

ConceptGraphs: Open-Vocabulary 3D Scene Graphs for Perception and Planning

Anonymous Author(s)

Affiliation

Address

email

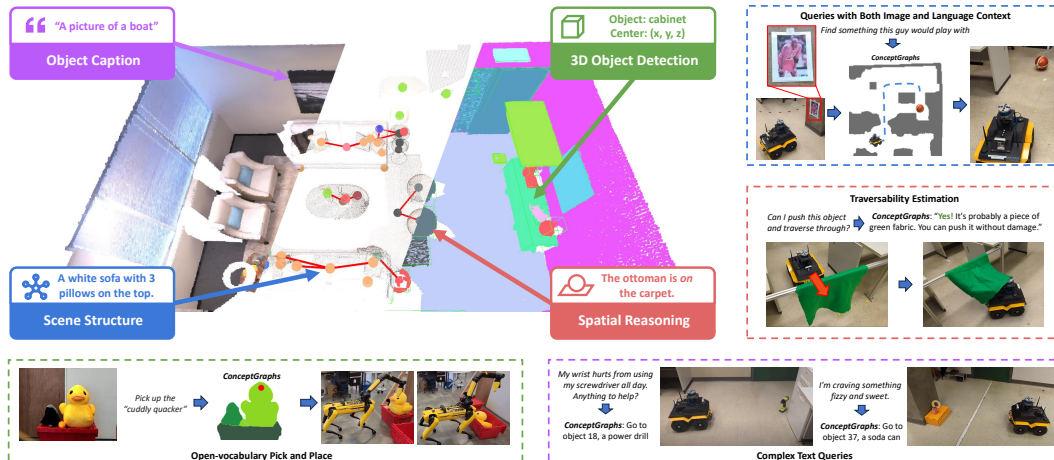


Figure 1: *ConceptGraphs* builds open-vocabulary 3D scene graphs. We design an object-based mapping system that (a) only assumes class-agnostic instance masks and fuses them to 3D, (b) extracts language tags for each mapped instance leveraging large vision-language models, and (c) builds a graph of object spatial relations. The object-centric nature of *ConceptGraphs* allows easy map maintenance and promotes scalability, and the graph structure provides relational information within the scene. Furthermore, our scene graph representations are easily mapped to natural language formats to interface with LLMs, enabling them to answer complex scene queries and granting robots access to useful facts about surrounding objects, such as traversability and utility. We implement and demonstrate *ConceptGraphs* on a number of real-world robotics tasks across wheeled and legged mobile robot platforms. ([Anonymized project page](#))

1 **Abstract:** For robots to perform a wide variety of tasks, they require a 3D repre-
2 sentation of the world that is semantically rich, yet compact and efficient for task-
3 driven perception and planning. Recent approaches have attempted to leverage
4 features from large vision-language models to encode semantics in 3D represen-
5 tations. However, these approaches tend to produce maps with per-point feature
6 vectors, which do not scale well in larger environments, nor do they contain se-
7 mantic spatial relationships between entities in the environment, which are useful
8 for downstream planning. In this work, we propose *ConceptGraphs*, an open-
9 vocabulary graph-structured representation for 3D scenes. *ConceptGraphs* is built
10 by leveraging 2D foundation models and fusing their output to 3D by multi-view
11 association. The resulting representations generalize to novel semantic classes,
12 without the need to collect large 3D datasets or finetune models. We demon-
13 strate the utility of this representation through a number of downstream planning
14 tasks that are specified through abstract (language) prompts and require complex
15 reasoning over spatial and semantic concepts. To explore the full scope of our
16 experiments and results, we encourage visiting our [anonymized project webpage](#).

17 1 Introduction

18 Scene representation is one of the key design choices that can facilitate downstream planning for a
19 variety of tasks, including mobility and manipulation. Robots need to build these representations
20 online from onboard sensors as they navigate through an environment. For efficient execution of
21 complex tasks such representations should be: *scalable and efficient to maintain*, as the volume of
22 the scene and the duration of the robot’s operation increases; *open-vocabulary*, not limited to making
23 inferences about a set of concepts that is predefined at training time, but capable of handling new
24 objects and concepts at inference time; and *with a flexible level of detail* to enable planning over a
25 range of tasks, from ones that require dense geometric information for mobility and manipulation,
26 to ones that need abstract semantic information and object-level affordance information for task
27 planning. We propose *ConceptGraphs*, a 3D scene representation method for robot perception and
28 planning that satisfies all the above requirements.

29 **Closed-vocabulary semantic mapping in 3D.** Early works reconstruct the 3D map through online
30 algorithms like simultaneous localization and mapping (SLAM) [1, 2, 3, 4, 5] or offline methods
31 like structure-from-motion (SfM) [6, 7]. Aside from reconstructing 3D geometry, recent works
32 also use deep learning-based object detection and segmentation models to reconstruct the 3D scene
33 representations with dense semantic mapping [8, 9, 10, 11] or object-level decomposition [12, 13,
34 14, 15]. While these methods achieve impressive results in mapping semantic information to 3D,
35 they are closed-vocabulary and their applicability is limited to object categories annotated in their
36 training datasets.

37 **3D scene representations using foundation models.** There have been significant recent efforts [16,
38 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30] focused on building 3D representations by
39 leveraging *foundation models* - large, powerful models that capture a diverse set of concepts and
40 accomplish a wide range of tasks [31, 32, 33, 34, 35]. Such models have excelled in tackling open-
41 vocabulary challenges in 2D vision. However, they require an “internet-scale” of training data, and
42 no 3D datasets exist yet of a comparable size. Recent works have therefore attempted to *ground*
43 the 2D representations produced by image and language foundation models to the 3D world and
44 show impressive results on open-vocabulary tasks, including language-guided object grounding [17,
45 18, 24, 26, 36], 3D reasoning [37, 38], robot manipulation [39, 40] and navigation [41, 42]. These
46 approaches project dense per-pixel features from images to 3D to build explicit representations such
47 as pointclouds [17, 18, 19, 20, 21] or implicit neural representations [16, 22, 23, 24, 25, 26, 27, 28,
48 29, 30].

49 However, such methods have two key limitations. First, assigning every point a semantic feature
50 vector is highly redundant and consumes more memory than necessary, greatly limiting scalability
51 to large scenes. Second, these dense representations do not admit an easy decomposition – this lack
52 of structure makes them less amenable to dynamic updates to the map (crucial for robotics).

53 **3D scene graphs.** 3D scene graphs (3DSGs) address the second limitation by compactly and effi-
54 ciently describing scenes with graph structures, with nodes representing objects and edges encoding
55 inter-object relationships [43, 44, 45, 46, 47]. These approaches have enabled building real-time
56 systems that can dynamically build up hierarchical 3D scene representations [48, 49, 50], and more
57 recently shown that various robotics planning tasks can benefit from efficiency and compactness of
58 3DSGs [51, 52]. However, existing work on building 3D scene graphs has been confined to the
59 closed-vocabulary setting, limiting their applicability to a small set of tasks.

60 1.1 Overview of Our Contribution

61 In this work, we mitigate all the aforementioned limitations and propose *ConceptGraphs*, an open-
62 vocabulary and object-centric 3D representation for robot perception and planning. In *Concept-*
63 *Graphs*, each object is represented as a node with geometric and semantic features, and relation-
64 ships among objects are encoded in the graph edges. At the core of *ConceptGraphs* is an object-
65 centric 3D mapping technique that integrates geometric cues from conventional 3D mapping sys-

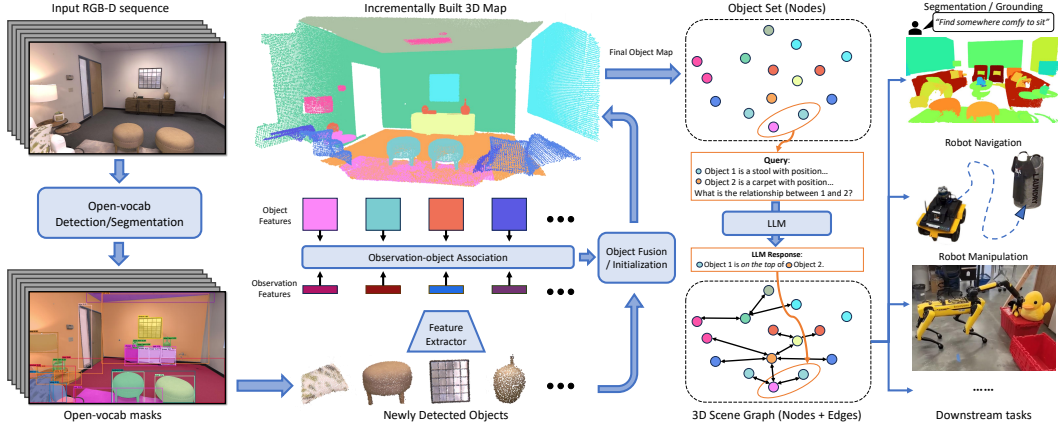


Figure 2: *ConceptGraphs* builds an open-vocabulary 3D scene graph from a sequence of posed RGB-D images. We use generic instance segmentation models to segment regions from RGB images, extract semantic feature vectors for each, and project them to a 3D point cloud. These regions are incrementally associated and fused from multiple views, resulting in a set of 3D objects and associated vision (and language) descriptors. Then large vision and language models are used to caption each mapped 3D objects and derive inter-object relations, which generates the edges to connect the set of objects and form a graph. The resulting 3D scene graph provides a structured and comprehensive understanding of the scene and can further be easily translated to a text description, useful for LLM-based task planning.

