TOPO-FIELD: TOPOMETRIC MAPPING WITH BRAIN INSPIRED HIERARCHICAL LAYOUT-OBJECT-POSITION FIELDS

Anonymous authors

Paper under double-blind review

Abstract

Mobile robots require comprehensive scene understanding to operate effectively in diverse environments, enriched with contextual information such as layouts, objects, and their relationships. While advancements like Neural Radiance Fields (NeRF) offer high-fidelity 3D reconstructions, they are computationally intensive and often lack efficient representations of traversable spaces essential for planning and navigation. In contrast, topological maps generated by LiDAR or visual SLAM methods are computationally efficient but lack the semantic richness necessary for a more complete understanding of the environment. Inspired by neuroscientific studies on spatial cognition, particularly the role of postrhinal cortex (POR) neurons that are strongly tuned to spatial layouts over scene content, this work introduces Topo-Field, a framework that integrates Layout-Object-Position (LOP) associations into a neural field and constructs a topometric map from this learned representation. LOP associations are modeled by explicitly encoding object and layout information, while a Large Foundation Model (LFM) technique allows for efficient training without extensive annotations. The topometric map is then constructed by querying the learned implicit neural representation, offering both semantic richness and computational efficiency. Empirical evaluations in multi-room apartment environments demonstrate the effectiveness of Topo-Field in tasks such as position attribute inference, query localization, and topometric planning, successfully bridging the gap between high-fidelity scene understanding and efficient robotic navigation.

006

007

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

028

029

1 INTRODUCTION

Mobile robots are rapidly moving from research labs to widespread use. For these robots to operate autonomously in complex environments, a deep understanding of their surroundings is crucial (Cadena et al., 2016). Efficient path planning and accurate identification of navigable spaces, along with detailed environmental reconstruction, will be key to enabling their deployment in real-world scenarios (Blochliger et al., 2018).

041 Recently, detailed environmental reconstruction has made great progress in producing lifelike 3D 042 reconstructions (Ullman, 1979; Forster et al., 2014; Dai et al., 2017; Tang & Tan, 2018), in which 043 NeRF (Mildenhall et al., 2020) is a prime instance. As improvements, works like (Zhi et al., 2021; 044 Fan et al., 2022; Xie et al., 2021) introduce semantic information for better scene understanding. 045 Further, features powered by Large-Foundation-Model (LFM)s, trained on massive datasets across various scenes, are employed with general knowledge for open scene understanding (Shafiullah 046 et al., 2022; Huang et al., 2023; Kerr et al., 2023). However, it is computationally demanding and 047 lacks global layout information using detailed neural fields for planning and navigation. 048

In contrast, existing topological maps for path planning and navigation in complex environments
 are often derived from LiDAR Simultaneous-Localization-and-Mapping (SLAM) using 3D dense
 submaps (Gomez et al., 2020) or visual SLAM by clustering free-space regions and extracting oc cupancy information from point clouds (Blochliger et al., 2018). While this approach increases path
 planning accuracy, computing topology with traditional methods comes with high computational
 costs and tends to strip away essential semantic information, reducing the robot's ability to fully



Figure 1: Illustration of the Topo-Field strategy and capabilities. Hierarchically dividing scene information into layout, object, and position to model them explicitly, layout-object-position associated knowledge enables robots with a topometric map representing the scene and planning navigable path to realize a more comprehensive spatial cognition.

understand and interpret the environment, which is critical for advanced autonomous functions such as language/image-prompted localization and navigation.

Neuroscientists have long discovered that animals process their surroundings using topological cod-079 ing, forming what is known as a "cognitive map" (Tolman, 1948), a concept embodied by place cells (O'Keefe & Dostrovsky, 1971). These place cells, along with spatial view cells (Rolls et al., 081 1998), respond to specific scene contents. More recently, research has shown that a population code 082 in the postrhinal cortex (POR) is strongly tuned to spatial layout rather than scene content (LaChance 083 et al., 2019), capturing spatial representations relative to environmental centers to form a high-level 084 cognitive map from egocentric perception to allocentric understanding (Zeng et al., 2022), unlike tra-085 ditional clustering from occupancy information (Blochliger et al., 2018) or Voronoi diagrams (Friedman et al., 2007). Inspired by this, we intuitively abstract the neural representations of space to build 087 topo-field in three key aspects: 1) The cognitive map corresponds to a topometric map, which uses 880 graph-like representations to encode relationships among its components, e.g. layouts and objects. 2) The population of place cells is analogous to a neural implicit representation with position en-089 coding, enabling location-specific responses. 3) POR, which prioritizes spatial layouts over content, 090 aligns with our spatial layout encoding of connected regions. We believe this approach makes a step 091 forward in applying mechanisms of spatial cognition in robotics. 092

