations. We focus on real-world scenarios and specifically formulate the problem as follows: given two 3D scene graphs that represent 3D scenes whose overlap can range from zero to partial or full and can contain changes, our goal is to find an alignment across nodes, if it exists. Interestingly enough, even though entity alignment (*i.e.*, node alignment) is used in knowledge graphs and in linguistics [26, 14, 28, 50, 8, 10, 40, 25], the task of aligning 3D scene graphs of environments has not been explored. An important note is that entity alignment in these domains assumes that there is overlap between graphs and that all inputs contain true information.

We propose **SGAligner**, the *first* method for aligning pairs of 3D scene graphs that is robust to in-the-wild scenarios (*i.e.*, unknown overlap – if any – and changes in the environment) (see Figure 1). We get inspired by entity alignment methods in multi-modality knowledge graphs [25] and redesign them for our setting. 3D scene graphs represent three main types of information [44, 2, 36]: semantic entities in the scene (*e.g.*, object instances), their attributes (*e.g.*, category, size, and material), and relationships between the entities (*e.g.*, relative position and attribute similarity). The main premise is to independently encode each of these modalities with the ultimate objective of learning a joint embedding that can reason how similar two nodes are. Given node matches, we perform the scene graph alignment using the matches with the highest similarity.

We additionally demonstrate our scene graph alignment method on the tasks of 3D point cloud registration and 3D alignment of a local 3D point cloud on a larger map that contains changes. Instead of directly computing 3D correspondences on the entire point clouds [16, 32, 47, 4], we use the alignment as coarse initialization for the registration. We further refine it by computing 3D correspondences [32] on the individual point clouds (*i.e.*, object instance point clouds) of each matched node pair. This is followed by robustly estimating [12] the rigid point cloud transformation using the correspondences from all matched nodes.

We evaluate all three tasks on the 3RScan [43] dataset, which contains 3D point clouds captured over time along with their 3D scene graph annotations [44]. Since 3RScan does not provide partial scene graphs of the same scene or a point cloud registration benchmark and there is no 3D scene graph alignment benchmark, we create the data, metrics, and evaluation needed for these tasks in 3RScan. Our experiments show that our approach reduces the relative translation error of state-of-the-art GeoTransformer [32] by 40% in point cloud registration, while being 3× faster during the overlap check, since it does not need to process the entire point clouds. Detailed ablation studies, along with experiments on the task of aligning a changed local 3D scene to a prior 3D map, demonstrate robustness of our approach.

We summarize the contributions of this paper as follows:

- We propose **SGAligner**, the first method for aligning pairs of 3D scene graphs whose overlap can range from zero to partial and that may contain changes.

- We demonstrate the potential of our method on the tasks of 3D point cloud registration, 3D point cloud mosaicking and 3D alignment of a point cloud in a larger map that contains changes.

- We create a scene graph alignment and 3D point cloud registration benchmark on the 3RScan [43, 44] dataset, with data, metrics, and evaluation procedure.

## 2. Related Work

**Multi-Modal Knowledge Graph Alignment.** There exists vast literature in the domain of graph matching and ontology alignment, with methods focusing on creating single feature vectors to compare on the entire graph-level (*e.g.*, [23]). Such methods are not applicable to our case, since the two graphs to align have partial semantic and geometric overlap. A more relevant task is that of multi-modal knowledge graph (KG) alignment [26, 14, 28, 50, 8, 10, 9, 40, 25], which refers to the task of aligning multiple knowledge graphs that represent information from different modalities (*e.g.*, text, images, and videos). The goal is to integrate the knowledge from different sources and provide a more comprehensive understanding of the world. EVA [26] leverages visual and auxiliary knowledge to achieve entity alignment in both supervised and unsupervised settings, using a loss formulation inspired by NCA-based text-image matching [27]. More recently, approaches like [8, 25] solve the task by learning a common embedding space for all modalities, where similar entities in the KG have similar embeddings. Both approaches use multiple individual encoders to obtain modality-specific representations for each entity in the KG. However, [25] introduces contrastive learning with intra-modal contrastive loss and inter-modal alignment loss to learn discriminative cross-modal embeddings, while [8] only performs common space learning to align the embeddings. Since 3D scene graphs can be also considered as containing multiple modalities (*i.e.*, object instances, their 3D geometry, attributes, and in-between relationships), we leverage the above architecture and adapt it for the setting of aligning 3D scene graphs of environments. A point to highlight is that KG alignment methods consider inputs as (partially) overlapping and true, something that in the real-world scenario of creating 3D maps does not hold due to arbitrary conditions and noise.

**3D Scene Graphs.** Following the success of utilizing 2D scene graphs [21], researchers introduced 3D scene graphs as structured and rich representations to describe real-world scenes [2, 20, 36, 44]. We follow the 3D scene graph structure and dataset presented in [44]. Existing work addresses the task of generating 3D scene graphs with a variety of approaches; from online incremental building [45, 17] to offline based on RGBD images and/or 3D reconstructions [2, 44, 37]. Similarly, [52] advocates a knowledge-inspired scene graph prediction method based on point clouds while [49] uses edge-assisted reasoning to bridge perception and reasoning in the context of 3D scene understanding. The 3D scene graph representation has been explored in embodied AI for tasks related to navigation and planning [39, 37, 35, 24, 7], task completion [13, 1, 34, 19, 22], variability estimation [29], and 3D scene manipulation [11]. Although prior work has demonstrated the applicability of 3D
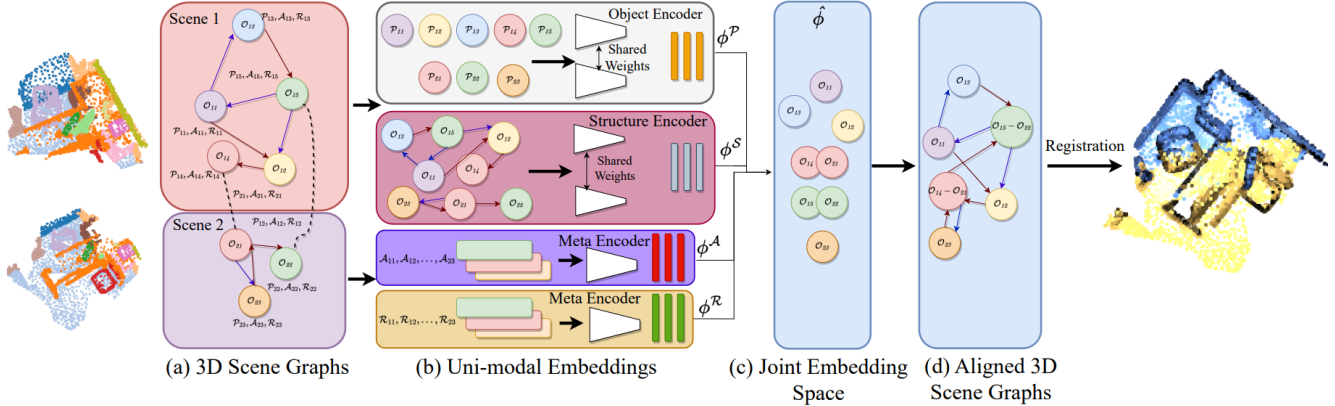
Figure 2. **Overview of SGAligner.** Given two 3D scene graphs with partial overlap (a), we create four uni-modal embeddings (b) that encode object, structure, attribute, and relationship information. These interact with each other to create a joint embedding space (c) where similar nodes are located closely. We use this similarity to match nodes and create the final aligned 3D scene graph (d). We further demonstrate the use of the alignment in 3D point cloud registration.

scene graphs in an increasing fashion, it has not been investigated for the task of creating 3D maps of scenes. Hydra [17] approaches the latter using hierarchical loop closure detection and 3D scene graph optimization to complete the non-optimised scene graph, that is built on the fly, into a globally consistent and persistent one as new information is captured. On the another hand, our contribution focuses on solving graph matching solely on the object-level. This would enable to leverage, share, and recycle the created graphs for general agent operation.

**3D Point Cloud Registration.** The field of 3D point cloud registration is well-established and active, with approaches mainly being feature-based and end-to-end. We focus our scope on feature-based approaches since they compute hard correspondences and perform more robustly in real-world scenes. Such methods [4, 16, 47, 32] consist of two steps: local (learned) feature extraction and pose estimation with RANSAC [12]. However, they focus on the in-vitro problem of aligning two input point clouds that have some degree of overlap. This does not always hold in a real-world scenario where there is no knowledge of whether there is any overlap or changes. Although some of these methods compute matchability scores per estimated 3D point correspondences (*e.g.*, [16, 32]), they assume overlapping input point clouds and do not have the mechanism to discard non-overlapping ones. In addition, when the number of individual point clouds to register becomes large, it requires $O(\mathcal{N}^2)$ complexity to process all possible pairs, while failing to recognize non-aligned pairs. Last, they also require large memory reserve to process input data if the point clouds are large. Our method allows processing all possible pairs faster while identifying non-overlapping pairs before performing the final registration. In addition, it is lightweight and can easily process large scenes, since, as shown in our experiments, our method can operate with a limited number of points per object instance. Please refer to

Sec. B in the supp. mat for more details on this.

## 3. SGAligner: 3D Scene Graph Alignment

Following standard formulation [44], we define a 3D scene graph $\mathcal{G}$ of a scene $s$ as a pair of sets $(\mathcal{N}, \mathcal{R})$ with nodes $\mathcal{N}$ and edges $\mathcal{R}$. The nodes represent 3D object instances $\mathcal{O}$ in the scan. Each node also contains a set of attributes $\mathcal{A}$ that characterizes the visual (*e.g.*, style and texture), geometric (*e.g.*, shape and symmetry), and state (*e.g.*, closed and empty) characteristics of the object instance, in addition to the object categories. Instance-level point clouds $\mathcal{P}$ are available per node. The edges define semantic relationships between the nodes such as `standing on` and `attached to`. Given two graphs $\mathcal{G}_1 = (\mathcal{N}_1, \mathcal{R}_1)$ and $\mathcal{G}_2 = (\mathcal{N}_2, \mathcal{R}_2)$ of scenes $s_1$ and $s_2$ respectively, we aim to find the set of objects in the overlapping regions of the two scenes, denoted as entity pairs, $\mathcal{F} = \{(n_1, n_2) \mid n_1 \equiv n_2, n_1 \in N_1, n_2 \in N_2\}$, forming node correspondences.

Our approach follows the network architecture proposed in [25], which we modify from the language domain for the 3D scene graph alignment task with varying degrees of overlap. The overall architecture is shown in Figure 2. We can formulate scene graphs as a multi-modal knowledge graph – commonly used in entity alignment – where the semantic and geometric information included in scene graphs is treated as the different modalities that are encountered in knowledge graphs. The goal is to learn a joint multi-modal embedding from individual encodings of each modality (uni-modal), in which nodes are closely located if they correspond to the same underlying object instance and belong to different graphs. Specifically, we create uni-modal embeddings using the three main 3D scene graph information types: *object* embedding that encodes $\mathcal{P}$, *structure* embedding $\mathcal{S}$ that encodes $\mathcal{R}$ in the form of a structured graph, and two *meta* modalities that encode $\mathcal{A}$ and $\mathcal{R}$

between objects in the form of a one-hot vector. To reason about the final entity alignment, similar to [25], we create a joint embedding by combining these uni-modal ones in a weighted manner and perform a joint optimization using knowledge distillation across all embeddings.