66 tems, and semantic cues from vision and language foundation models [31, 33, 34, 53, 54, 55].
 67 Objects are assigned language tags by leveraging large language models (LLMs) [32] and large
 68 vision-language models (LVLMs) [55], which provide semantically rich descriptions and enable
 69 free-form language querying, all while using off-the-shelf models (no training/finetuning). The 3D
 70 scene graph [43, 44, 45, 46, 47] structure allows us to efficiently represent large scenes with a low
 71 memory footprint and makes for efficient task planning.

72 In experiments, we demonstrate that *ConceptGraphs* is able to discover, map, and caption a large
 73 number of objects in a scene. Further, we conduct real-world trials on multiple robot platforms
 74 over a wide range of downstream tasks, including manipulation, navigation, localization, and map
 75 updates. To summarize, our key **contributions** are:

- 76 • We propose a novel object-centric mapping system that integrates geometric cues from tra-
 77 ditional 3D mapping systems and semantic cues from 2D foundation models.
- 78 • We construct open-vocabulary 3D scene graphs; efficient and structured semantic abstrac-
 79 tions for perception and planning.
- 80 • We implement *ConceptGraphs* on real-world wheeled and legged robotic platforms,
 81 demonstrating downstream perception and planning for complex/abstract language queries.

82 2 Method

83 *ConceptGraphs* builds a compact, semantically rich representation of a 3D environment. Given
 84 a set of posed RGB-D frames, we run a class-agnostic segmentation model to obtain candidate
 85 objects, associate them across multiple views using geometric and semantic similarity measures,
 86 and instantiate nodes in a 3D scene graph. We then use an LVLM to caption each node and an
 87 LLM to infer relationships between adjoining nodes, which results in edges in the scene graph.
 88 This resultant scene graph is open-vocabulary, encapsulates object properties, and can be used for a
 89 multitude of downstream tasks including segmentation, object grounding, navigation, manipulation,
 90 localization, and remapping. The approach is illustrated in Fig. 2. For implementation details, see
 91 appendix

92 2.1 Object-based 3D Mapping

93 **Object-centric 3D representation:** Given a sequence of RGB-D observations $\mathcal{I} = \{I_1, I_2, \dots, I_t\}$,
94 *ConceptGraphs* constructs a map, a 3D scene graph, $\mathcal{M}_t = \langle \mathbf{O}_t, \mathbf{E}_t \rangle$, where $\mathbf{O}_t = \{\mathbf{o}_j\}_{j=1\dots J}$ and
95 $\mathbf{E}_t = \{\mathbf{e}_k\}_{k=1\dots K}$ represent the sets of objects and edges, respectively. Each object \mathbf{o}_j is charac-
96 terized by a 3D point cloud $\mathbf{p}_{\mathbf{o}_j}$ and a semantic feature vector $\mathbf{f}_{\mathbf{o}_j}$. This map is built incrementally,
97 incorporating each incoming frame $I_t = \langle I_t^{\text{rgb}}, I_t^{\text{depth}}, \theta_t \rangle$ (color image, depth image, pose) into the
98 existing object set \mathbf{O}_{t-1} , by either adding to existing objects or instantiating new ones.

99 **Class-agnostic 2D Segmentation:** When processing frame I_t , a class-agnostic segmentation model
100 $\text{Seg}(\cdot)$ is used to obtain a set of masks $\{\mathbf{m}_{t,i}\}_{i=1\dots M} = \text{Seg}(I_t^{\text{rgb}})$ corresponding to candidate ob-
101 jects¹. Each extracted mask $\mathbf{m}_{t,i}$ is then passed to a visual feature extractor (CLIP [31], DINO [53])
102 to obtain a visual descriptor $\mathbf{f}_{t,i} = \text{Embed}(I_t^{\text{rgb}}, \mathbf{m}_{t,i})$. Each masked region is projected to 3D, de-
103 noised using DBSCAN clustering, and transformed to the map frame. This results in a pointcloud
104 $\mathbf{p}_{t,i}$ and its corresponding unit-normalized semantic feature vector $\mathbf{f}_{t,i}$.

105 **Object Association:** For every newly detected object $\langle \mathbf{p}_{t,i}, \mathbf{f}_{t,i} \rangle$, we compute semantic and geo-
106 metric similarity with respect to all objects $\mathbf{o}_{t-1,j} = \langle \mathbf{p}_{\mathbf{o}_j}, \mathbf{f}_{\mathbf{o}_j} \rangle$ in the map that shares any partial
107 geometric overlap. The geometric similarity $\phi_{\text{geo}}(i, j) = \text{nrratio}(\mathbf{p}_{t,i}, \mathbf{p}_{\mathbf{o}_j})$ is the proportion of
108 points in point cloud $\mathbf{p}_{t,i}$ that have nearest neighbors in point cloud $\mathbf{p}_{\mathbf{o}_j}$, within a distance threshold
109 of δ_{nn} . The semantic similarity $\phi_{\text{sem}}(i, j) = \mathbf{f}_{t,i}^T \mathbf{f}_{\mathbf{o}_j} / 2 + 1/2$ is the normalized cosine distance be-
110 tween the corresponding visual descriptors.² The overall similarity measure $\phi(i, j)$ is a sum of both:
111 $\phi(i, j) = \phi_{\text{sem}}(i, j) + \phi_{\text{geo}}(i, j)$. We perform object association by a greedy assignment³ strategy
112 where each new detection is matched with an existing object with the highest similarity score. If no
113 match is found with a similarity higher than δ_{sim} , we initialize a new object.

114 **Object Fusion:** If a detection $\mathbf{o}_{t-1,j}$ is associated with a mapped object \mathbf{o}_j , we fuse the detection with
115 the map. This is achieved by updating the object semantic feature as $\mathbf{f}_{\mathbf{o}_j} = (n_{\mathbf{o}_j} \mathbf{f}_{\mathbf{o}_j} + \mathbf{f}_{t,i}) / (n_{\mathbf{o}_j} + 1)$,
116 where $n_{\mathbf{o}_j}$ is the number of detections that have been associated to \mathbf{o}_j so far; and updating the
117 pointcloud as $\mathbf{p}_{t,i} \cup \mathbf{p}_{\mathbf{o}_j}$, followed by down-sampling to remove redundant points.

118 **Node Captioning:** Once the entire image sequence has been processed, a vision-language model,
119 denoted LVL $\text{LM}(\cdot)$, is used to generate object captions. For each object, the associated image crops
120 from the *best*⁴ 10 views are passed to the language model with the prompt “*describe the central*
121 *object in the image*” to generate a set of initial rough captions $\hat{\mathbf{c}}_j = \{\hat{\mathbf{c}}_{j,1}, \hat{\mathbf{c}}_{j,2}, \dots, \hat{\mathbf{c}}_{j,10}\}$ for each
122 detected object \mathbf{o}_j . Each set of captions is then refined to the final caption by passing $\hat{\mathbf{c}}_j$ to another
123 language model LLM (\cdot) with a prompt instruction to summarize the initial captions into a coherent
124 and accurate final caption \mathbf{c}_j .

125 2.2 Scene Graph Generation

126 Given the set of 3D objects \mathbf{O}_T obtained from the previous step, we estimate their spatial relation-
127 ships, i.e., the edges \mathbf{E}_T , to complete the 3D scene graph. We do this by first estimating potential
128 connectivity among object nodes based on their spatial overlaps. We compute the 3D bounding box
129 IoU between every pair of object nodes to obtain a similarity matrix (i.e., a dense graph), which we
130 prune by estimating a minimum spanning tree (MST), resulting in a refined set of potential edges
131 among the objects. To further determine the semantic relationships, for each edge in the MST,
132 we input the information about the object pair, consisting of object captions and 3D location, to a

¹Without loss of generality, $\text{Seg}(\cdot)$ may be replaced by open-/closed-vocabulary models to build category-specific mapping systems.