093 To this end, this work proposes a Topo-Field, integrating the Layout-Object-Position (LOP) association into neural field training and constructing a topometric map based on the learned neural implicit 094 representation for hierarchical robotic scene understanding. By inputting RGB-D sequences, ob-095 jects and background contexts are encoded separately as contents and layout information to train a 096 neural field, forming a detailed scene representation. A contrast loss against features from LFMs 097 is employed, resulting in little need for annotation. Further, a topometric map is built based on co-098 observation relationship among frames, sampling points, and querying the learned field, which is efficient for navigable path planning. To validate the effectiveness of Topo-Field, we conduct quan-100 titative and qualitative experiments on several multi-room apartment scenes evaluating the abilities 101 including position attributes inference, text/image query localization, and planning. 102

- 103 Our contributions can be listed as follows:
- 104

075

We develop a brain-inspired Topo-Field, which combines detailed neural scene representation with high-level efficient topometric mapping for hierarchical robotic scene understanding and navigable path planning. Various quantitative and qualitative experiments on real-world datasets are conducted, showing high accuracy and low error in position at-

and path planning are also employed. 110 111 • We explain the theoretical basis and neuroscience reference to manage the hierarchical 112 encoding of spatial layouts and contents in the form of objects and connected regions, according to the spatial mechanism of cognitive map with POR population and place cells. 113 114 • We propose to learn a Layout-Object-Position associated implicit neural representation 115 with target features from separately encoded object instances and background contexts as 116 objects and layouts. The process is explicitly supervised by LFM-powered strategy with 117 little human labor. 118 119 • We propose a topometric map construction pipeline by querying the learned neural repre-120 sentation in a two-stage mapping and updating approach, leveraging LLM to validate edges 121 conducted among vertices. 122 123 124 2 **RELATED WORKS** 125 126 2.1 DENSE REPRESENTATION WITH NEURAL RADIANCE FIELD 127 128 Detailed 3D scene reconstruction has made great efforts in producing lifelike results, among which

tributes inference and multi-modal localization tasks. Examples of topometric construction

129 NeRF (Neural Radiance Fields) (Mildenhall et al., 2020) has widely attracted researchers' attention. 130 While numerous efforts improve the NeRF (Yu et al., 2021; Martin-Brualla et al., 2021; Zhu et al., 131 2022), a popular research direction is to integrate semantics with NeRF to achieve a more com-132 prehensive understanding of scenes (Zhi et al., 2021; Fan et al., 2022; Xie et al., 2021). Recently, 133 several robotic works have demonstrated that features from LFMs can be used for self-supervised 134 learning, which reduces the costly manual annotation (Shafiullah et al., 2022; Huang et al., 2023; 135 Kerr et al., 2023). However, the semantic feature fields learned in the above methods focus on object semantics but do not include layout-level features. RegionPLC (Yang et al., 2023) considered region 136 information by fusing multi-model features but with no explicit representation of layout features. In 137 contrast, in our work, CLIP (Radford et al., 2021) and Sentence-BERT (Reimers & Gurevych, 2019) 138 are employed to generate vision-language and semantic features for objects and layout respectively. 139 In addition to using object semantics, we annotate the belonging regions based on spatial layout and 140 regional division of scenes. Such annotations incur minimal cost but establish connections between 141 the position of 3D points, object semantics, and scene regions. 142