## 3.1. Uni- and Multi-Modal Embeddings

To leverage rich information in 3D scene graphs, we process each modality separately in our framework and create uni-modal embeddings, which are later processed to model complex inter-modal interactions in the joint embedding.

**Object Embedding.** Point clouds contain rich geometric information about objects. Each of the individual point cloud $\mathcal{P}_i$ of $\mathcal{O}_i$ is the input to the object encoder. We employ a point cloud feature extractor backbone architecture such as [31, 15] as the object encoder to extract the visual features $\phi_i^{\mathcal{P}}$ for every node.

**Structure Embedding.** 3D Scene Graphs contain information on relationships between $\mathcal{O}$, which we leverage to encode the layout of objects in $s$. We represent this information in the form of a *structure* graph: node features are the relative translation between object instances, and edges are the aforementioned relationships. We calculate relative translation by taking the distance between the object instance consisting of the highest number of relationships and that of any other object instance in the scene. Specifically, we compute distances using the barycenter of the convex hull of the point clouds. We use a Graph Attention Network (GAT) [42] to model the structural information in $\mathcal{G}_1$ and $\mathcal{G}_2$ via the structure graph. We limit the weight matrix to a diagonal matrix, as suggested by [25], to minimize computations and improve the scalability of the model. As per [25], the neighborhood structure embedding $\phi_i^{\mathcal{S}}$ is produced by the last GAT layer, using a two-layer GAT model to aggregate the neighborhood information over several hops.

**Meta Embeddings.** Along with modeling the geometric and structural properties of the scene, we model the attributes and corresponding relationships per object $\mathcal{O}_i$ in two separate embeddings. We regard the relationship of $\mathcal{O}_i$ with other objects in the input scene graph as a one-hot encoded feature vector and pass it through a single-layer MLP to obtain the relational embedding $\phi_i^{\mathcal{R}}$. For instance, consider the set of relationships `standing on`, `built in`, `attached to`. If an object (node) in the input scene graph has a single edge associated with it and this represents the `built in` relationship, the input vector for the relationship encoder is $[0, 1, 0]$. We adopt the same approach for the attributes of $\mathcal{O}_i$ for simplicity to get the $\phi_i^{\mathcal{A}}$. These single-layer networks called meta encoders provide valuable insights into scene composition and facilitate straightforward extension to new data.

**Joint Embedding.** We concatenate each of the previously discussed uni-modal features to a single compact represen-

tation $\hat{\phi}_i$ for each object $\mathcal{O}_i$ as follows:

$$\hat{\phi}_i = \oplus_{k \in \mathcal{K}} \left[ \frac{\exp(w_k)}{\sum_{j \in \mathcal{K}} \exp(w_j)} h_i^m \right], \quad (1)$$

where $\oplus$ denotes concatenation, $\mathcal{K} = \{\mathcal{P}, \mathcal{S}, \mathcal{R}, \mathcal{A}\}$, and $w_m$ is a trainable attention weight for each modality $k$. We apply $L_2$ normalization to each uni-modal feature before the final concatenation. These embeddings are coarse-grained without any interaction between different modalities.

## 3.2. Contrastive Learning

To model interaction between modalities, we formulate a representation learning framework. A contrastive loss function encourages comparable samples, or aligned entities, to be closer together and dissimilar samples to be farther away in the learned representation space. Using a cross-modal contrastive loss is a typical strategy when working with many modalities, such as in the case of 3D scene graphs, as it encourages samples from various modalities that are semantically related to be closer together in the joint representation space. Inspired by [25], we use the Intra-Modal Contrastive Loss (ICL) and Inter-modal Alignment Loss (IAL) and formulate them similarly.

During training, we assume that $E \subset \mathcal{F}$ is available to us as seed-aligned entity pairs. Formally, for the $i^{th}$ object node $n_1 \in \mathcal{N}_1$, we define $E = \{n_1^i \mid n_2^i \in \mathcal{N}_2\}$, where $(n_1^i, n_2^i)$ is an aligned pair. We define the unaligned pairs within the same graph as $H_1^i = \{n_1^j \mid \forall n_1^j \in \mathcal{N}_1, j \neq i\}$, and aligned pairs across graphs as $H_2^i = \{n_2^j \mid \forall n_2^j \in \mathcal{N}_2, j \neq i\}$ (Figure 3). These two samples define the constrained joint embedding space. We model $\mathcal{L}_k^{ICL}$ to learn intra-modal dynamics for more discriminative boundaries for each modality $k$ in the embedding space. We apply ICL separately on each uni-modal embedding and on the joint concatenated embedding, after $L_2$ normalization. Each uni-modal embedding is trained individually using ICL and is not intended to interact with others.

Our complete representation is the joint embedding space, and our goal is to learn proper uni-modal encodings that enable node alignment in this joint space. To achieve this, we minimize the bi-directional KL-divergence loss between joint embedding and uni-modal embeddings as the Inter-modal Alignment Loss (IAL), thereby, emphasizing on aggregating the distribution of various modalities, which narrows the modality gap by learning interactions between various modalities inside each entity. We train our model end-to-end, optimizing both losses as follows:

$$\mathcal{L} = \mathcal{L}_o^{ICL} + \sum_{k \in \mathcal{K}} \alpha_k \mathcal{L}_k^{ICL} + \sum_{k \in \mathcal{K}} \beta_k \mathcal{L}_k^{IAL}, \quad (2)$$

where $\mathcal{K} = \{\mathcal{P}, \mathcal{S}, \mathcal{R}, \mathcal{A}\}$ and $o$ is the joint embedding. Variables $\alpha_k$ and $\beta_k$ are hyper-parameters that are automat-
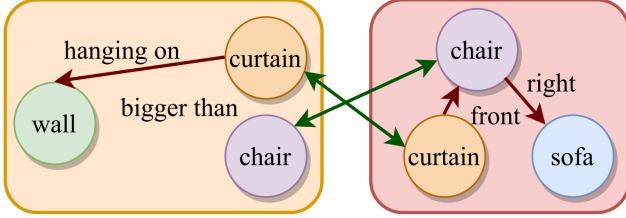
Figure 3. **An example of contrastive learning pairs.** Aligned pairs across graphs are marked with green edges and unaligned pairs within the same graph are marked with brown edges.

ically learned during training. We direct the readers to [25] for a deeper understanding of the loss functions.

### 3.3. 3D Point Cloud Registration

In this section, we describe how we approach the task of 3D point cloud registration by leveraging our scene graph alignment results. The output of the previously described scene alignment method is the set of matched entity pairs $n_1$ and $n_2$, in the scene graphs $\mathcal{G}_1$ and $\mathcal{G}_2$, respectively. For each entity pair $n_1^i$ and $n_2^i$, we extract 3D point correspondences from $\mathcal{P}_1^i$ and $\mathcal{P}_2^i$. The correspondences are estimated by an off-the-shelf correspondence extraction algorithm (*e.g.*, [32]) by running it on node pairs independently. We collect the point correspondences across all matched entities and then use the robust estimator, *e.g.* RANSAC [12] or one of its recent variants [33, 5, 6], to get the transformation $\mathbf{T} \in \mathrm{SE}(3)$ between the point clouds of the two scenes.

Performing registration on 3D correspondences that stem from node-to-node matches allows for being less sensitive to changes in the point clouds and incorrect matching than state-of-the-art techniques. Such an approach has two major advantages: (i) it filters non-overlapping scene pairs, which should not be registered, **faster** than state-of-the-art point cloud registration methods and without any need for registration. This is enabled by obtaining object-level (node-to-node) correspondences instead of performing the registration directly on a large-scale 3D point cloud. (ii) It performs better than standard registration methods even on point clouds with low overlap, which we showcase with experiments in the following section.

## 4. Experiments

We evaluate **SGAligner** on the task of 3D scene graph alignment (Section 4.1) and on downstream applications, namely 3D point cloud registration (Section 4.2), 3D point cloud mosaicking (Section 4.3) and 3D scene alignment with changes present in the data (Section 4.4). We provide additional ablation studies to further understand the performance on node matching in the supp. mat.

In our experiments we use the 3RScan dataset [43], which consists of 1335 indoor scenes, of which 1178 are used for training and 157 for testing. The dataset contains

semantically annotated 3D point clouds per scene, some of which depict the same environment over time. 3D scene graph annotations for 3R Scan are provided in [44]. Although this allows to evaluate the robustness of our method in the case of changed environments, there are no annotations to evaluate in static environments. To enable a thorough evaluation with scenes that range in overlap (from zero to partial), we create sub-scene graphs in single scenes given their full 3D scene graph provided in the annotations, following standard point cloud registration benchmarks [48]. The generated data is made public to maximize reproducibility. In our experiments, we use PointNet [31] as the object encoder, a comprehensive evaluation with multiple backbones is provided in the supp. mat.

**Dataset Generation.** To create the sub-scene graphs, we generate sub-scenes per scene on the geometry level. Point clouds in [43] are generated from RGBD frames. To imitate a realistic setting, we create sub-scenes by selecting groups of sequential frames and reconstructing the point cloud of the depicted scene. Frames across groups are unique and there is a possibility of 3D spatial overlap in the generated point clouds. We generate a total of 6731 sub-scenes from the training scenes and 843 sub-scenes from the validation scenes. We create pairs using the sub-scenes generated from the same scene, such that the spatial overlap in a pair ranges from 10-90%. This results to 15102 pairs for training and 1932 pairs for testing. More details, as well as statistics on the generated data are in the supp. mat. For simplicity, hereafter we refer to sub-scenes as *scenes* ($s$ in Section 3).

The individual point clouds of object instances in the generated data will have a varying number of points. To use them as input to the object encoder (Section 3.1), we require them to have the same size (512). We use farthest point sampling to downsample each $\mathcal{P}_i$ of object $\mathcal{O}_i$ as

$$\mathcal{P}_i = \{\delta_{k^i} \odot p_k\}_{k=1,|P|}, \qquad (3)$$

where $\delta$ represents the Kronecker delta, $p$ is a point in $\mathcal{P}$, and $|.|$ is the cardinality of $\mathcal{P}_\rangle$, i.e, the number of points.

### 4.1. 3D Scene Graph Alignment

In this section, we evaluate **SGAligner** on the 3D scene graph alignment task. We first evaluate how our method performs on node alignment. We utilize cosine similarity on the joint embedding to calculate the similarity between two matched entities and employ Mean Reciprocal Rank (MRR) and Hits@K, where K=$\{1, 2, \ldots, 5\}$. MRR denotes the mean reciprocal rank of correct matches. Hits@K denotes the ratio of correct matches appearing within top K, based on their cosine similarity ranking.