²For the sake of brevity, we only describe the best-performing geometric and semantic similarity measures. For an exhaustive list of alternatives, please see our project website and code.

³While we also experimented with optimal assignment strategies such as the Hungarian algorithm, we experimentally determined them to be slower and offer only a minuscule improvement over greedy association.

⁴We maintain a running index of the number of noise-free points each view contributes to the object point cloud.

133 language model LLM. The prompt instructs the model to describe the likely spatial relationship
 134 between the objects, such as “*a on b*” or “*b in a*”, along with the underlying reasoning. The model
 135 outputs a relationship label with an explanation detailing the rationale. The use of an LLM allows
 136 us to extend the nominal edge type defined above to other output relationships a language model can
 137 interpret, such as “a backpack *may be stored in* a closet” and “sheets of paper *may be recycled in* a
 138 trash can”. This results in an open-vocabulary 3D scene graph $\mathcal{M}_T = (\mathbf{O}_T, \mathbf{E}_T)$, a compact and
 139 efficient representation for use in downstream tasks.

140 2.3 Robotic Task Planning through LLMs

141 To enable users to carry out tasks described in natural language queries, we interface the scene
 142 graph \mathcal{M}_T with an LLM. For each object in \mathbf{O}_T , we construct JSON-structured text descriptions
 143 that include information about its 3D location (bounding box) and its node caption. Given a text
 144 query, we task the LLM to identify the most relevant object in the scene. We then pass the 3D
 145 pose of this object to the appropriate pipeline for the downstream task (e.g., grasping, navigation).
 146 This integration of *ConceptGraphs* with an LLM is easy to implement, and enables a wide range
 147 of open-vocabulary tasks by giving robots access to the semantic properties of surrounding objects⁵
 148 (see Sec. 3).

149 2.4 Implementation Details

150 The modularity of *ConceptGraphs* enables any appropriate open/closed-vocabulary segmentation
 151 model, LLM, or LVLm to be employed. Our experiments use Segment-Anything (SAM) [33] as
 152 the segmentation model $\text{Seg}(\cdot)$, and the CLIP image encoder [31] as the feature extractor $\text{Embed}(\cdot)$.
 153 We use LLaVA [55] as the vision-language model LVLm and GPT-4 [32] (gpt-4-0613) for our
 154 LLM. The voxel size for point cloud downsampling and nearest neighbor threshold δ_{nn} are both
 155 2.5cm. We use 1.1 for the association threshold δ_{sim} .

156 We also develop a variant of our system, *ConceptGraphs-Detector* (CG-D), where we employ an im-
 157 age tagging model (RAM [54]) to list the object classes present in the image and an open-vocabulary
 158 2D detector (Grounding DINO [34]) to obtain object bounding boxes⁶. In this variant, we need to
 159 separately handle detected background objects (wall, ceiling, floor) by merging them regardless of
 160 their similarity scores.

161 3 Experiments

162 3.1 Scene Graph Construction

163 We first evaluate the accuracy of the 3D scene graphs
 164 output by the *ConceptGraphs* system in Table 1.
 165 For each scene in the Replica dataset [56], we re-
 166 port scene graph accuracy metrics for both CG and
 167 the detector-variant CG-D. The open-vocabulary na-
 168 ture of our system makes automated evaluation of
 169 the quality of nodes and edges in the scene graph
 170 challenging. We instead evaluate the scene graph by
 171 engaging human evaluators on Amazon Mechanical
 172 Turk (AMT). For each node, we compute precision
 173 as the fraction of nodes for which at least 2 of 3 hu-
 174 man evaluators deem the node caption correct. We
 175 also report the number of valid objects retrieved by

	scene	node prec.	valid objects	duplicates	edge prec.
CG	room0	0.78	54	3	0.91
	room1	0.77	43	4	0.93
	room2	0.66	47	4	1.0
	office0	0.65	44	2	0.88
	office1	0.65	23	0	0.9
	office2	0.75	44	3	0.82
	office3	0.68	60	5	0.79
	Average	0.71	-	-	0.88
CG-D	room0	0.56	60	4	0.87
	room1	0.70	40	3	0.93
	room2	0.54	49	2	0.93
	office0	0.59	35	0	1.0
	office1	0.49	24	2	0.9
	office2	0.67	47	3	0.88
	office3	0.71	59	1	0.83
	Average	0.61	-	-	0.91

Table 1: **Accuracy of constructed scene graphs:** node precision indicates the accuracy of the label for each node (as measured by a human evaluator); valid objects is the number of human-recognizable objects (mturkers used) discovered by our system; duplicates are the number of redundant detections; edge precision indicates the accuracy of each estimated spatial relationship (again, as evaluated by an mturker)

⁵For large scenes where the description length of the scene graph exceeds the context length of the LLM, one can easily substitute in alternative (concurrent) LLM planners [52].

⁶We discard the (often noisy) *tags* produced by the image tagging model, relying instead on our node captions.

176 each variant by asking evaluators whether they deem each node a valid object. Both CG and CG-D
 177 identify a number of valid objects in each scene, and incur only a small number (0-5) of duplicate
 178 detections. The node labels are accurate about 70% of the time; most of the errors are incurred due
 179 to errors made by the LVM employed (LLaVA [55]). The edges (spatial relationships) are labeled
 180 with a high degree of accuracy (90% on average).

181 3.2 3D Semantic Segmentation

182 *ConceptGraphs* focuses on the construction of the open-
 183 vocabulary 3D scene graphs for scene understanding and
 184 planning. For completeness, in this section, we also use
 185 an open-vocabulary 3D semantic segmentation task to
 186 evaluate the quality of the obtained 3D maps. To generate
 187 the semantic segmentation, given a set of class names,
 188 we compute the similarity between the fused semantic
 189 feature of each object node and the CLIP text embed-
 190 dings of the phrase *an image of {class}*. Then
 191 the points associated with each object are assigned to the
 192 class with the highest similarity, which gives a point cloud
 193 with dense class labels. In Table 2, we report the semantic
 194 segmentation results on the Replica [56] dataset, follow-
 195 ing the evaluation protocol used in ConceptFusion [17].
 196 ConceptFusion+SAM, by replacing the Mask2Former used in ConceptFusion with the more perform-
 197 ant SAM [33] model. As shown in Table 2, the proposed *ConceptGraphs* performs comparably
 198 with or better than ConceptFusion, which has a much larger memory footprint.

	Method	mAcc	F-mIoU
Privileged	CLIPSeg (rd64-uni) [57]	28.21	39.84
	LSeg [58]	33.39	51.54
	OpenSeg [59]	41.19	53.74
Zero-shot	MaskCLIP [60]	4.53	0.94
	Mask2former + Global CLIP feat	10.42	13.11
	ConceptFusion [17]	24.16	31.31
	ConceptFusion [17] + SAM [33]	31.53	38.70
	<i>ConceptGraphs</i> (Ours)	40.63	35.95
	<i>ConceptGraphs-Detector</i> (Ours)	38.72	35.82

Table 2: Open-vocabulary semantic segmentation on the Replica [56] dataset. **Privileged** methods specifically finetune the pre-trained models for semantic segmentation. **Zero-shot** approaches do not need any finetuning and are evaluated off the shelf.

We also provide an additional baseline, ConceptFusion+SAM, by replacing the Mask2Former used in ConceptFusion with the more performant SAM [33] model. As shown in Table 2, the proposed *ConceptGraphs* performs comparably with or better than ConceptFusion, which has a much larger memory footprint.

199 3.3 Object Retrieval based on Text Queries

200 We assess the capability of *ConceptGraphs* to handle
 201 complex semantic queries, focusing on three key types.

- 202 • Descriptive: E.g., *A potted plant*.
- 203 • Affordance: E.g., *Something to use for temporarily securing a broken zipper*.
- 204
- 205 • Negation: E.g., *Something to drink but not soda*.

206 We evaluate on the Replica dataset [56] and a real-world
 207 scan of a lab, where we staged a number of items includ-
 208 ing clothes, tools, and toys. For Replica, human evalua-
 209 tors on AMT annotate captions for SAM mask proposals,
 210 which serve as both ground truth labels and descriptive queries. We created 5 affordance and nega-
 211 tion queries for each scene type (office & room) in Replica and 10 queries of each type for the
 212 lab scan, ensuring that each query corresponds to at least one relevant object. We manually select
 213 relevant objects as ground truth for each query.