143 144

145

- 2.2 TOPOMETRIC MAP FOR SCENE STRUCTURE UNDERSTANDING
- 146 Using detailed neural fields for planning and navigation is computationally demanding, on the other 147 hand, hybrid topometric mapping has been known for its efficiency in terms of managing the information and being queried for downstream tasks (Zhang, 2015; Zhang et al., 2015; Garrote et al., 148 2018). It takes advantage of both metric maps and topological maps. Metric maps could refine the 149 local scale geometry accuracy and navigation plans while topological maps provide reliable global 150 topological cues and large-scale plans (Oleynikova et al., 2018; Badino et al., 2012). However, 151 most topological maps have not introduced information such as semantics. This makes it unsuit-152 able for language/image-guided planning tasks, which is a growing trend in scene representation 153 applications. Concept-graph (Gu et al., 2024) makes a step forward utilizing LFM to model the 154 object structure with a topo map. CLIO Maggio et al. (2024)built a task-driven scene graph inspired 155 by Information Bottleneck (IB) principle to form task-relevant clusters of primitives. At the same 156 time, HOV-SG Werby et al. (2024) proposed a hierarchical scene understanding pipeline, using fea-157 ture point cloud clustering of zero-shot embeddings in a fusion scheme and realizing the mapping 158 in an incremental approach. Unlike the incremental mapping and clustering-based graph construc-159 tion method, we propose to build the topometric map based on querying the trained neural field which serves as knowledge-like memory base, whose nodes and edges include attributes represent-160 ing object and layout information explicitly learned when training the specific neural representation 161 encoding.

162 2.3 SPATIAL UNDERSTANDING WITH LAYOUT INFORMATION 163

164 Generally, topology is built based on traditional clustering from occupancy information or Voronoi 165 diagrams (He et al., 2021), regardless of the contents and layout relationship. However, neuroscience findings suggest a mechanism to form a high-level cognitive map from egocentric perception to 166 allocentric representation (Zeng et al., 2022). Neuroscientists have long discovered that animals 167 process their surroundings using topological coding, forming a "cognitive map" (Tolman, 1948). 168 Place cells (O'Keefe & Dostrovsky, 1971), as the embodiment, together with spatial view cells show activity to contents (Rolls et al., 1998). Recently, Patrick et al. (LaChance et al., 2019) showed 170 that a population code in the POR is more strongly tuned to the spatial layout than to the content in 171 a scene. This suggests that there are specialized cells and signaling mechanisms to process layout 172 in the process of scene understanding, which captures the spatial layout of complex environments 173 to rapidly form a high-level cognitive map representation (Zeng et al., 2022). Inspired by the above 174 research, we propose that the spatial layout connected by regions, as a high-level abstract feature, 175 is closely related to the object contents and purposes of the scene. We mimic the neural scene 176 understanding mechanism by employing egocentric neural field knowledge to construct allocentric topometric map. 177

OVERVIEW 3

We propose to learn an implicit representation of a scene with the neural encoding approach by establishing associations between 3D positions and their corresponding layout and object features as the scene knowledge. Then, a topometric map is built with the learned neural field to form an efficient and queriable representation with a comprehensive understanding of the scene. Therefore, 185 we need to train a scene-dependent implicit function, denoted as

178 179

180 181

182

183

184

187

188

189 190

206

207 208

209

 $F: \mathbb{R}^3 \to \mathbb{R}^n,$ where for any 3D point P in space, F(P) is supposed to match with

$$\mathcal{E}\{(e_v, e_s)\} \in \mathbb{R}^n,\tag{2}$$

(1)

representing the layout-object-position associated embedding of that point. e_v and e_s are vision-191 language embedding and semantic embedding of image point where P is back-projected from. 192 CLIP (Radford et al., 2021) image encoder is introduced to encode e_v integrating the vision and 193 language feature space. Besides, the Sentence-BERT (Reimers & Gurevych, 2019) feature is also 194 introduced to encode e_s in this work. Because intuitively, unlike objects that can have similar ap-195 pearances within a certain category, region information often lacks specific visual appearances and 196 is closely related to semantic representations like the integration purpose of the scene and object semantics. Models trained on large-scale question-answering datasets can aid in understanding the 197 semantic relationships between regions and objects. Target feature processing and training strategy to match the embeddings to targets are described in Section 4.1 and 4.4. Applications utilizing the 199 learned field are discussed in Section 4.3.1. Based on the trained F, we aim to build a topometric 200 map denoted as 201

$$G = (V, E), \tag{3}$$

202 where vertices V include object vertices \mathbf{v}_o and region vertices \mathbf{v}_r and edges E include edges 203 between objects e_{o-o} , edges between regions e_{r-r} , and edges between object and region e_{o-r} . 204 The topological map architecture and construction pipeline are described in Section 4.3.2. 205

4 METHOD

4.1 TARGET FEATURE PROCESSING

210 (We clarify the formulation and polish for better understanding. The ground-truth label of layout regions is 211 described. More details are described for pixel-wise encoding of image and information in the target supervising 212 embeddings.) 213