We compare **SGAligner** to using different modalities, and as a result embeddings, as well as with a baseline from entity alignment in the language domain. For the latter purpose, we adapt the Entity Visual Alignment method (EVA)
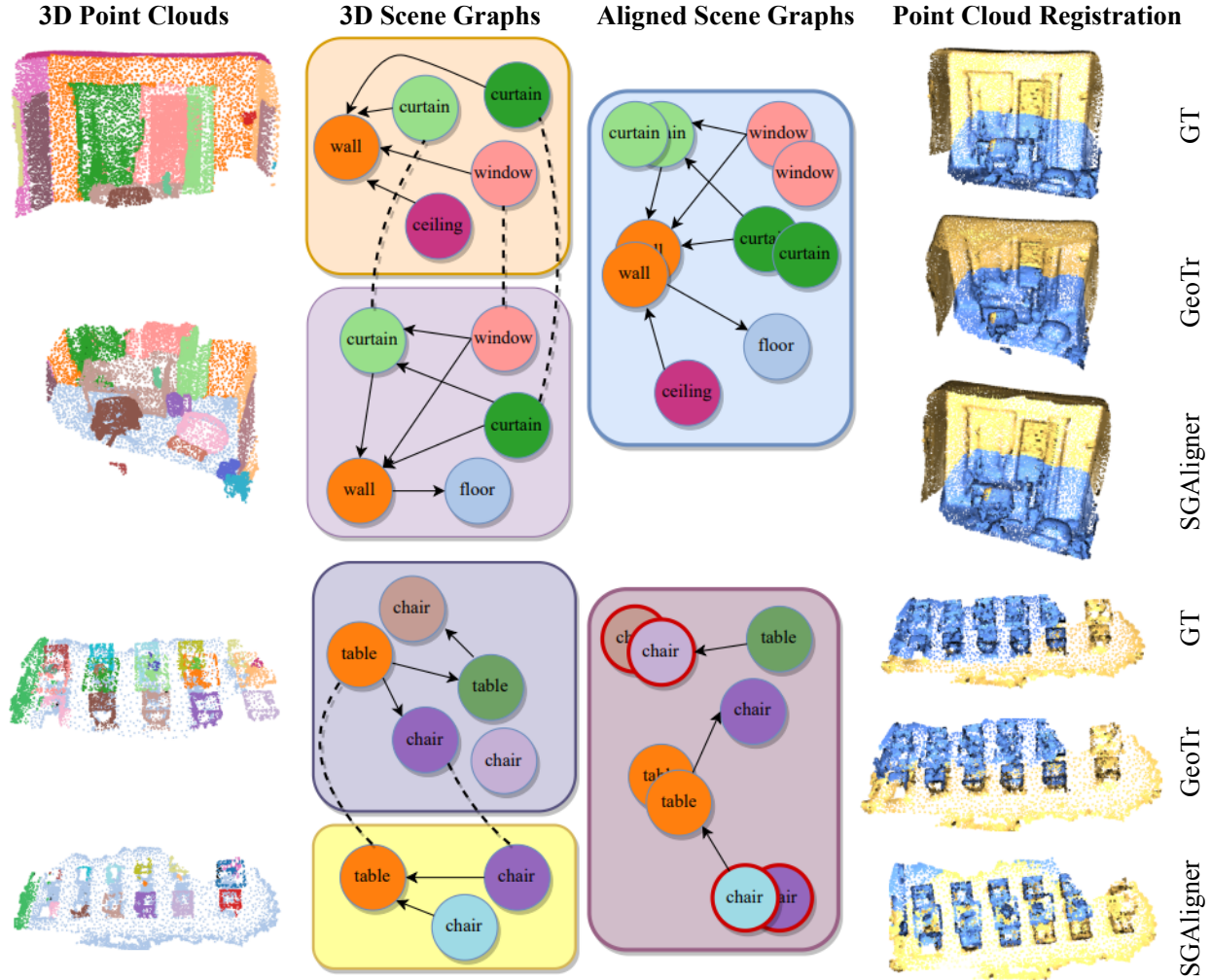
Figure 4. **Qualitative Results.** Given input 3D scene graphs (left), we showcase a success (top) and a failure (bottom) case of **SGAligner** for alignment (middle) and registration (right). The existence of the same exact object in multiple replicas creates erroneous matches (nodes outlined in red) which cannot be recovered in registration. When objects are more discriminative, **SGAligner** provides accurate matching.

| Modalities | Mean RR ↑ | Hits @ ↑ | | | | |
|---|---|---|---|---|---|---|
| | | K = 1 | K = 2 | K = 3 | K = 4 | K = 5 |
| *w/ Ground Truth 3D Scene Graphs* | | | | | | |
| EVA [26] | 0.867 | 0.790 | 0.884 | 0.938 | 0.963 | 0.977 |
| $\mathcal{P}$ | 0.884 | 0.835 | 0.886 | 0.921 | 0.938 | 0.951 |
| $\mathcal{P} + \mathcal{S}$ | 0.897 | 0.852 | 0.899 | 0.931 | 0.945 | 0.955 |
| $\mathcal{P} + \mathcal{S} + \mathcal{R}$ | 0.911 | 0.861 | 0.916 | 0.947 | 0.961 | 0.970 |
| **SGAligner** | **0.950** | **0.923** | **0.957** | **0.974** | **0.982** | **0.987** |
| *w/ Predicted 3D Scene Graphs* | | | | | | |
| **SGAligner** | 0.882 | 0.833 | 0.881 | 0.918 | 0.937 | 0.951 |

Table 1. **Evaluation on node matching.** We compare the performance of **SGAligner** for different modality combinations, as well as for ground truth and predicted scene graphs.

[26] for 3D scene graphs by replacing the visual encoder with PointNet architecture [31], same as $\mathcal{P}$ in our approach. Results are in Table 1. Please note that when employing

only the instance level point clouds $\mathcal{P}$ there is no IAL used. As evident, our method, even when using a single modality, outperforms EVA [26] with a margin of approximately 2%. Furthermore, and as expected, each modality in our method contributes to improved performance in all metrics. Interestingly, using all modalities provides at K=2 better results than only $\mathcal{P}$ at k=5, and at k=3 it is already better than any of the other combinations at k=5. To further verify the robustness of our method and hence its suitability for real-world applications, we also compare the performance of our method using both ground truth and predicted 3D scene graphs as input during inference (in both cases the network has been trained on ground truth data). We compute the scene graph predictions using the pre-trained network for 3D scene graph generation given 3D point cloud of [44][1].

---

[1]We refer the reader to Table 2 in [44] for an evaluation on scene graph prediction of this method.

**SGAligner** is still able to provide reasonable matches despite the noise in the input, although with an expected performance drop. The success and failure case presented in Figure 4, shows that the existence of the same exact object in multiple replicas creates erroneous matches; when the objects are more discriminative, **SGAligner** provides accurate results. For more examples please see supp. mat.

We further ablate the results of our method (on ground truth data) for scenes of different overlap to evaluate the robustness in cases where few nodes can be matched. Results are reported in Table 2. As expected, the performance drops with lower overlap, however, the gap between the very high and very low overlap is not drastic.

| Overlap (%) | Mean RR ↑ | Hits @ ↑ | | | | |
|---|---|---|---|---|---|---|
| | | K = 1 | K = 2 | K = 3 | K = 4 | K = 5 |
| 10-30 | 0.872 | 0.806 | 0.886 | 0.927 | 0.948 | 0.962 |
| 30-40 | 0.908 | 0.859 | 0.917 | 0.949 | 0.968 | 0.978 |
| 40-50 | 0.941 | 0.908 | 0.950 | 0.973 | 0.980 | 0.985 |
| 50-60 | 0.960 | 0.937 | 0.967 | 0.982 | 0.989 | 0.994 |
| 60- | **0.979** | **0.966** | **0.982** | **0.990** | **0.994** | **0.995** |
| All data | 0.950 | 0.923 | 0.957 | 0.974 | 0.982 | 0.987 |

Table 2. **Evaluation on node matching per overlap range.**

In order to deeper understand the effect of semantic noise to **SGAligner**'s performance, we provide experiments with controlled semantic noise on the test data. Specifically, we consider the following scenarios of noise: (i) only relationships are removed[2]; (ii) only object instances are removed – their corresponding attributes and any relationships that include them are also removed; (iii) both relationships and object instances are removed; (iv) object instances are wrong (*i.e.* they have the wrong semantic label); and (v) both relationships and objects are wrong. Noise is applied to each input scene graph randomly to 15-40% of the modified semantic. For (iv) and (v), we randomly assign any other semantic label, *i.e.*, a chair could be labeled as floor. Results, including those on the full ground truth dataset (GT) and with predicted 3D scene graphs (Pred.) for reference, are in Table 3. Note that the training set remains the same. As shown, the noise that our method can handle the best is that of missing objects (ii). This means that the relationships between objects can be more important than having structured information, however not by a large margin. What drops the performance drastically is wrong semantic labels for both object and relationship labels. Comparing the results for scenarios (iii) and (iv) with the use of predicted 3D scene graphs, we observe that for the former the values are significantly lower. This has to be taken into account in case of real-world applications, especially if generalization of the scene graph prediction algorithm is unknown.

---

[2]Ground truth annotations do not offer an exhaustive list of relationships and attributes per node. We remove edges from these annotations.

| | Affected Modules | Mean RR ↑ | Hits @ ↑ | | | | |
|---|---|---|---|---|---|---|---|
| | | | K = 1 | K = 2 | K = 3 | K = 4 | K = 5 |
| (i) | $\mathcal{S}, \mathcal{R}$ | **0.906** | **0.858** | 0.915 | **0.949** | **0.964** | 0.974 |
| (ii) | All | <u>**0.924**</u> | <u>**0.878**</u> | <u>**0.942**</u> | <u>**0.968**</u> | <u>**0.977**</u> | <u>**0.985**</u> |
| (iii) | All | 0.902 | 0.848 | **0.918** | **0.949** | **0.964** | **0.975** |
| (iv) | $\mathcal{P}$ | 0.661 | 0.563 | 0.643 | 0.699 | 0.743 | 0.776 |
| (v) | $\mathcal{P}, \mathcal{S}, \mathcal{R}$ | 0.587 | 0.479 | 0.573 | 0.632 | 0.674 | 0.709 |
| GT | None | 0.950 | 0.923 | 0.957 | 0.974 | 0.982 | 0.987 |
| Pred. | None | 0.882 | 0.833 | 0.881 | 0.918 | 0.937 | 0.951 |

Table 3. **Evaluation on node matching with controlled semantic noise.** *Best* values are in <u>**underlined bold**</u>. *Second best* in **bold**.

**3D Scene Graph Alignment.** Entity alignment provides pairs of matched nodes. Here, we evaluate how well two scene graphs can be aligned given the predicted node matching. In theory, since the scenes are rigid, two matched nodes would be enough to perform the alignment, if there was no noise in the matches. Since in reality matches are noisy, we evaluate 3D scene graph alignment with K equals top-2, top-50%, and all of the matched nodes. To measure this performance we introduce the *Scene Graph Alignment Recall (SGAR)* metric. Similar to the standard recall metric, we calculate the amount of correct alignments of **SGAligner**. We provide these results in Table 4 for both ground truth and predicted scene graphs. As shown, results with K=2 perform the best in both cases. This means that **SGAligner** can identify at least two matches that are very closely located in the joint embedding space. In addition, SGAR for the top-2 matches is approximately the same for both ground truth and predicted scene graphs, which further shows that our approach can robustly align them.

| Input Scene Graphs | Scene Graph Alignment Recall ↑ | | |
|---|---|---|---|
| | R@top-2 | R@top-50% | R@All |
| Ground Truth | **0.964** | 0.948 | 0.738 |
| Predicted | **0.963** | 0.856 | 0.450 |

Table 4. **Evaluation on 3D scene graph alignment.** We report for both ground truth and predicted scene graphs.