214 We use two object retrieval strategies: CLIP-based and LLM-based. CLIP selects the object with
 215 the highest similarity to the query’s embedding, while the LLM goes through the scenegraph nodes
 216 to identify the object with the most relevant caption. Table 3 shows that CLIP excels with descrip-
 217 tive queries but struggles with complex affordance and negation queries [61]. For example, CLIP
 218 inaccurately retrieves a backpack for the broken zipper query, whereas the LLM correctly identifies
 219 a roll of tape. The LLM performs well across the board, but is limited by the accuracy of the node
 220 captions, as discussed in Section 3.1. Since the lab has a larger variety of objects to choose from,
 221 the LLM finds compatible objects for complex queries more reliably there.

Dataset	Query Type	Model	R@1	R@2	R@3	# Queries
Replica	Descriptive	CLIP	0.59	0.82	0.86	20
		LLM	0.61	0.64	0.64	
	Affordance	CLIP	0.43	0.57	0.63	5
		LLM	0.57	0.63	0.66	
	Negation	CLIP	0.26	0.60	0.71	5
		LLM	0.80	0.89	0.97	
Lab	Descriptive	CLIP	1.00	–	–	10
		LLM	1.00	–	–	
	Affordance	CLIP	0.40	0.60	0.60	10
		LLM	1.00	–	–	
	Negation	CLIP	0.00	–	–	10
		LLM	1.00	–	–	

Table 3: Object retrieval from text queries on the Replica and REAL Lab scenes. We measure the top-1, top-2, and top-3 recall. CLIP refers to object retrieval using cosine similarity, whereas LLM refers to having an LLM parse the scene graph and return the most relevant object.

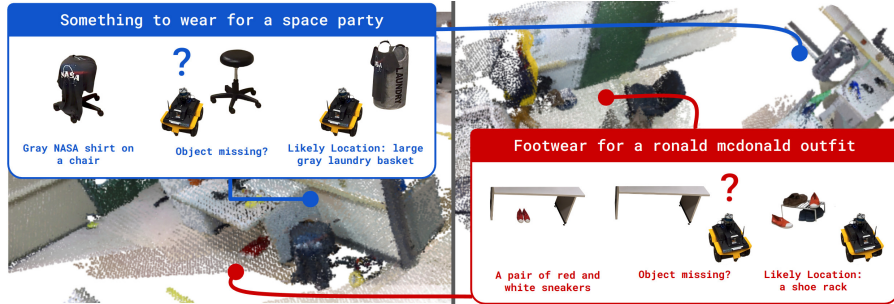


Figure 4: A Jackal robot answering user queries using the *ConceptGraphs* representation of a lab environment. We first query an LLM to identify the most relevant object given the user query, then validate with an LVLM if the target object is at the expected location. If not, we query the LLM again to find a likely location or container for the missing object. (Blue) When prompted with something to wear for a space party, the Jackal attempts to find a grey shirt with a NASA logo. After failing to detect the shirt at the expected location, the LLM reasons that it could likely be in the laundry bag. (Red) The Jackal searches for red and white sneakers after receiving the user query footwear for a Ronald McDonald outfit. The LLM redirects the robot to a shoe rack after failing to detect the sneakers at their initial mapped position.

222 3.4 Complex Visual-Language Queries

223 To assess the performance of *ConceptGraphs* in a real-world environment, we carry out navigation
 224 experiments in a lab environment with a Clearpath Jackal UGV. The robot is equipped with a VLP-
 225 16 LiDAR and a forward-facing Realsense D435i camera.

226 The Jackal needs to respond to abstract user queries and navigate to the most relevant object (Fig-
 227 ure 1). By using an LVLM [55] to add a description of the current camera image to the text prompt,
 228 the robot can also answer visual queries. For example, when shown a picture of Michael Jordan and
 229 prompted with `Something this guy would play with`, the robot finds a basketball.

230 3.5 Object Search and Traversability Estimation

231 In this section, we showcase how the interaction between the *ConceptGraphs* representation and an
 232 LLM can enable a mobile robot to access a vast knowledge base of everyday objects. Specifically,
 233 we prompt an LLM to infer two additional object properties from *ConceptGraphs* captions: i) the
 234 location where a given object is typically found, and ii) if the object can be safely pushed or traversed
 235 by the Jackal robot. We design two tasks around the LLM predictions.

236 **Object search:** The robot receives an abstract user query
 237 and must navigate to the most relevant object in the *Con-*
 238 *ceptGraphs* map. Using an LVLM [55], the robot then
 239 checks if the object is at the expected location. If not, it
 240 queries an LLM to find a new plausible location given the
 241 captions of the other objects in the representation. In our
 242 prompt, we nudge the LLM to consider typical containers
 243 or storage locations. We illustrate two such queries where
 244 the target object is moved in Figure 4.

245 **Traversability estimation:** As shown in Fig. 3, we de-
 246 sign a real-world scenario where the robot finds itself en-
 247 caged by objects. In this scenario, the robot must push
 248 around multiple objects and create a path to the goal
 249 state. While traversability can be learned through experi-
 250 ence [62], we show that grounding LLM knowledge in
 251 a 3D map can grant similar capabilities to robotic agents.



Figure 3: The Jackal robot solving a traversability challenge. All paths to the goal are obstructed by objects. We query an LLM to identify which objects can be safely pushed or traversed by the robot (green) and which objects would be too heavy or hinder the robot’s movement (red). The LLM relies on the *ConceptGraphs* node captions to make traversability predictions and we add the non-traversable objects to the Jackal costmap for path planning. The Jackal successfully reaches the goal by going through a curtain and pushing a basketball, while also avoiding contact with bricks, an iron dumbbell, and a flower pot.

252 3.6 Open-Vocabulary Pick and Place

253 To illustrate how *ConceptGraphs* can act as the perception backbone for open-vocabulary mobile
254 manipulation, we conducted a series of experiments with a Boston Dynamics Spot Arm robot. Us-
255 ing an onboard RGBD camera and a *ConceptGraphs* representation of the scene, the Spot robot
256 responds to the query `cuddly quacker` by grabbing a duck plush toy and placing it in a nearby
257 box (Figure 1). In the supplementary video, Spot completes a similar grasping maneuver with a
258 mango when prompted with the query `something healthy to eat`.

259 3.7 Localization and Map Updates

260 *ConceptGraphs* can also be used for object-based localization and map updates. We showcase this
261 with a 3-DoF (x , y and yaw) localization and remapping task in the AI2Thor [63, 64] simulation
262 environment, where a mobile robot uses a particle filter to localize in a pre-built *ConceptGraphs* map
263 of the environment. During the observation update step of particle filtering, the robot’s detections
264 are matched against the objects in the map based on the hypothesized pose, in a similar way as
265 described in Section 2.1. Refer to our video explainer for a demonstration.

266 3.8 Limitations

267 Despite its impressive performance, *ConceptGraphs* has failure modes that remain to be addressed
268 in future work. First, node captioning incurs errors due to the current limitations of LVLMS like
269 LLaVA [55]. Second, our 3D scene graph occasionally misses small or thin objects and makes
270 duplicate detections. This impacts downstream planning, particularly when the incorrect detection
271 is crucial to planning success. Additionally, the computational and economic costs of our system
272 include multiple LVLMS (LLaVA [55]) and one or more proprietary LLM inference(s) when building
273 and querying the scenegraph, which may be significant.

274 4 Concurrent Work

275 We briefly review recent and unpublished pre-prints that are exploring themes related to open-
276 vocabulary object-based factorization of 3D scenes. Concurrently to us, [65, 66] have explored
277 open-vocabulary object-based factorization of 3D scenes. Where [65] assumes a pre-built point
278 cloud map of the scene, [66] builds a map on the go. Both approaches associate CLIP descriptors to
279 the reconstruction, resulting in performance comparable to our system’s CLIP variant, which strug-
280 gles with queries involving complex affordances and negation, as shown in Table 3. OGSV [67]
281 is closer to our setting, building an open-vocabulary 3D scene graph from RGB-D images. How-
282 ever, [67] relies on a (closed-set) graph neural network to predict object relationships; whereas
283 *ConceptGraphs* relies on LLMs, eliminating the need to train an object relation model.