RGB-D image sequences with poses are accepted as input to get the target layout-object-position 214 features for training F. For pure RGB image sequences, depth point clouds and camera poses can 215 also be estimated through methods like COLMAP (Schönberger & Frahm, 2016) or simultaneous

233

234 235 236



Figure 2: Pipeline of the Topo-Field. (a) The ground truth generation of layout-object-position vision-language and semantic embeddings for weakly-supervising. (b) The neural implicit network mapping 3D positions to target feature space. A contrastive loss is optimized against each other. (c) Topometric mapping process with trained neural field. (Formulation in the figure has been clarified.)

localization and mapping (SLAM). The only employed GT annotation is the layout distribution of 237 environment where the region of each 3D point P is denoted as $r_P \in R = \{r_1, r_2, \dots, r_q\}$, where 238 q is the number of regions. Such information is available in datasets like Matterport3D Chang et al. 239 (2017). However, in fact, partitioning the buildings needs little human labor, where in most human-240 made buildings spatial layouts are easily available divided by straight walls if not provided. As in 241 our practice, region annotation of a house with 8 rooms only takes 3 min by drawing lines from 242 top-down view according to walls to form a rule to separate (x, y) coordinates, bounding 3D points 243 to different regions.

244 For each image I, we employ Detic (Zhou et al., 2022) D to generate object instance patches with 245 number i, including bounding-boxes $B = \{b_1, b_2, \ldots, b_i\}$, masks $M = \{m_1, m_2, \ldots, m_i\}$, and 246 labels $L = \{l_1, l_2, \dots, l_i\}.$ 247

For object pixels p_o in instance mask j, CLIP (Radford et al., 2021) C is employed to compute 248 per-pixel features in mask b_i and Sentence-BERT (Reimers & Gurevych, 2019) S is employed to 249 process the semantic feature of l_j , prompted in the form of " l_j in r_{p_o} ". Given the related region r_{p_o} 250 of p_o , embedding of p_o can be denoted as $e_{p_o} = \{C(b_j), S(l_j, r_{p_o})\}$. 251

What's more, the background appearance is also considered which we proposed to include context 252 information for region layout. For background pixels p_b out of masks, per-pixel feature of the whole 253 image I is encoded. Its related region $r_{p_b} \in R = \{r_1, r_2, \dots, r_m\}$ is regarded as the text label and 254 embedding of p_b can be calculated as $e_{p_b} = \{C(I), S(r_{p_b})\}.$ 255

256 Then, pixel-wise embeddings are back-projected to 3D space based on depth and pose and averagely 257 counted to form a distilled 3D feature point cloud. Consequently, the target feature space $\mathcal{E}\{(e_n, e_s)\}$ consists of object and layout features, where (e_v, e_s) directs from $\{e_{p_o}, e_{p_b}\}_{p_o, p_b \in P}$. The pipeline 258 is shown in Fig. 2 259

260 Compared with previous implicit neural field methods, (e_v, e_s) includes (1) separately encoded 261 vision-language and semantic information by supervising embeddings from object and background 262 pixels. (2) region information consisted of vision-language embeddings from per-pixel image en-263 coding and semantic embeddings from region text labels. (3) context included object label in the form of " l_p in r_p ", where l_p and r_p is object label and region label at point p (e.g., cup in the kitchen). 264 Ablation studies of these improvements are conducted in Section 4 with more details. 265

266

- 267 4.2 Scene Neural Encoding
- (We clarify the formulation and include more details of the MHE and feature mapping head to get the high 269 dimension feature in neural representation.)



Figure 3: **Capabilities of the learned neural field**. (a) The attributes inference using position input. (b) The LOP association helped localization of text and image queries. (Formulation in the figure has been clarified.)