### 4.2. 3D Point Cloud Registration

In this section, we evaluate the performance of **SGAligner** for 3D point cloud registration on the same data. We employ the state-of-the-art Geotransformer [32] as a 3D correspondence extraction method, as per Section 3.3. We compare the performance of our approach to standard approaches and use Geotransformer directly on the point clouds of the two scenes to register. Please note that in our case, we use Geotransformer to extract correspondences from individual point clouds on the level of object instances only for the matched nodes. In all cases, we use the Geotransformer model trained on the 3DMatch dataset [48].

**Metrics:** We compute the standard metrics of Chamfer distance (CD) as in [46], relative rotation error (RRE), rel-

ative translation error (RTE), feature match recall (FMR), and registration recall (RR). RR is calculated with the standard threshold of RMSE = 0.2.

Results for 3D point cloud registration for *overlapping scenes* are shown in Table 5. When using predicted data from SceneGraphFusion [45], our method not only remains robust but also shows even higher gains, as we leverage node-to-node alignment for improved local context compared to global registration on noisy point cloud predictions. Our method is consistently providing best results in all metrics w.r.t. the standard registration approach, despite the less geometric information in the point clouds we use; specifically on CD, we provide a 49.4% improvement (K=2). This is an expected behavior since our alignment method, even when it contains incorrect matches, is still providing an initialization to the task and narrows down the search space for correspondences. With respect to self baselines, our method with K=2 performs the best. The drop in K=3 can be attributed to the following: more matches per each node of which one is correct, means that more outliers are provided to RANSAC which it cannot easily remove. Since the gap of Hits@K between K=2 and K=1 is larger than that between K=3 and K=2, K=2 can still provide a boost of inliers to RANSAC even though it will also increase outliers with respect to K=1. We include examples on in Figures 4 and 5.

| Methods | CD ↓ | RRE (°) ↓ | RTE (cm) ↓ | FMR (%) ↑ | RR(%) ↑ |
|---|---|---|---|---|---|
| *w/ Ground Truth 3D Scene Graphs* | | | | | |
| GeoTr | 0.02247 | 1.813 | 2.79 | **98.94** | 98.49 |
| Ours K=1 | 0.01677 | **1.425** | 2.88 | <u>**99.85**</u> | 98.79 |
| Ours K=2 | <u>**0.01111**</u> | <u>**1.012**</u> | <u>**1.67**</u> | <u>**99.85**</u> | <u>**99.40**</u> |
| Ours K=3 | **0.01525** | 1.736 | **2.55** | <u>**99.85**</u> | 98.81 |
| *w/ Predicted 3D Scene Graphs* | | | | | |
| GeoTr | 0.06643 | 5.697 | 9.54 | 92.23 | 93.15 |
| Ours K=1 | 0.05041 | 2.49 | 3.86 | 95.25 | 94.95 |
| Ours K=2 | 0.04251 | 1.725 | 3.36 | 97.12 | 98.33 |
| Ours K=3 | 0.04863 | 2.194 | 2.55 | 96.83 | 97.96 |

Table 5. **3D Point Cloud Registration.** *Best* values are in <u>**underlined bold**</u>. *Second best* in **bold**.

We further ablate the results on this task for scenes of different overlap. As shown in Table 6, our approach outperforms the standard one per overlap, and even provides better results for scenes with 30% overlap or higher than the standard approach can do on 60% or higher.

| Overlap (%) | CD ↓ | RRE ↓ (°) | RTE (cm) ↓ | FMR (%) ↑ | RR (%) ↑ |
|---|---|---|---|---|---|
| GeoTr. 10-30 | 0.09788 | 8.830 | 13.56 | 94.57 | 92.25 |
| GeoTr. 30-60 | 0.00584 | 0.156 | 0.24 | 97.28 | 97.36 |
| GeoTr. 60- | 0.00177 | 0.048 | 0.07 | 99.47 | 99.31 |
| Ours 10-30 | 0.05160 | 5.660 | 8.48 | 99.23 | 95.35 |
| Ours 30-60 | **0.00127** | **0.045** | **0.05** | **99.68** | **98.34** |
| Ours 60- | <u>**0.00046**</u> | <u>**0.018**</u> | <u>**0.02**</u> | <u>**99.92**</u> | <u>**99.93**</u> |

Table 6. **3D Point Cloud Registration per overlap.** For our method we use the best performing (K=2). *Best* values are in <u>**underlined bold**</u>. *Second best* in **bold**.
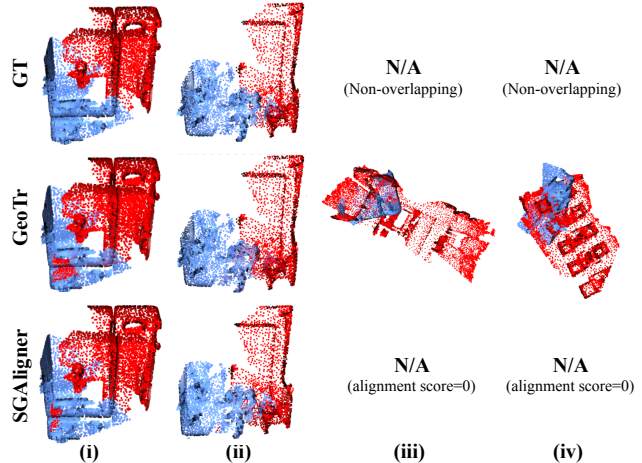


Figure 5. **Registration Results for SGAligner and [32]**. Columns (i)-(ii) represent registration on overlapping pairs. (iii)-(iv) show incorrect results of [32] for **non-overlapping** pairs, where our *alignment score* is zero (*i.e.*, there is no registration output to report for **SGAligner**). Since these pairs are not overlapping, there is also no ground truth registration.

**Overlapping vs Non-Overlapping Scenes.** In practical and real-world applications, we do not always know if two scenes are overlapping or not. While standard point cloud registration approaches compute a matchability score, they typically train and test only on *overlapping* scenes and do not have the mechanism to discard non-overlapping ones. Here, we demonstrate how our approach can indirectly provide a solution to this problem, without any additional supervision. We formulate it as to how well a method can identify whether a pair of scenes overlaps. For the state-of-the-art Geotransformer, we compute a scene-level *matchability* $\mu$ by averaging over all correspondence matchability scores and consider two scenes as overlapping if $\mu \geq 0.2$. For our approach, we compute a scene-level *alignment* score $\xi$ by averaging the total number of matched nodes (for K=1) and consider two scenes as overlapping if $\xi \geq 0.2$.

To evaluate this task, we create a new test set that includes overlapping and non-overlapping scenes from the data we generate. Specifically, we use 1932 of the overlapping pairs from the former, and sample from the rest of the scenes 1932 non-overlapping pairs. Please note that the training set remains the same and contains *only* overlapping pairs. We compute the precision, recall, F1-score, and time in milliseconds required to make the overlap decision for all pairs. Results are in Table 7. Our approach runs 3 times faster than Geotransfomer since it does not process the entire point cloud, which can be computationally demanding. In addition, it leads to comparative performance while being able to identify more overlapping pairs correctly.

We choose $\mu$ and $\xi$ based on the best performing value of $\mu$ for Geotransformer. Specifically, we observed that a higher value (0.4) leads to no true positives and a lower one

(0.1) leads to lower precision (66.94%) since it identifies more false positives. This is understandable given that it was not trained on this data. However, (0.2) for a metric measuring average point correspondence similarity is low. For $\xi$ we observe only a slight decrease in precision instead of the best performing value (0.4). In contrast, our approach does not require such difficult-to-set thresholding scheme.

| Method | Prec. (%) ↑ | Recall (%) ↑ | F1 (%) ↑ | Average Time Per Scene (ms) ↓ |
|---|---|---|---|---|
| Geotr | **99.63** | 80.98 | 89.34 | 442.50 |
| Ours | 92.03 | **90.94** | **91.48** | **139.64** |

Table 7. **Overlap Check for Point Cloud Registration.**

### 4.3. 3D Point Cloud Mosaicking

In this section, we aim at reconstructing the 3D scene from a set of partial point clouds observing parts of the scene. We proceed by selecting one of the point clouds as the origin and then estimating the absolute pose (*i.e.*, 3D translation and rotation) between the origin and each remaining point cloud in the set. One can imagine this procedure as 3D point cloud mosaicking. Please note that there is no guarantee that all partial clouds overlap with the one chosen as origin. While this could be alleviated by forming all possible pairs and forcing global consistency, solving the 3D mosaicking problem falls outside the scope of this paper. We only aim to demonstrate the potential of the proposed algorithm for this problem.

We perform the pairwise registration for all pairs using the method described in Section 3.3. In Table 8, we report results and compare with Geotransformer [32], using the same reconstruction paradigm. The evaluation metrics we use focus on the geometric aspects of accuracy and completeness [41] of the resulting reconstruction, as well as on precision, recall, and F1-score of registered point clouds. **SGAligner** has higher performance on 3 out of 5 metrics, and is mainly affected in completion and precision. The performance drop is due to the fact that for some scenes with low overlap, **SGAligner** fails to perform node alignment and, hence, registration with the incorrect alignments fails. In these scenes, GeoTr has a better complete context of the entire scene unlike the only object-based context in our approach. Furthermore, it is interesting to note that our K=1 method performs better than K=2. This is expected since there are more spurious matches from K=2 that lead to worse performance. We show qualitative results of success and failure in supp. mat.