284 5 Conclusion

285 In this paper, we introduced *ConceptGraphs*, a novel approach to open-vocab object-centric 3D
286 scene representation that addresses key limitations in the existing landscape of dense and implicit
287 representations. Through effective integration of foundational 2D models, *ConceptGraphs* signifi-
288 cantly mitigates memory constraints, provides relational information among objects, and allows for
289 dynamic updates to the scene—three pervasive challenges in current methods. Experimental evi-
290 dence underscores *ConceptGraphs*’ robustness and extensibility, highlighting its superiority over
291 existing baselines for a variety of real-world tasks including manipulation and navigation. The ver-
292 satility of our framework also accommodates a broad range of downstream applications, thereby
293 opening new avenues for innovation in robot perception and planning. Future work may delve into
294 integrating temporal dynamics into the model and assessing its performance in less structured, more
295 challenging environments.

References

- 296
- 297 [1] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohi, J. Shotton,
298 S. Hodges, and A. Fitzgibbon. Kinectfusion: Real-time dense surface mapping and track-
299 ing. In *IEEE international symposium on mixed and augmented reality*, pages 127–136. IEEE,
300 2011. 2
- 301 [2] T. Whelan, S. Leutenegger, R. Salas-Moreno, B. Glocker, and A. Davison. Elasticfusion:
302 Dense slam without a pose graph. In *Robotics Science and Systems*, 2015. 2
- 303 [3] B. Wen, J. Tremblay, V. Blukis, S. Tyree, T. Müller, A. Evans, D. Fox, J. Kautz, and S. Birch-
304 field. BundleSDF: Neural 6-dof tracking and 3d reconstruction of unknown objects. In *Pro-
305 ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
306 606–617, 2023. 2
- 307 [4] E. Sucar, S. Liu, J. Ortiz, and A. J. Davison. iMAP: Implicit mapping and positioning in real-
308 time. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages
309 6229–6238, 2021. 2
- 310 [5] Z. Zhu, S. Peng, V. Larsson, W. Xu, H. Bao, Z. Cui, M. R. Oswald, and M. Pollefeys. Nice-
311 slam: Neural implicit scalable encoding for slam. In *Proceedings of the IEEE/CVF Conference
312 on Computer Vision and Pattern Recognition*, pages 12786–12796, 2022. 2
- 313 [6] J. L. Schönberger and J.-M. Frahm. Structure-from-motion revisited. In *Proceedings of the
314 IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016. 2
- 315 [7] J. L. Schönberger, E. Zheng, J.-M. Frahm, and M. Pollefeys. Pixelwise view selection for
316 unstructured multi-view stereo. In *European Conference on Computer Vision*, pages 501–518.
317 Springer, 2016. 2
- 318 [8] J. McCormac, A. Handa, A. Davison, and S. Leutenegger. Semanticfusion: Dense 3d semantic
319 mapping with convolutional neural networks. In *IEEE International Conference on Robotics
320 and automation (ICRA)*, pages 4628–4635. IEEE, 2017. 2
- 321 [9] M. Runz, M. Buffier, and L. Agapito. Maskfusion: Real-time recognition, tracking and re-
322 construction of multiple moving objects. In *IEEE International Symposium on Mixed and
323 Augmented Reality (ISMAR)*, pages 10–20. IEEE, 2018. 2
- 324 [10] J. McCormac, R. Clark, M. Bloesch, A. Davison, and S. Leutenegger. Fusion++: Volumetric
325 object-level slam. In *international conference on 3D vision (3DV)*, pages 32–41. IEEE, 2018.
326 2
- 327 [11] G. Narita, T. Seno, T. Ishikawa, and Y. Kaji. Panopticfusion: Online volumetric semantic map-
328 ping at the level of stuff and things. In *2019 IEEE/RSJ International Conference on Intelligent
329 Robots and Systems (IROS)*, pages 4205–4212. IEEE, 2019. 2
- 330 [12] J. Qian, V. Chatrath, J. Yang, J. Servos, A. P. Schoellig, and S. L. Waslander. POCD: probabilis-
331 tic object-level change detection and volumetric mapping in semi-static scenes. In K. Hauser,
332 D. A. Shell, and S. Huang, editors, *Robotics Science and Systems*, 2022. 2
- 333 [13] J. Qian, V. Chatrath, J. Servos, A. Mavrinas, W. Burgard, S. L. Waslander, and A. P. Schoellig.
334 POV-SLAM: probabilistic object-aware variational SLAM in semi-static environments. In
335 K. E. Bekris, K. Hauser, S. L. Herbert, and J. Yu, editors, *Robotics Science and Systems*, 2023.
336 2
- 337 [14] K. Li, D. DeTone, Y. F. S. Chen, M. Vo, I. Reid, H. Rezatofghi, C. Sweeney, J. Straub, and
338 R. Newcombe. O DAM: Object detection, association, and mapping using posed rgb video. In
339 *Proceedings of International Conference on Computer Vision*, 2021. 2

- 340 [15] M. Zins, G. Simon, and M.-O. Berger. OA-SLAM: Leveraging objects for camera relocal-
341 ization in visual slam. In *IEEE International Symposium on Mixed and Augmented Reality*
342 (*ISMAR*), pages 720–728. IEEE, 2022. 2
- 343 [16] K. Liu, F. Zhan, J. Zhang, M. Xu, Y. Yu, A. E. Saddik, C. Theobalt, E. Xing, and S. Lu.
344 3d open-vocabulary segmentation with foundation models. *arXiv preprint arXiv:2305.14093*,
345 2023. 2
- 346 [17] K. M. Jatavallabhula, A. Kuwajerwala, Q. Gu, M. Omama, G. Iyer, S. Saryazdi, T. Chen,
347 A. Maalouf, S. Li, N. V. Keetha, A. Tewari, J. B. Tenenbaum, C. M. de Melo, K. M. Krishna,
348 L. Paull, F. Shkurti, and A. Torralba. ConceptFusion: Open-set multimodal 3d mapping. In
349 *Robotics: Science and Systems*, 2023. doi:10.15607/RSS.2023.XIX.066. 2, 6
- 350 [18] S. Peng, K. Genova, C. Jiang, A. Tagliasacchi, M. Pollefeys, T. Funkhouser, et al. Openscene:
351 3d scene understanding with open vocabularies. In *Proceedings of the IEEE/CVF Conference*
352 *on Computer Vision and Pattern Recognition*, pages 815–824, 2023. 2
- 353 [19] R. Ding, J. Yang, C. Xue, W. Zhang, S. Bai, and X. Qi. Lowis3d: Language-driven open-world
354 instance-level 3d scene understanding. *arXiv preprint arXiv:2308.00353*, 2023. 2
- 355 [20] R. Ding, J. Yang, C. Xue, W. Zhang, S. Bai, and X. Qi. PLA: Language-driven open-vocabulary
356 3d scene understanding. In *Proceedings of Computer Vision and Pattern Recognition*, 2023. 2
- 357 [21] J. Zhang, R. Dong, and K. Ma. CLIP-FO3D: Learning free open-world 3d scene representa-
358 tions from 2d dense clip. *arXiv preprint arXiv:2303.04748*, 2023. 2
- 359 [22] V. Tschernezki, I. Laina, D. Larlus, and A. Vedaldi. Neural feature fusion fields: 3d distillation
360 of self-supervised 2d image representations. In *International Conference on 3D Vision (3DV)*.
361 IEEE, 2022. 2
- 362 [23] S. Kobayashi, E. Matsumoto, and V. Sitzmann. Decomposing nerf for editing via feature field
363 distillation. *Neural Information Processing Systems*, 35:23311–23330, 2022. 2
- 364 [24] N. M. M. Shafiullah, C. Paxton, L. Pinto, S. Chintala, and A. Szlam. Clip-fields: Weakly
365 supervised semantic fields for robotic memory. In K. E. Bekris, K. Hauser, S. L. Herbert, and
366 J. Yu, editors, *Robotics: Science and Systems*, 2023. doi:10.15607/RSS.2023.XIX.074. 2
- 367 [25] N. Tsagkas, O. Mac Aodha, and C. X. Lu. VI-fields: Towards language-grounded neural
368 implicit spatial representations. *arXiv preprint arXiv:2305.12427*, 2023. 2
- 369 [26] J. Kerr, C. M. Kim, K. Goldberg, A. Kanazawa, and M. Tancik. LERF: Language embedded
370 radiance fields. In *International Conference on Computer Vision (ICCV)*, 2023. 2
- 371 [27] C. Huang, O. Mees, A. Zeng, and W. Burgard. Audio visual language maps for robot naviga-
372 tion. *arXiv preprint arXiv:2303.07522*, 2023. 2
- 373 [28] W. Shen, G. Yang, A. Yu, J. Wong, L. P. Kaelbling, and P. Isola. Distilled feature fields
374 enable few-shot manipulation. In *International Conference on Robot Learning*, 2023. URL
375 https://openreview.net/forum?id=Rb0nGIt_kh5. 2
- 376 [29] F. Engelmann, F. Manhardt, M. Niemeyer, K. Tateno, M. Pollefeys, and F. Tombari. Open-set
377 3d scene segmentation with rendered novel views, 2023. 2
- 378 [30] K. Mazur, E. Sucar, and A. J. Davison. Feature-realistic neural fusion for real-time, open set
379 scene understanding. In *IEEE International Conference on Robotics and Automation (ICRA)*.
380 IEEE, 2023. 2
- 381 [31] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,
382 P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervi-
383 sion. In *International Conference on Machine Learning*. PMLR, 2021. 2, 3, 4, 5