### 4.4. Aligning 3D Scenes with Changes

In this section, we investigate the task of aligning a new 3D scene (target) on a prior 3D map (source), where the new scene can overlap fully or partially with the prior

| Methods | Acc ↓ | Comp ↓ | Precision ↑ | Recall ↑ | F1-Score ↑ |
|---|---|---|---|---|---|
| GeoTr [32] | 0.20721 | **0.03835** | **0.94180** | 0.79118 | 0.83667 |
| Ours K=1 | **0.00944** | 0.09347 | 0.90873 | **0.97444** | **0.93575** |
| Ours K=2 | 0.01215 | 0.10641 | 0.89234 | 0.97042 | 0.92345 |

Table 8. **Point cloud mosaicking from multiple point clouds.** SGAligner performs best in 3 out of the 5 metrics, with K = 1. The drop in the other 2 is due to a few low overlap pairs where node alignment fails, hence, registration too. *Best* values in **bold**.

map and may contain changes (both geometric and semantic). Specifically, we investigate the following scenarios: (i) aligning a *local* 3D scene on a *larger* prior map that contains no changes – here overlap of the local scene with the map is 100%; (ii) aligning a 3D scene on a prior map that contains changes; and (iii) aligning a *local* 3D scene on a *local* prior map that contains changes. We approach this as a 3D scene graph alignment task. For (i) and (ii), we use as large maps the 3D scene graphs of the entire scenes offered in [43, 44]. For local 3D scenes we use the data generated above. The results in Table 9 show that performance depends on the size of the prior map, whether there is full or partial overlap, and on the existence of temporal differences. They also demonstrate that our method can handle the alignment of maps that showcase temporal changes, even if not explicitly trained for this purpose.

| | Mean RR ↑ | Hits @ ↑ | | | | | No. of Pairs |
|---|---|---|---|---|---|---|---|
| | | K = 1 | K = 2 | K = 3 | K = 4 | K = 5 | |
| (i) | 0.970 | 0.952 | 0.976 | 0.989 | 0.993 | 0.995 | 827 |
| (ii) | 0.934 | 0.907 | 0.933 | 0.960 | 0.966 | 0.972 | 110 |
| (iii) | 0.886 | 0.833 | 0.894 | 0.928 | 0.946 | 0.957 | 2262 |

Table 9. **Alignment of a local 3D scene to a prior 3D map with differences in overlap and changes.**

## 5. Conclusion

We presented SGAligner, the first method capable of aligning 3D scene graphs directly on the graph level, that is robust to the in-the-wild scenarios, such as unknown overlap between scenes or changing environments. We demonstrated that aligning the scenes directly on the scene graph level can improve downstream tasks (*e.g.*, point cloud alignment) in terms of accuracy and speed. We believe our work could unlock agents to leverage this emerging scene representation for creating 3D maps of the environment, further using it for and sharing it with downstream tasks.

## 6. Acknowledgement

# Supplementary Material
# SGAligner: 3D Scene Alignment with Scene Graphs

## Abstract

*In the supplemental material, we provide additional details about the following:*

1. *Visualisation on point cloud mosaicking given multiple individual observations (Section A),*

2. *Comparison to a retrieval-based approach (Section B),*

3. *Additional ablation on SGAligner to further understand the performance of node matching (Section C),*

4. *Information on the SGAligner benchmark, including details on the generated data, evaluation protocol, and metrics (Section D), and*

5. *Details on implementation (Section E).*

## A. Application: Point Cloud Mosaicking

In Section 4.3 of the main paper, we demonstrate the potential of **SGAligner** on 3D point cloud mosaicking. Here, two success cases are shown in Figure 6 and a failure in Figure 7 (the graphs are shown simplified for visualisation purposes and do not represent the entire available graph).

## B. Application: Finding Overlapping Scenes

In the main paper, we discussed that **SGAligner** provides less than $O(N^2)$ computation complexity when addressing the task of registering multiple 3D scenes for which we have no knowledge of whether they overlap or not. Another approach to avoid full registration on all pairs (standard registration methods), is to use a retrieval-based approach. We consider the following approach as baseline: (i) extract local 3D keypoints for all available 3D point clouds [53]; (ii) generate a 3D descriptor per extracted keypoint [38]; (iii) accumulate the 3D keypoint descriptors into a global descriptor for each point cloud [18], and (iii) perform kNN search to rank global descriptors based on the queried one. Similarly here, this experiment serves as a demonstration of the potential of **SGAligner** and does not aim to solve the task.

Specifically, given a point cloud, we extract keypoints from the entire scene based purely on geometry and without any notion of object-ness or semantics. We randomly select 500 keypoints and their descriptors per scene, which

we use to train [18] so as to generate optimized global descriptors. During inference and given a query point cloud and its corresponding global descriptor, we perform a kNN search to identify the closest neighbors of it in the rest of the point clouds. We evaluate on Mean Reciprocal Rank (MRR) and compare with **SGAligner**.

Results are shown in Table 10. We evaluate on different point cloud densities, ranging from using the full point cloud density offered in [43] to random subsampling for 10, 20, 30, and 50%. Please note that in **SGAligner**, we do not use the entire scene, only objects in the scene graph. As described in the main paper, we downsample object point clouds using farthest point sampling to 512 points. We follow the same protocol here and perform this operation per subsampled scene level.

As expected, the retrieval-based method is performing well when there is a dense point cloud, since our method employs a limited amount of points per object instance. However, VLAD+KNN cannot retain a robust performance when density decreases, already reaching lower performance than **SGAligner** at 10% subsampling. In contrast, our approach is barely affected by a changing density since it already operates on lower-resolution point clouds. This showcases that the topological information encoded in 3D scene graphs can lead to more robust results when dealing with common failure cases in global descriptors (*i.e.*, changes in point cloud density).

| Subsampling % | VLAD + KNN | SGAligner |
|---|---|---|
| 0 | **<u>0.557</u>** | 0.383 |
| 10 | 0.316 | **0.356** |
| 20 | 0.276 | **0.343** |
| 30 | 0.222 | **0.339** |
| 50 | 0.162 | **0.312** |

Table 10. **Mean Reciprocal Rank (↑) comparison with a retrieval-based approach.** *Best* results per subsampling level are in **bold**. *Overall best* in **<u>underlined bold</u>**.

## C. Additional Ablation Studies

### C.1. Analysis with Various Object Encoders

In our experiments, we choose PointNet [31] as our encoder because it is commonly employed in most scene graph methods [45], [49], [52]. Pointnet has been shown to
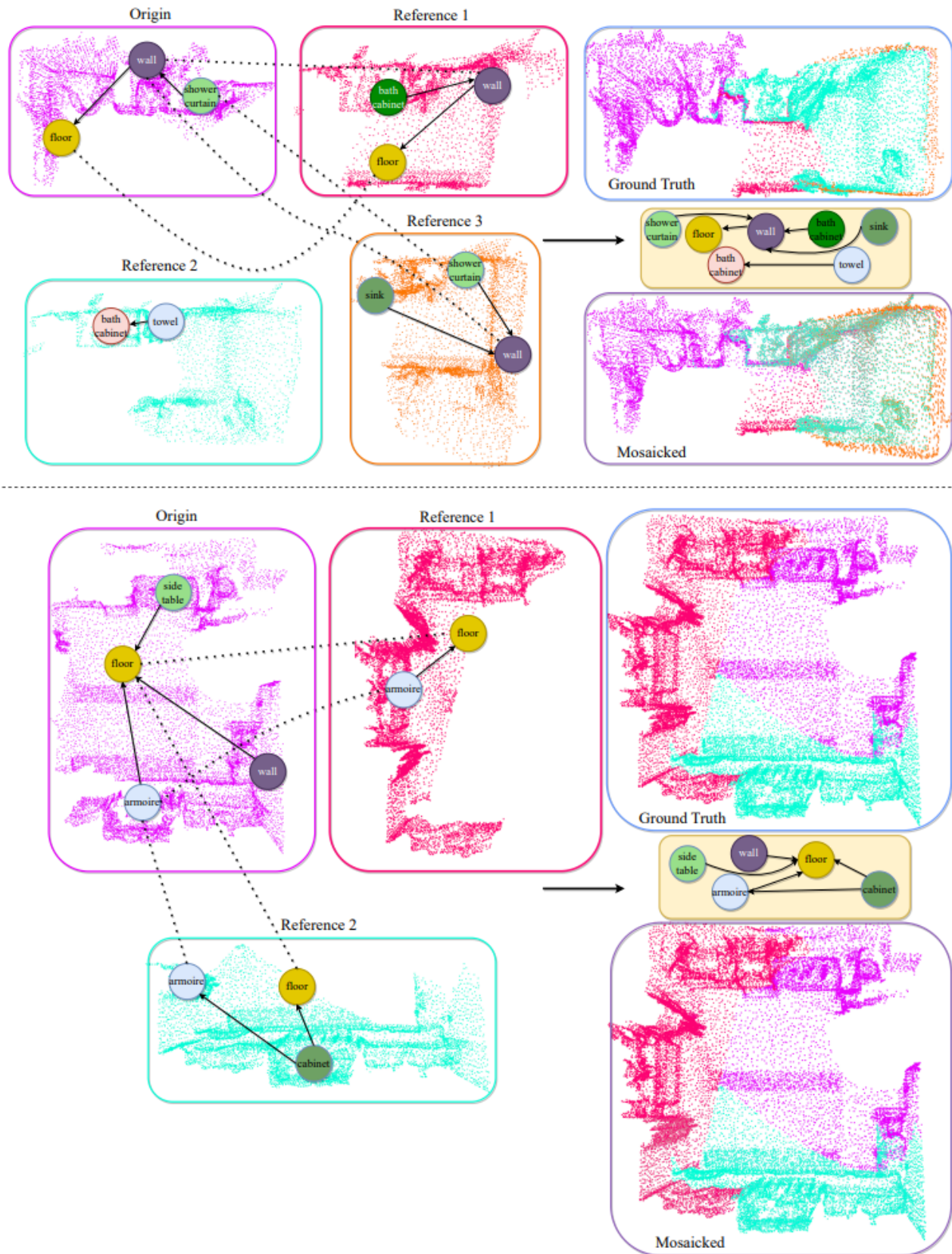
Figure 6. **Qualitative Results on Point Cloud Mosaicking.** Given partial point clouds of a scene and the corresponding 3D scene graphs, we showcase two example results on Point Cloud Mosaicking and the creation of a unified scene graph using the methodology discussed in Section A. The solid lines show the relationships between objects and dashed lines represent the ground truth entity pairs $\mathcal{F}$.
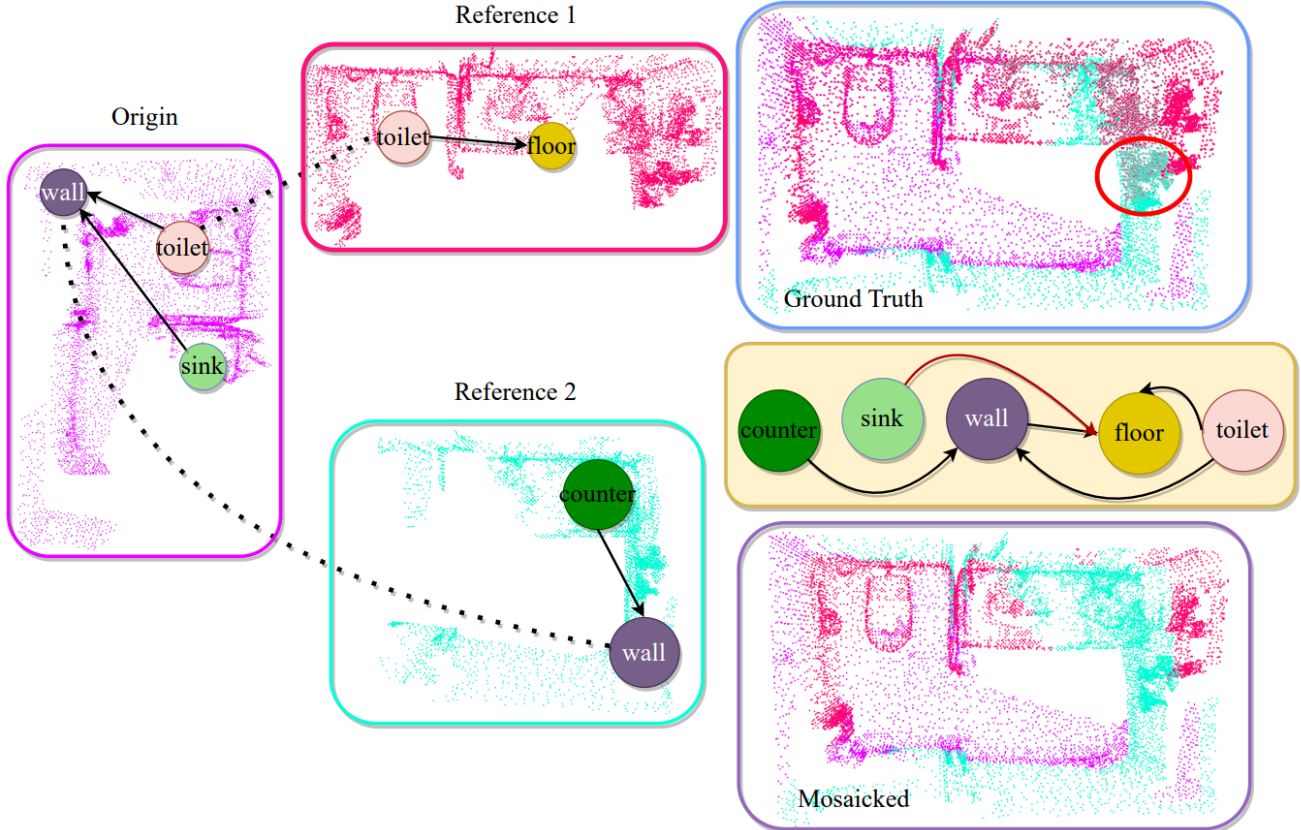
Figure 7. **Failure on Point Cloud Mosaicking.** Similar to the success cases in Fig. 6, we also showcase a failure case of our approach on point cloud mosaicking. The solid lines show the relationships between objects and dashed lines represent the ground truth entity pairs $\mathcal{F}$. The red circle on the ground truth point cloud depicts the area where the failure is the most visible and red arrow demonstrates the misalignment of **SGAligner** between a *sink* and a *floor*.