- 384 [32] OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 2, 3, 5
- 385 [33] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead,
386 A. C. Berg, W.-Y. Lo, et al. Segment anything. *Proceedings of International Conference on*
387 *Computer Vision*, 2023. 2, 3, 5, 6
- 388 [34] S. Liu, Z. Zeng, T. Ren, F. Li, H. Zhang, J. Yang, C. Li, J. Yang, H. Su, J. Zhu, et al. Grounding
389 dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint*
390 *arXiv:2303.05499*, 2023. 2, 3, 5
- 391 [35] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image syn-
392 thesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer*
393 *vision and pattern recognition*, pages 10684–10695, 2022. 2
- 394 [36] Y. Hong, Y. Du, C. Lin, J. Tenenbaum, and C. Gan. 3d concept grounding on neural fields.
395 *Neural Information Processing Systems*, 2022. 2
- 396 [37] Y. Hong, C. Lin, Y. Du, Z. Chen, J. B. Tenenbaum, and C. Gan. 3d concept learning and rea-
397 soning from multi-view images. In *Proceedings of Computer Vision and Pattern Recognition*,
398 2023. 2
- 399 [38] Y. Hong, H. Zhen, P. Chen, S. Zheng, Y. Du, Z. Chen, and C. Gan. 3d-llm: Injecting the 3d
400 world into large language models. *Neural Information Processing Systems*, 2023. 2
- 401 [39] M. Shridhar, L. Manuelli, and D. Fox. CLIPort: What and where pathways for robotic manip-
402 ulation. In *Conference on Robot Learning*, volume 164, pages 894–906. PMLR, 2021. 2
- 403 [40] S. Sharma, A. Rashid, C. M. Kim, J. Kerr, L. Y. Chen, A. Kanazawa, and K. Goldberg.
404 Language embedded radiance fields for zero-shot task-oriented grasping. In *International*
405 *Conference on Robot Learning*, 2023. URL <https://openreview.net/forum?id=k-Fg8JDQmc>. 2
406
- 407 [41] S. Y. Gadre, M. Wortsman, G. Ilharco, L. Schmidt, and S. Song. Clip on wheels: Zero-shot
408 object navigation as object localization and exploration. *arXiv preprint arXiv:2203.10421*,
409 2022. 2
- 410 [42] D. Shah, B. Osiński, S. Levine, et al. Lm-nav: Robotic navigation with large pre-trained mod-
411 els of language, vision, and action. In *International Conference on Robot Learning*. PMLR,
412 2023. 2
- 413 [43] M. Fisher, M. Savva, and P. Hanrahan. Characterizing structural relationships in scenes using
414 graph kernels. *ACM Trans. Graph.*, 30(4):34, 2011. 2, 3
- 415 [44] P. Gay, J. Stuart, and A. Del Bue. Visual graphs from motion (vgfm): Scene understanding
416 with object geometry reasoning. In *Asian Conference on Computer Vision*. Springer, 2019. 2,
417 3
- 418 [45] I. Armeni, Z.-Y. He, J. Gwak, A. R. Zamir, M. Fischer, J. Malik, and S. Savarese. 3d scene
419 graph: A structure for unified semantics, 3d space, and camera. In *Proceedings of International*
420 *Conference on Computer Vision*, October 2019. 2, 3
- 421 [46] U.-H. Kim, J.-M. Park, T.-J. Song, and J.-H. Kim. 3-d scene graph: A sparse and semantic rep-
422 resentation of physical environments for intelligent agents. *IEEE transactions on cybernetics*,
423 50(12):4921–4933, 2019. 2, 3
- 424 [47] J. Wald, H. Dhano, N. Navab, and F. Tombari. Learning 3d semantic scene graphs from 3d
425 indoor reconstructions. In *Proceedings of Computer Vision and Pattern Recognition*, 2020. 2,
426 3

- 427 [48] A. Rosinol, A. Violette, M. Abate, N. Hughes, Y. Chang, J. Shi, A. Gupta, and L. Carlone.
428 Kimera: From slam to spatial perception with 3d dynamic scene graphs. *The International*
429 *Journal of Robotics Research*, 40(12-14):1510–1546, 2021. 2
- 430 [49] N. Hughes, Y. Chang, and L. Carlone. Hydra: A real-time spatial perception system for 3d
431 scene graph construction and optimization. *arXiv preprint arXiv:2201.13360*, 2022. 2
- 432 [50] S.-C. Wu, J. Wald, K. Tateno, N. Navab, and F. Tombari. Scenegrphfusion: Incremental 3d
433 scene graph prediction from rgb-d sequences. In *Proceedings of Computer Vision and Pattern*
434 *Recognition*, 2021. 2
- 435 [51] C. Agia, K. M. Jatavallabhula, M. Khodeir, O. Miksik, V. Vineet, M. Mukadam, L. Paull,
436 and F. Shkurti. Taskography: Evaluating robot task planning over large 3d scene graphs. In
437 *International Conference on Robot Learning*. PMLR, 2022. 2
- 438 [52] K. Rana, J. Abou-Chakra, S. Garg, J. Haviland, I. Reid, and N. Suenderhauf. Sayplan: Ground-
439 ing large language models using 3d scene graphs for scalable task planning. In *International*
440 *Conference on Robot Learning*, 2023. URL <https://openreview.net/forum?id=wMpOM00Ss7a>. 2, 5, 19
- 442 [53] M. Oquab, T. Darcet, T. Moutakanni, H. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. Haz-
443 iza, F. Massa, A. El-Nouby, et al. Dinov2: Learning robust visual features without supervision.
444 *arXiv preprint arXiv:2304.07193*, 2023. 3, 4
- 445 [54] Y. Zhang, X. Huang, J. Ma, Z. Li, Z. Luo, Y. Xie, Y. Qin, T. Luo, Y. Li, S. Liu, et al. Recognize
446 anything: A strong image tagging model. *arXiv preprint arXiv:2306.03514*, 2023. 3, 5
- 447 [55] H. Liu, C. Li, Q. Wu, and Y. J. Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*,
448 2023. 3, 5, 6, 7, 8, 13, 18, 19
- 449 [56] J. Straub, T. Whelan, L. Ma, Y. Chen, E. Wijmans, S. Green, J. J. Engel, R. Mur-Artal, C. Ren,
450 S. Verma, A. Clarkson, M. Yan, B. Budge, Y. Yan, X. Pan, J. Yon, Y. Zou, K. Leon, N. Carter,
451 J. Briales, T. Gillingham, E. Mueggler, L. Pesqueira, M. Savva, D. Batra, H. M. Strasdat, R. D.
452 Nardi, M. Goesele, S. Lovegrove, and R. Newcombe. The Replica dataset: A digital replica of
453 indoor spaces. *arXiv preprint arXiv:1906.05797*, 2019. 5, 6, 14
- 454 [57] T. Lüddecke and A. Ecker. Image segmentation using text and image prompts. In *Proceedings*
455 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7086–7096,
456 2022. 6
- 457 [58] B. Li, K. Q. Weinberger, S. Belongie, V. Koltun, and R. Ranftl. Language-driven semantic
458 segmentation. In *International Conference on Learning Representations*, 2022. 6
- 459 [59] G. Ghiasi, X. Gu, Y. Cui, and T.-Y. Lin. Scaling open-vocabulary image segmentation with
460 image-level labels. In *European Conference on Computer Vision*, pages 540–557. Springer,
461 2022. 6
- 462 [60] C. Zhou, C. C. Loy, and B. Dai. Extract free dense labels from clip. In *European Conference*
463 *on Computer Vision (ECCV)*, 2022. 6
- 464 [61] Y. Du, S. Li, and I. Mordatch. Compositional visual generation with energy based models. In
465 *Neural Information Processing Systems*, 2020. 6
- 466 [62] S. Levine and D. Shah. Learning robotic navigation from experience: principles, methods and
467 recent results. *Philosophical Transactions of the Royal Society B*, 378(1869):20210447, 2023.
468 7
- 469 [63] E. Kolve, R. Mottaghi, W. Han, E. VanderBilt, L. Weihs, A. Herrasti, M. Deitke, K. Ehsani,
470 D. Gordon, Y. Zhu, et al. Ai2-thor: An interactive 3d environment for visual ai. *arXiv preprint*
471 *arXiv:1712.05474*, 2017. 8

- 472 [64] M. Deitke, E. Vanderbilt, A. Herrasti, L. Weihs, K. Ehsani, J. Salvador, W. Han, E. Kolve,
473 A. Kembhavi, and R. Mottaghi. Proctor: Large-scale embodied ai using procedural genera-
474 tion. *Advances in Neural Information Processing Systems*, 35:5982–5994, 2022. 8
- 475 [65] A. Takmaz, E. Fedele, R. W. Sumner, M. Pollefeys, F. Tombari, and F. Engelmann. Open-
476 mask3d: Open-vocabulary 3d instance segmentation. *arXiv preprint arXiv:2306.13631*, 2023.
477 8
- 478 [66] S. Lu, H. Chang, E. P. Jing, A. Boularias, and K. Bekris. OVIR-3d: Open-vocabulary 3d
479 instance retrieval without training on 3d data. In *International Conference on Robot Learning*,
480 2023. URL https://openreview.net/forum?id=gVBvtRqU1_. 8
- 481 [67] H. Chang, K. Boyalakuntla, S. Lu, S. Cai, E. P. Jing, S. Keskar, S. Geng, A. Abbas,
482 L. Zhou, K. Bekris, and A. Boularios. Context-aware entity grounding with open-vocabulary
483 3d scene graphs. In *International Conference on Robot Learning*, 2023. URL <https://openreview.net/forum?id=cjEI5qXoT0>. 8
484
- 485 [68] E. Jelavic, J. Nubert, and M. Hutter. Open3d slam: Point cloud based mapping and localiza-
486 tion for education. In *Robotic Perception and Mapping: Emerging Techniques, ICRA 2022*
487 *Workshop*, page 24. ETH Zurich, Robotic Systems Lab, 2022. 18
- 488 [69] M. Labbé and F. Michaud. Rtab-map as an open-source lidar and visual simultaneous local-
489 ization and mapping library for large-scale and long-term online operation. *Journal of Field*
490 *Robotics*, 36(2):416–446, 2019. 18
- 491 [70] Z. Sun, S. Shen, S. Cao, H. Liu, C. Li, Y. Shen, C. Gan, L.-Y. Gui, Y.-X. Wang, Y. Yang,
492 K. Keutzer, and T. Darrell. Aligning large multimodal models with factually augmented rlhf.
493 2023. 19

494 Appendix

495 A1 3D Scene Graph: Generating Node Captions

496 Once we build an object-level map of the scene using the methodology described in Sec. 2.1, we ex-
497 tract and summarize captions for each object. We first extract upto the 10 *most-informative* views for
498 each object, by tracking the number of (noise-free) 3D points that each image segment contributes
499 to an object in the map⁷. Intuitively, these views offer the best image views for the object. We run
500 each view through an LVLMM, here LLaVA-7B [55], to generate an image caption. We use the same
501 prompt across all images: *describe the central object in this image*.

502 We found the captions generated by LLaVA-7B to be incoherent or unreliable across all viewpoints.
503 To alleviate this, we employed GPT-4 as a caption summarizer, to map all of the LLaVA-7B cap-
504 tions to a coherent object tag (or optionally, declare the object as an invalid detection). We use the
505 following GPT-4 system prompt:

⁷We track these statistics throughout the mapping lifecycle; meaning that we do not impose any additional computational overhead to determine the 10 best views per object

Identify and describe objects in scenes. Input and output must be in JSON format. The input field 'captions' contains a list of image captions aiming to identify objects. Output 'summary' as a concise description of the identified object(s). An object mentioned multiple times is likely accurate. If various objects are repeated and a container/surface is noted such as a shelf or table, assume the (repeated) objects are on that container/surface. For unrelated, non-repeating (or empty) captions, summarize as 'conflicting (or empty) captions about [objects]' and set 'object_tag' to 'invalid'. Output 'possible_tags' listing potential object categories. Set 'object_tag' as the conclusive identification. Focus on indoor object types, as the input captions are from indoor scans.

Listing 1: GPT-4 system prompt used for caption summarization

A2 LLM Planner: Implementation details

For task planning over 3D scene graphs, we use GPT-4 (gpt-4-0613) with a context length of 8K tokens⁸. We first convert each node in the 3D scene graph into a structured text format (here, a JSON string). Each entry in the JSON list corresponds to one object in the scene, and contains the following attributes:

1. *object id*: a unique (numerical) object identifier
2. *bounding box extents*: dimensions of each side of the bounding cuboid
3. *bounding box center*: centroid of the object bounding cuboid
4. *object tag*: a brief tag describing the object
5. *caption*: a one-sentence caption (possibly encoding more details than present in the object tag)

Here is a sample snippet from the scene graph for the room0 scene of the Replica [56] dataset.

```
[
  {
    id: 2,
    bbox_extent: [2.0, 0.7, 0.6],
    bbox_center: [-0.6, 1.1, -1.2],
    object_tag: wooden dresser or chest of drawers,
    caption: A wooden dresser or chest of drawers
  },
  {
    id: 3,
    bbox_extent: [0.6, 0.5, 0.4],
    bbox_center: [2.8, -0.4, -0.8],
    object_tag: vase,
    caption: a white, floral-patterned vase (or possibly a ceramic bowl)
  },
  ...
  ...
  {
    id: 110,
    bbox_extent: [1.2, 0.6, 0.0],
    bbox_center: [2.2, 2.1, 1.2],
    object_tag: light fixture,
    caption: a light fixture hanging from the ceiling
  }
]
```

Listing 2: Sample text entries in the 3D scene graph

⁸We also prototyped variants of this approach on off-the-shelf LLMs with larger context lengths, such as Claude-2 with a context length of 32K tokens, and found it to work reliably.

544 For planning using the above text description of the scene, we instantiate the system prompt for
545 GPT-4 as follows:

```
546 The input to the model is a 3D scene described in a JSON format. Each
547 entry in the JSON describes one object in the scene, with the
548 following five fields:
549 1. "id": a unique object id
550 2. "bbox_extent": extents of the 3D bounding box for the object
551 3. "bbox_center": centroid of the 3D bounding box for the object
552 4. "object_tag": a brief (but sometimes inaccurate) tag categorizing the
553 object
554 5. "caption": a brief caption for the object
555
556 Once you have parsed the JSON and are ready to answer questions about the
557 scene, say "I'm ready".
558
559 The user will then begin to ask questions, and the task is to answer
560 various user queries about the 3D scene. For each user question,
561 respond with a JSON dictionary with the following fields:
562 1. "inferred_query": your interpretation of the user query in a succinct
563 form
564 2. "relevant_objects": list of relevant object ids for the user query (if
565 applicable)
566 3. "query_achievable": whether or not the user-specified query is
567 achievable using the objects and descriptions provided in the 3D
568 scene.
569 4. "final_relevant_objects": A final list of objects relevant to the user
570 -specified task. As much as possible, sort all objects in this list
571 such that the most relevant object is listed first, followed by the
572 second most relevant, and so on.
573 5. "explanation": A brief explanation of what the most relevant object(s)
574 is(are), and how they achieve the user-specified task.
```

Listing 3: System prompt for the LLM Planner interface

575 Here are some example GPT-4 responses on our queries

```
576 User: I need a quick gift. Help!
577
578 LLM-Planner: {
579   inferred_query: Find suitable object for a gift.,
580   relevant_objects: [0, 6, 7, 23, 25, 31],
581   query_achievable: true,
582   final_relevant_objects: [6],
583   explanation: The most suitable object for a gift could be the ceramic
584     vase (id 6). Vases are commonly gifted items and this one could
585     potentially be filled with a plant or flower arrangement, making a
586     lovely present.
587 }
588
589 User: Need something to use as a paperweight
590
591 LLM-Planner: {
592   inferred_query: Find object to use as a paperweight.,
593   relevant_objects: [0, 6, 7, 8, 15, 23, 29, 31],
594   query_achievable: true,
595   final_relevant_objects: [6],
596   explanation: The ceramic vase (id 6) could potentially be used as a
597     paperweight due to its small size and assumed weight.
598 }
```

Listing 4: Sample queries and raw outputs from the LLM Planner

599 **A3 Example Text Queries**

600 Here we provide some of the text queries used in our experiments.

601 Replica Dataset Descriptive Queries: For each room, we randomly selected 20 ground truth annotations collected via Amazon Mechanical Turk (AMT). Here is a sample from `room0` and `office0`.