perform in real-time scenarios [3], which makes it suitable for mobile robot applications. In Table 11, we present a comparison of object encoders, on the node matching task. Point Cloud Transformer (PCT) [15], being inherently permutation invariant to an unordered point cloud, shows an improvement in the metrics. These results also showcase that our method is robust and agnostic to the 3D visual encoder.

| Encoder | Mean RR ↑ | Hits @ ↑ | | | | |
|---------|-----------|----------|---|---|---|---|
| | | K = 1 | K = 2 | K = 3 | K = 4 | K = 5 |
| PointNet [31] | 0.950 | 0.923 | 0.957 | 0.974 | 0.982 | 0.987 |
| PCT [15] | **0.965** | **0.947** | **0.968** | **0.983** | **0.988** | **0.991** |

Table 11. **Comparison on node matching of SGAligner using various object encoders.** Best values are in **bold**.

## C.2. Intra-Graph Alignment Recall

To further validate how our model performs on aligning nodes between two 3D scene graphs (source-reference) with no/partial overlap, we formulate *Intra-Graph Alignment Recall (IGAR)* metric. It measures what fraction of the nodes in the source graph are aligned (K=1), with nodes in the same source graph or, in other words, how many node matches out of total are self-aligned. We provide these results in Table 12. We do not explicitly model **not** self-matching nodes within the same graph, yet, *IGAR* values stand to show that our method rarely performs this.

$$IGAR = \frac{1}{M} \sum_{M} \frac{|pred\{n^i \equiv n^j\}|}{|\mathcal{F}|}, n \in \mathcal{N} \quad (4)$$

where, $i \neq j$, $pred\{n^i \equiv n^j\}$ is the set of nodes in the graph which **SGAligner** aligned with nodes in the same graph, $\mathcal{F}$ is the set of ground truth anchor pairs and $\mathcal{N}$ denotes the set of objects in a single graph and $M$ is the total number of graphs.

## C.3. Confusion Matrix

We compute a confusion matrix to identify which object categories are most frequently misaligned during entity alignment and if our method fails on certain semantic classes (*e.g.*, *chair*, *table*, etc). In Figure 8, we show the confusion matrix on all 4 module combinations of
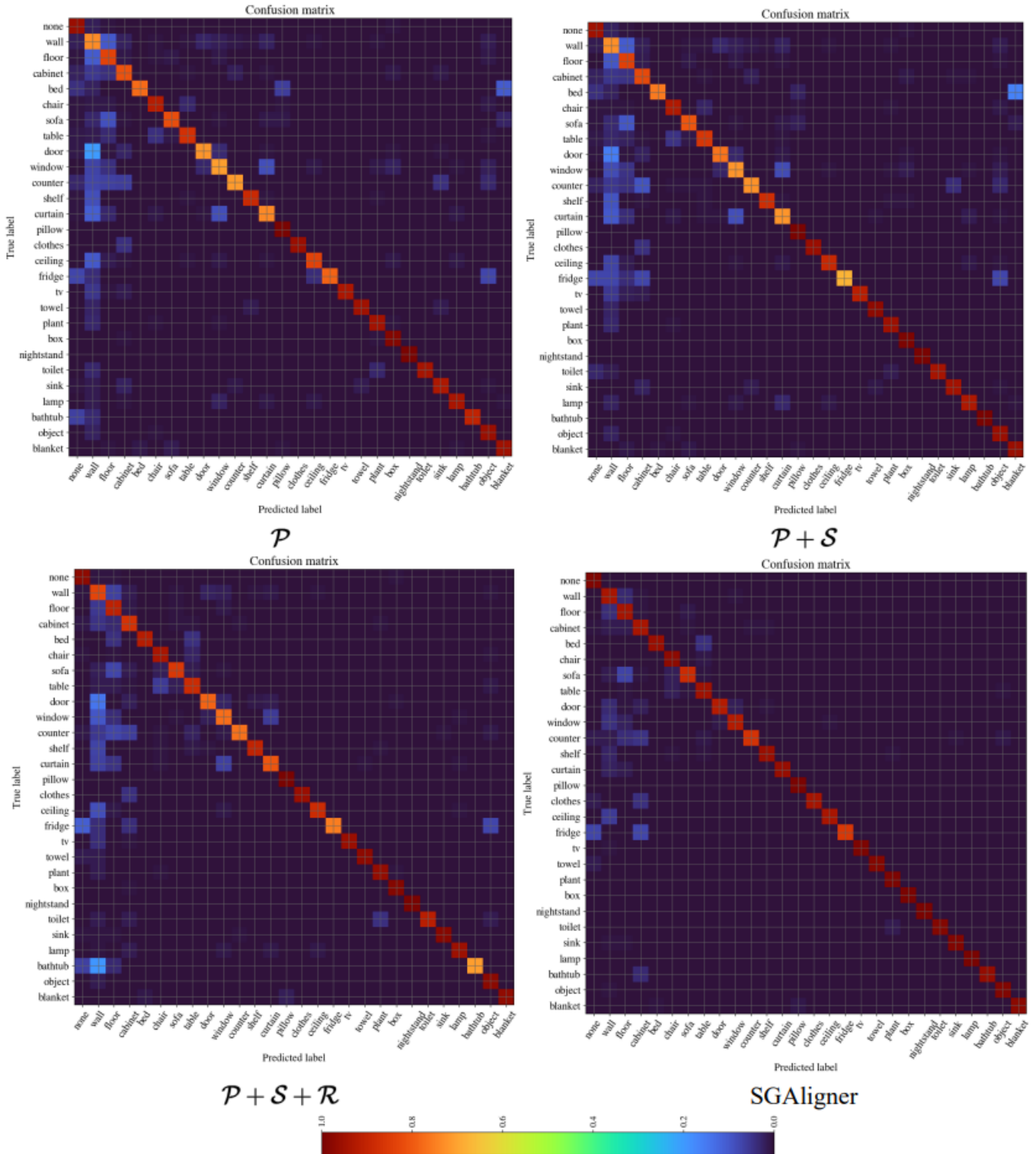
Figure 8. **Object Confusion Matrices** of the 4 module combinations of **SGAligner**: object encoder ($\mathcal{P}$), object and structure encoders ($\mathcal{P} + \mathcal{S}$), object, structure and relationship encoders ($\mathcal{P} + \mathcal{S} + \mathcal{R}$), and the proposed method with all modules (**SGAligner**). High values indicate that an object (denoted on $y$-axis) is often recognized as the object denoted on $x$-axis – everything but the diagonal should be 0.

| Method | IGAR ↓ (%) |
|---|---|
| $\mathcal{P}$ | 16.9 |
| $\mathcal{P} + \mathcal{S}$ | 16.5 |
| $\mathcal{P} + \mathcal{S} + \mathcal{R}$ | 13.1 |
| **SGAligner** | **8.2** |

Table 12. **Evaluation on node self-alignment. SGAligner** has not been explicitly modeled to not create self-matches but still is able to differentiate between nodes from the same and different graphs.

**SGAligner.** As expected, the object encoder module $\mathcal{P}$, although performing well, confuses the most the *wall* and *floor* classes. This is due to the fact that purely on a semantic level, without encoding any positional/structural information, these classes are similar. We can further observe that **SGAligner** is robust to certain classes like *pillow*, *tv*, *lamp*, etc. Classes like *wall*, *floor*, and *fridge* are the ones easily susceptible to misalignment on the tested dataset, albeit less than in $\mathcal{P}$.

## C.4. Robustness to Missing Geometric Information

In this section, we provide an ablation study of **SGAligner** on the 3D Scene Graph alignment task and evaluate how it performs on node alignment, when all the **geometric relationships** encoding positional information between the nodes such as `left` and `standing on` are missing. Results are in Table 13. As expected, the structure module $\mathcal{S}$ suffers from this compared to the full ground-truth experiment, since the number of edges encoded in the neighborhood of an entity gets reduced. However, overall, our method does not show a drastic drop in node alignment metrics due to the absence of geometric relationships. This can be attributed to the fact that we do not discriminate between different types of relationships in our encoders, however, this is a very important robustness characteristic, especially, while working with predicted scene graphs where the relationships could be missing or incorrectly labelled. This also shows that once trained with full ground truth, our method is able to handle missing data during inference which would be useful for a navigation agent.

| Modalities | Mean RR ↑ | Hits @ ↑ | | | | |
|---|---|---|---|---|---|---|
| | | K = 1 | K = 2 | K = 3 | K = 4 | K = 5 |
| $\mathcal{P}$ | 0.884 | 0.835 | 0.886 | 0.921 | 0.938 | 0.951 |
| $\mathcal{P} + \mathcal{S}$ | 0.880 | 0.830 | 0.882 | 0.918 | 0.936 | 0.948 |
| $\mathcal{P} + \mathcal{S} + \mathcal{R}$ | 0.893 | 0.844 | 0.898 | 0.933 | 0.949 | 0.959 |
| **SGAligner** | **0.948** | **0.921** | **0.952** | **0.971** | **0.979** | **0.985** |

Table 13. **Evaluation on node matching.** We compare the performance of **SGAligner** for different modality combinations, when all geometric edges are missing.

## D. SGAligner Benchmark

In this section, we offer qualitative explanations on our dataset generation procedure and discuss the evaluation metrics used to asses performance with respect to the various tasks we reported.
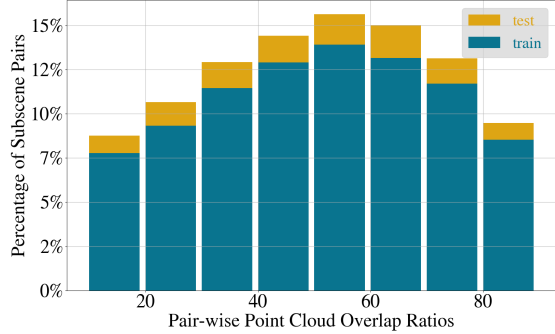


Figure 9. **Overlap statistics for the generated sub-scenes.**

## D.1. Dataset

In Sec. 4 of the main paper, we provide a description of the data generation procedure. In Figure 9, we report statistics on the spatial overlap of the generated pairs. We show examples of sub-scenes generated using this approach in Figure 10, alongside the camera trajectory used to capture the corresponding scan in [43]. We will make our code and benchmark public.

## D.2. Evaluation Metrics

The evaluation metrics that we use to assess performance in Section 4.1 of the main paper, as well as in this supplementary material, are formally defined in this section.

### D.2.1 Alignment Metrics

Inspired by works in multi-modal entity alignment [25], we define our alignment metrics as follows :

**Hits @ K** represents the fraction of true anchor entities present in the top k predictions :

$$H_k(r_1, ..., r_n) = \frac{1}{n} \sum_{i=1}^{n} I[r_i \leq k] \qquad (5)$$

where, $I[x \leq y] = 1$ when $x \leq y$ else 0 and $k \in [1, 2, 3, 4, 5]$.

**Mean Reciprocal Rank (MRR)** corresponds to the arithmetic mean over the reciprocals of ranks of true triples.

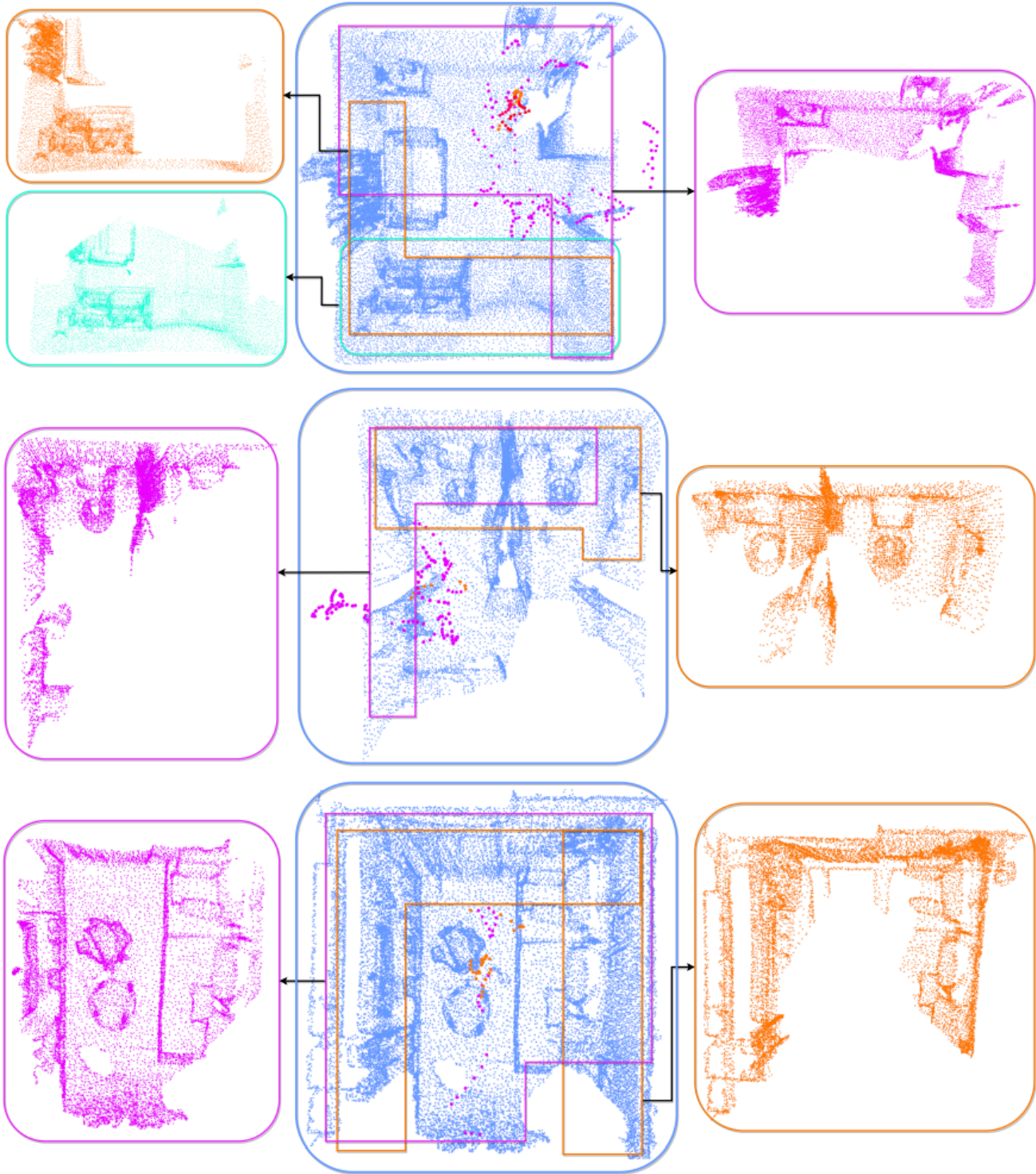$$MRR(r_1, ..., r_n) = \frac{1}{n} \sum_{i=1}^{n} r_i^{-1} \qquad (6)$$

Figure 10. **Visualization of data generation process.** We visualise the creation of subscenes using our approach from Section 4 of the main paper. Given a point cloud from [43] (middle column), we create multiple sub-scenes (left and right columns) and showcase the used camera trajectory. The colors depict the camera trajectory and 3D spatial area of each sub-scene in the parent scene.

### D.2.2 Registration Metrics

**Feature Matching Recall (FMR)** [16][32] measures the fraction of point cloud pairs for which, based on the number

of inlier correspondences, it is likely that accurate transformation parameters can be recovered with a robust estimator such as RANSAC. It should be noted that FMR simply verifies whether the inlier ratio (IR) is higher than a threshold $\mathcal{T} = 0.05$. It does not examine if the transformation can actually be inferred from those correspondences, which is not always the case because of the possibility that their geometric arrangement is (almost) degenerate, such as when they are situated closely together or along a straight edge.

$$FMR = \frac{1}{M} \sum_{i=1}^{M} [\![ IR_i > \mathcal{T} ]\!] \qquad (7)$$

$$MRR(r_1, ..., r_n) = \frac{1}{n} \sum_{i=1}^{n} r_i^{-1} \qquad (8)$$

where $M$ is the number of all point cloud pairs.

**Registration Recall (RR)** is the fraction of registered point cloud pairs for which the transformation error is smaller than 0.2m. The transformation error is the root mean squared error of the ground truth correspondence $\mathcal{H}^*$ after applying the predicted transformation $\mathbf{T}_{P \to Q}$.

$$RMSE = \sqrt{\frac{1}{|\mathcal{H}^*|} \sum_{(p^*_{x_i}, q^*_{y_i}) \in \mathcal{H}^*} ||T_{P \to Q}(p^*_{x_i}) - q^*_{y_i}||_2^2} \qquad (9)$$

$$RR = \frac{1}{M} \sum_{i=1}^{M} [\![ RMSE_i < 0.2m ]\!] \qquad (10)$$

**Relative Rotation Error (RRE)** is the geodesic distance in degrees between estimated and ground-truth rotation matrices.

$$RRE = arccos(\frac{trace(R^T \cdot \bar{R} - 1)}{2}) \qquad (11)$$

**Relative Translation Error (RTE)** is the the euclidean distance between estimated and ground-truth translation vectors.

$$RTE = ||t - \bar{t}|| \qquad (12)$$

We compute mean RRE and RTE between all the registered point cloud pairs.

**Chamfer Distance** measures the quality of registration. Following [46], [16], we use the *modified* Chamfer distance metric :

$$\bar{C}D(P, Q) = \frac{1}{|P|} \sum_{p \in P} min_{q \in Q_{raw}} ||T_P{}^Q(p) - q||_2^2 + $$
$$\frac{1}{|Q|} \sum_{q \in Q} min_{p \in P_{raw}} ||q - T_P{}^Q(p)||_2^2 \qquad (13)$$

where, $P_{raw}$ and $Q_{raw}$ are $raw/clean$ source and target point clouds respectively.

### D.2.3   Reconstruction Metrics

The definition of full 3D reconstruction metrics is provided in Table 14.

| Metric | Definition |
|---|---|
| Acc | $mean_{p \in P}(min_{p^* \in P^*}||p - p^*||)$ |
| Comp | $mean_{p^* \in P^*}(min_{p \in P}||p - p^*||)$ |
| Precision | $mean_{p \in P}(min_{p^* \in P^*}||p - p^*|| < 0.05)$ |
| Recall | $mean_{p^* \in P^*}(min_{p \in P}||p - p^*|| < 0.05)$ |
| F1-Score | $\frac{2*precision*recall}{precision+recall}$ |

Table 14. **3D Reconstruction Metric Definitions.** $p$ and $p^*$ are ground truth and predicted point clouds respectively.

## E. Implementation Details

Inspired by MCLEA [25], we use a multi-layered GAT with 2 layers and each hidden unit being 128-dimensional. All the modules output a 100-dimensional embedding and the joint embedding, being a weighted concatenation, is 400-dimensional. We use $\mathcal{T}_1$ for ICL loss as 0.1 and $\mathcal{T}_2$ for IAL loss as 1.0. We train our model for 50 epochs on a NVIDIA GeForce RTX 3060 Ti 8GB GPU with a batch size of 4 using AdamW [30] optimizer and a learning rate of 0.001.