603 **office0 Descriptive Queries:**

- 604 1. *This is a trash can against the wall next to a sofa.*
- 605 2. *A chaise lounge right next to a small table.*
- 606 3. *This is a television.*
- 607 4. *This is a dropped, tiled ceiling in what appears to be a classroom for children.*
- 608 5. *This is a plant and it is next to the screens.*
- 609 6. *This is the back of a chair in front of a screen.*
- 610 7. *A small table in front of a large gray sectional couch.*
- 611 8. *This is an armless chair and it's opposite a coffee table by the sofa.*
- 612 9. *This is a plug-in and it is on the floor.*
- 613 10. *These are table legs and they are underneath the table.*
- 614 11. *These are chairs and they are next to a table.*
- 615 12. *A diner style table in front of two chairs.*
- 616 13. *These are rocks and they are on the wall.*
- 617 14. *This is the right panel of a lighted display screen.*
- 618 15. *This is a planet and it is on the wall.*
- 619 16. *This is an electronic display screen showing a map, on the wall.*
- 620 17. *This is a couch and it is between a table and the wall.*
- 621 18. *This is a garbage can and it is in front of the wall.*
- 622 19. *This is a rug and it is on the floor.*
- 623 20. *This is a table that is above the floor.*

624 **room0 Descriptive Queries:**

- 625 1. *This is a pillow and this is on top of a couch.*
- 626 2. *A pillow on top of a white couch.*
- 627 3. *This is a couch and it is under a window.*
- 628 4. *This is a stool and it is on top of a rug.*
- 629 5. *This is a side table under a lamp.*
- 630 6. *This is a ceiling light next to the window.*
- 631 7. *This is an end table and it is below a lamp.*
- 632 8. *These are books and they are on the table.*
- 633 9. *This is a couch and it is in front of the wall.*
- 634 10. *White horizontal blinds in a well lit room.*
- 635 11. *This is a striped throw pillow on the loveseat.*
- 636 12. *The pillow is on top of the chair.*
- 637 13. *This is a window and it is next to a door.*

- 638 14. *This is a hurricane candle and it is on top of a cabinet.*
639 15. *This is a vase and it is on top of the table.*
640 16. *This is a vent in the ceiling.*
641 17. *This is a fish and it is on top of a cabinet.*
642 18. *This is a window behind a chair.*
643 19. *This is a trash can against a wall.*
644 20. *Two cream colored cushioned chairs with blue pillows adjacent to each other.*

645 **Replica Dataset Affordance Queries for Office Scenes:**

- 646 1. *Something to watch the news on*
647 2. *Something to tell the time*
648 3. *Something comfortable to sit on*
649 4. *Something to dispose of wastepaper in*
650 5. *Something to add light into the room*

651 **Replica Dataset Affordance Queries for Room Scenes:**

- 652 1. *Somewhere to store decorative cups*
653 2. *Something to add light into the room*
654 3. *Somewhere to set food for dinner*
655 4. *Something I can open with my keys*
656 5. *Somewhere to sit upright for a work call*

657 **Replica Dataset Negation Queries for Office Scenes:**

- 658 1. *Something to sit on other than a chair*
659 2. *Something very heavy, unlike a clock*
660 3. *Something rigid, unlike a cushion*
661 4. *Something small, unlike a couch*
662 5. *Something light, unlike a table*

663 **Replica Dataset Negation Queries for Room Scenes:**

- 664 1. *Something small, unlike a cabinet*
665 2. *Something light, unlike a table*
666 3. *Something soft, unlike a table*
667 4. *Something not transparent, unlike a window*
668 5. *Something rigid, unlike a rug*

669 **REAL Lab Scan Descriptive queries:**

- 670 1. *A pair of red and white sneakers*
671 2. *A NASA t-shirt*
672 3. *A Rubik's cube*
673 4. *A basketball*
674 5. *A toy car*
675 6. *A backpack*

- 676 7. *An office chair*
- 677 8. *A pair of headphones*
- 678 9. *A yellow jacket*
- 679 10. *A laundry bag*

680 **REAL Lab Affordance Queries:**

- 681 1. *Something to use to disassemble or take apart a laptop*
- 682 2. *Something to use for cooling a CPU*
- 683 3. *Something to use for carrying books day to day*
- 684 4. *Something to use for temporarily securing a broken zipper*
- 685 5. *Something to use to help a student understand how a computer works*
- 686 6. *An object that is used in a sport involving rims and nets*
- 687 7. *Something to keep myself from getting distracted by loud noises*
- 688 8. *Something to help explain math proofs to a student*
- 689 9. *Something I can use to protect myself from the harsh winter in Canada*
- 690 10. *Something fun to pass the time with*

691 **REAL Lab Negation Queries:**

- 692 1. *A toy for someone who dislikes basketball*
- 693 2. *Shoes that you wouldn't wear to something formal*
- 694 3. *Something to protect me from the rain that's not an umbrella*
- 695 4. *Shoes that are not red and white*
- 696 5. *Something to make a cape with that's not green*
- 697 6. *Something to drink other than soda*
- 698 7. *Something to use for exercise other than weights*
- 699 8. *Something to wear unrelated to space or science*
- 700 9. *Something light to store belongings, not a backpack*
- 701 10. *Something to play with that's not a puzzle or colorful*

702 **A4 Navigation Experiments**

703 For our navigation experiments with the Jackal robot. Our robot is equipped with a VLP-16 lidar
704 and a forward-facing Realsense D435i camera. We begin by building a pointcloud of a lab space
705 using the onboard VLP-16 and Open3d SLAM [68]. The initial Jackal pointcloud does not include
706 task-relevant objects and is downprojected to a 2D costmap for navigation using the base Jackal
707 ROS stack.

708 We then stage two separate scenes with different objects: one for object search and another for
709 traversability estimation. In both cases, we map the scene with an Azure Kinect Camera and rely
710 on RTAB-Map [69] to obtain camera poses and the scene point cloud. We proceed to build a *Con-*
711 *ceptGraphs* representation and register the scene point cloud with the initial Jackal map. For our
712 navigation experiments, we only use the objects \mathbf{O}_T .

713 For object search queries, we use the LLM Planner described in Section A2 as part of a simple
714 state machine. The robot first attempts to go look at the 3D coordinates of the most relevant object
715 identified in \mathbf{O}_T by the LLM Planner. We then pass the onboard camera image to LLaVA [55] and
716 ask if the target object is in view. If not, we remove the target object from the scene graph and

717 ask the LLM Planner to provide a new likely location for the object in the scene with the following
718 GPT-4 system prompt:

```
The object described as 'description' is not in the scene. Perhaps
someone has moved it, or put it away. Let's try to find the object by
visiting the likely places, storage or containers that are appropriate
719 for the missing object (eg: a cabinet for a wineglass, or closet for a
broom). The new query is: find a likely container or storage space
where someone typically would have moved the object described as
'description'?
```

Listing 5: GPT system prompt for object localization.

720 For traversability estimation, we task GPT to classify a given object as traversable or non-traversable
721 based on its description and possible tags. The system prompt is:

```
You are a wheeled robot that can push a maximum of 5 pounds or 2.27 kg.
Can you traverse through or push an object identified as 'description'
722 with possible tags 'possible.tags'? Specifically, is it possible for you
to push the object out of its path without damaging yourself?
```

Listing 6: GPT-4 system prompt for traversability estimation.

723 We then take the pointclouds of each non-traversable objects and downproject them in the Jackal
724 costmap before launching the navigation episode. The goal is provided in this case as a specific pose
725 in the room.

726 For all experiments in this section, we run a local instance of LLaVA offboard on a desktop when
727 needed and otherwise use the GPT-4 API for LLM queries.

728 A5 Limitations

729 As indicated in Sec. 3.8, there are a few failure modes of *ConceptGraphs* that remain to be addressed
730 in subsequent work. In particular, the LLaVA-7B [55] model used for node captioning misclassifies
731 a non-negligible number of small objects as *toothbrushes* or *pairs of scissors*. We believe that using
732 more performant vision-language models, including instruction-finetuned variants of LLaVA [70]
733 can alleviate this issue to a large extent. This will, in turn, improve the node and edge precisions of
734 3D scene graphs beyond what we report in Table 1.

735 In this work, we do not explicitly focus on improving LLM-based planning over 3D scene graphs.
736 We refer the interested reader to concurrent work, SayPlan [52], for insights into how one might
737 leverage the hierarchy inherent in 3D scene graphs, for efficient planning.