## References

[1] Christopher Agia, Krishna Murthy Jatavallabhula, Mohamed Khodeir, Ondrej Miksik, Vibhav Vineet, Mustafa Mukadam, Liam Paull, and Florian Shkurti. Taskography: Evaluating robot task planning over large 3d scene graphs. In *Conference on Robot Learning*, pages 46–58. PMLR, 2022. 1, 2

[2] Iro Armeni, Zhi-Yang He, JunYoung Gwak, Amir R Zamir, Martin Fischer, Jitendra Malik, and Silvio Savarese. 3d scene graph: A structure for unified semantics, 3d space, and camera. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5664–5673, 2019. 1, 2

[3] Lin Bai, Yecheng Lyu, and Xinming Huang. Pointnet on fpga for real-time lidar point cloud processing. 10 2020. 12

[4] Xuyang Bai, Zixin Luo, Lei Zhou, Hongbo Fu, Long Quan, and Chiew-Lan Tai. D3feat: Joint learning of dense detection and description of 3d local features. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6359–6367, 2020. 2, 3

[5] Daniel Barath and Jiří Matas. Graph-cut ransac. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6733–6741, 2018. 5

[6] Daniel Barath, Jana Noskova, Maksym Ivashechkin, and Jiri Matas. Magsac++, a fast, reliable and accurate robust estimator. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1304–1312, 2020. 5

[7] Yun Chang, Luca Ballotta, and Luca Carlone. D-lite: Navigation-oriented compression of 3d scene graphs

under communication constraints. *arXiv preprint arXiv:2209.06111*, 2022. 1, 2

[8] Liyi Chen, Zhi Li, Yijun Wang, Tong Xu, Zhefeng Wang, and Enhong Chen. Mmea: entity alignment for multi-modal knowledge graph. In *Knowledge Science, Engineering and Management: 13th International Conference, KSEM 2020, Hangzhou, China, August 28–30, 2020, Proceedings, Part I 13*, pages 134–147. Springer, 2020. 1, 2

[9] Liyi Chen, Zhi Li, Tong Xu, Han Wu, Zhefeng Wang, Nicholas Jing Yuan, and Enhong Chen. Multi-modal siamese network for entity alignment. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 118–126, 2022. 2

[10] Bo Cheng, Jia Zhu, and Meimei Guo. Multijaf: Multi-modal joint entity alignment framework for multi-modal knowledge graph. *Neurocomputing*, 500:581–591, 2022. 1, 2

[11] Helisa Dhamo, Fabian Manhardt, Nassir Navab, and Federico Tombari. Graph-to-3d: End-to-end generation and manipulation of 3d scenes using scene graphs. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16352–16361, 2021. 2

[12] Martin A. Fischler and Robert C. Bolles. Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography. *Communication of ACM*, 1981. 2, 3, 5

[13] Samir Yitzhak Gadre, Kiana Ehsani, Shuran Song, and Roozbeh Mottaghi. Continuous scene representations for embodied ai. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14849–14859, June 2022. 1, 2

[14] Hao Guo, Jiuyang Tang, Weixin Zeng, Xiang Zhao, and Li Liu. Multi-modal entity alignment in hyperbolic space. *Neurocomputing*, 461:598–607, 2021. 1, 2

[15] Meng-Hao Guo, Jun-Xiong Cai, Zheng-Ning Liu, Tai-Jiang Mu, Ralph R. Martin, and Shi-Min Hu. Pct: Point cloud transformer. *Computational Visual Media*, 7(2):187–199, Apr 2021. 4, 12

[16] Shengyu Huang, Zan Gojcic, Mikhail Usvyatsov, Andreas Wieser, and Konrad Schindler. Predator: Registration of 3d point clouds with low overlap. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021. 2, 3, 15, 16

[17] N. Hughes, Y. Chang, and L. Carlone. Hydra: A real-time spatial perception system for 3D scene graph construction and optimization. 2022. 2, 3

[18] Hervé Jégou, Matthijs Douze, Cordelia Schmid, and Patrick Pérez. Aggregating local descriptors into a compact image representation. In *2010 IEEE computer society conference on computer vision and pattern recognition*, pages 3304–3311. IEEE, 2010. 10

[19] Ziyuan Jiao, Yida Niu, Zeyu Zhang, Song-Chun Zhu, Yixin Zhu, and Hangxin Liu. Sequential manipulation planning on scene graph. In *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 8203–8210. IEEE, 2022. 1, 2

[20] Ue-Hwan Kim, Jin-Man Park, Taek-Jin Song, and Jong-Hwan Kim. 3-d scene graph: A sparse and semantic representation of physical environments for intelligent agents.

[21] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123:32–73, 2017. 2

[22] Xinghang Li, Di Guo, Huaping Liu, and Fuchun Sun. Embodied semantic scene graph generation. In *Conference on Robot Learning*, pages 1585–1594. PMLR, 2022. 1, 2

[23] Yujia Li, Chenjie Gu, Thomas Dullien, Oriol Vinyals, and Pushmeet Kohli. Graph matching networks for learning the similarity of graph structured objects. In *International conference on machine learning*, pages 3835–3845. PMLR, 2019. 2

[24] Yongwei Li, Yalong Ma, Xiang Huo, and Xinkai Wu. Remote object navigation for service robots using hierarchical knowledge graph in human-centered environments. *Intelligent Service Robotics*, 15(4):459–473, 2022. 1, 2

[25] Zhenxi Lin, Ziheng Zhang, Meng Wang, Yinghui Shi, Xian Wu, and Yefeng Zheng. Multi-modal contrastive representation learning for entity alignment. *arXiv preprint arXiv:2209.00891*, 2022. 1, 2, 3, 4, 5, 14, 16

[26] Fangyu Liu, Muhao Chen, Dan Roth, and Nigel Collier. Visual pivoting for (unsupervised) entity alignment. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 4257–4266, 2021. 1, 2, 6

[27] Fangyu Liu, Rongtian Ye, Xun Wang, and Shuaipeng Li. Hal: Improved text-image matching by mitigating visual semantic hubs. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34:11563–11571, 04 2020. 2

[28] Ye Liu, Hui Li, Alberto Garcia-Duran, Mathias Niepert, Daniel Onoro-Rubio, and David S Rosenblum. Mmkg: multi-modal knowledge graphs. In *The Semantic Web: 16th International Conference, ESWC 2019, Portorož, Slovenia, June 2–6, 2019, Proceedings 16*, pages 459–474. Springer, 2019. 1, 2

[29] Samuel Looper, Javier Rodriguez-Puigvert, Roland Siegwart, Cesar Cadena, and Lukas Schmid. 3d vsg: Long-term semantic scene change prediction through 3d variable scene graphs. *arXiv preprint arXiv:2209.07896*, 2022. 2

[30] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. 16

[31] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 652–660, 2017. 4, 5, 6, 10, 12

[32] Zheng Qin, Hao Yu, Changjian Wang, Yulan Guo, Yuxing Peng, and Kai Xu. Geometric transformer for fast and robust point cloud registration. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022. 2, 3, 5, 7, 8, 9, 15

[33] Rahul Raguram, Ondrej Chum, Marc Pollefeys, Jiri Matas, and Jan-Michael Frahm. Usac: A universal framework for

*IEEE transactions on cybernetics*, 50(12):4921–4933, 2019. 1, 2

random sample consensus. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):2022–2038, 2012. 5

[34] Zachary Ravichandran, J Daniel Griffith, Benjamin Smith, and Costas Frost. Bridging scene understanding and task execution with flexible simulation environments. *arXiv preprint arXiv:2011.10452*, 2020. 1, 2

[35] Zachary Ravichandran, Lisa Peng, Nathan Hughes, J. Daniel Griffith, and Luca Carlone. Hierarchical representations and explicit memory: Learning effective navigation policies on 3d scene graphs using graph neural networks. In *2022 International Conference on Robotics and Automation (ICRA)*, pages 9272–9279, 2022. 1, 2

[36] Antoni Rosinol, Arjun Gupta, Marcus Abate, Jingnan Shi, and Luca Carlone. 3d dynamic scene graphs: Actionable spatial perception with places, objects, and humans. *arXiv preprint arXiv:2002.06289*, 2020. 1, 2

[37] Antoni Rosinol, Andrew Violette, Marcus Abate, Nathan Hughes, Yun Chang, Jingnan Shi, Arjun Gupta, and Luca Carlone. Kimera: From slam to spatial perception with 3d dynamic scene graphs. *The International Journal of Robotics Research*, 40(12-14):1510–1546, 2021. 1, 2

[38] Radu Bogdan Rusu, Nico Blodow, and Michael Beetz. Fast point feature histograms (fpfh) for 3d registration. In *2009 IEEE international conference on robotics and automation*, pages 3212–3217. IEEE, 2009. 10

[39] Gabriel Sepulveda, Juan Carlos Niebles, and Alvaro Soto. A deep learning based behavioral approach to indoor autonomous navigation. In *2018 IEEE international conference on robotics and automation (ICRA)*, pages 4646–4653. IEEE, 2018. 1, 2

[40] Yinghui Shi, Meng Wang, Ziheng Zhang, Zhenxi Lin, and Yefeng Zheng. Probing the impacts of visual context in multimodal entity alignment. In *Web and Big Data: 6th International Joint Conference, APWeb-WAIM 2022, Nanjing, China, November 25–27, 2022, Proceedings, Part II*, pages 255–270. Springer, 2023. 1, 2

[41] Jiaming Sun, Yiming Xie, Linghao Chen, Xiaowei Zhou, and Hujun Bao. NeuralRecon: Real-time coherent 3D reconstruction from monocular video. *CVPR*, 2021. 9

[42] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. In *International Conference on Learning Representations*, 2018. 4

[43] Johanna Wald, Armen Avetisyan, Nassir Navab, Federico Tombari, and Matthias Nießner. Rio: 3d object instance relocalization in changing indoor environments. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7658–7667, 2019. 2, 5, 9, 10, 14, 15

[44] Johanna Wald, Helisa Dhamo, Nassir Navab, and Federico Tombari. Learning 3d semantic scene graphs from 3d indoor reconstructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3961–3970, 2020. 1, 2, 3, 5, 6, 9

[45] Shun-Cheng Wu, Johanna Wald, Keisuke Tateno, Nassir Navab, and Federico Tombari. Scenegraphfusion: Incremental 3d scene graph prediction from rgb-d sequences. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7515–7525, 2021. 2, 8, 10

[46] Zi Jian Yew and Gim Hee Lee. Rpm-net: Robust point matching using learned features. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11824–11833, 2020. 7, 16

[47] Zi Jian Yew and Gim Hee Lee. Regtr: End-to-end point cloud correspondences with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022. 2, 3

[48] Andy Zeng, Shuran Song, Matthias Nießner, Matthew Fisher, Jianxiong Xiao, and Thomas Funkhouser. 3dmatch: Learning local geometric descriptors from rgb-d reconstructions. In *CVPR*, 2017. 5, 7

[49] Chaoyi Zhang, Jianhui Yu, Yang Song, and Weidong Cai. Exploiting edge-oriented reasoning for 3d point-based scene graph analysis. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9705–9715, June 2021. 2, 10

[50] Qingheng Zhang, Zequn Sun, Wei Hu, Muhao Chen, Lingbing Guo, and Yuzhong Qu. Multi-view knowledge graph embedding for entity alignment. *arXiv preprint arXiv:1906.02390*, 2019. 1, 2

[51] Ruohan Zhang, Dhruva Bansal, Yilun Hao, Ayano Hiranaka, Jialu Gao, Chen Wang, Roberto Martín-Martín, Li Fei-Fei, and Jiajun Wu. A dual representation framework for robot learning with human guidance. In *6th Annual Conference on Robot Learning*. 1

[52] Shoulong Zhang, Shuai Li, Aimin Hao, and Hong Qin. Knowledge-inspired 3d scene graph prediction in point cloud. In Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan, editors, *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 18620–18632, 2021. 2, 10

[53] Yu Zhong. Intrinsic shape signatures: A shape descriptor for 3d object recognition. In *2009 IEEE 12th international conference on computer vision workshops, ICCV workshops*, pages 689–696. IEEE, 2009. 10