289 Our proposed Topo-Field involves an implicit mapping function to encode the 3D position into a 290 spatial vector representation $g: \mathbb{R}^3 \to \mathbb{R}^d$ and separate heads $h: \mathbb{R}^d \to \mathbb{R}^n$ processing encodings to 291 match the target feature space $\mathcal{E}\{(e_v, e_s)\}$. To select an appropriate implicit function, considering 292 that the target feature space includes object-level local features and layout-level region feature rep-293 resentation, we employ the Multi-scale Hash Encoding (MHE) introduced in Instant-NGP (Müller 294 et al., 2022) as g with d = 144. The feature pyramid structure used in MHE allows for consid-295 ering structural features ranging from coarse to fine in the spatial domain. Additionally, MHE has a faster training speed compared to traditional NeRF (Mildenhall et al., 2020) network structures. 296 For mapping the position encodings to the target feature space, we employ a unified and simple 297 Multi-Layer Perceptron (MLP) network structure. It includes heads $h_v : \mathbb{R}^d \to f_v$ for obtaining vision-language features and $h_s : \mathbb{R}^d \to f_s$ for semantic features, which together form the high 298 299 dimension embeddings $\{f_v, f_s\} \in \mathbb{R}^n$. The model is shown in Fig. 2. 300

In this way, given a posed RGB-D image, the target feature of each pixel is processed as mentioned in Section 4.1 denoted as $\mathcal{E}\{(e_v, e_s)\}$. At the same time the related pixel in depth image is backprojected into 3D space according to depth and pose value and processed by the above mentioned g, h to form $\{f_v, f_s\}$. A contrastive loss is conducted between $\{(e_v, e_s)\}$ and $\{f_v, f_s\}$ to train the neural representation. Training details are declared in Section 4.4.

306 307

308

285

286

287 288

4.3 TOPOMETRIC MAPPING

With the function and feature representation mentioned above, we can integrate 3D positions with the object and region information and construct a topometric map. The topo map construction process is formed in a mapping and updating strategy, while the implicit neural representation is introduced and queried as scene knowledge in this process. Detailed pipeline is introduced as follows.

313 314

315

4.3.1 KNOWLEDGE FROM LEARNED NEURAL FIELD

316 (We clarify the formulation and definition for better understanding.)

Position Attributes Inference. Using spatial 3D point *P* as input, assuming a collection of space regions *R* (e.g., "living room""bathroom""bedroom"...), we compute the vision-language features $C_R = \{C(r_1), C(r_2), \dots, C(r_m)\}$ and semantic features $S_R = \{S(r_1), S(r_2), \dots, S(r_m)\}$ using CLIP Radford et al. (2021) encoder *C* and Sentence-BERT Reimers & Gurevych (2019) encoder *S*, where *m* is the number of rooms. Then the cosine similarity between $F(P) = \{(f_v, f_s)\}_P$ and $\{C_R, S_R\}$ is calculated to find the most likely region to which *P* belongs. The inference process is shown in Fig. 3 (a). Similarly, the object information of *P* can be inferred with the same approach replacing the region set *R* with object set *O*.



Figure 4: **Qualitative comparison of text query localization** results among state-of-the-art methods and our method with text input in the form of "*object in the region*". Blue box shows the ground truth bounding box of object. Red box means miss-predicted box, while green box means the correctly predicted results.

Mathada	Sce	ne1	Sce	ne2	Sce	ne3	Scene4		
Methous	Dist.	Acc.	Dist.	Acc.	Dist.	Acc.	Dist.	Acc.	
CLIP-Field(2022)	2.97	0.24	3.35	0.21	2.98	0.20	3.06	0.17	
VLMaps(2023)	2.78	0.28	3.63	0.16	3.05	0.24	3.12	0.12	
LERF(2023)	2.86	0.32	2.82	0.11	3.49	0.17	3.04	0.20	
Topo-Field	0.92	0.85	0.86	0.84	0.36	0.95	0.27	0.97	
Text queries	1(00	1(.00 60			60		

Table 1: **Quantitative comparison of text query localization** results on different scenes from the Matterport3D dataset. The average distance (m) from the target to the localized point cloud and the accuracy evaluating whether predicted positions are in the correct region are used as metrics.

Localization with Text/Image Query. For natural language text input t (e.g., "cup in the bed-358 room"), most existing robotic scene representations struggle to locate specific objects of interest 359 (e.g., differentiating between cups in the living room and the bedroom). However, with our pro-360 posed Topo-Field that includes region information, we can calculate the cosine similarity between 361 $\{\mathcal{C}_t, \mathcal{S}_t\}$ and the embeddings $F(P^*) = \{(f_v, f_s)\}_{P^*}$ to find the most likely position of queries, 362 where P^* are sampled from 3D points set to train F. As for image input I, we can calculate the 363 cosine similarity of $\{C_I, S_I\}$ with $F(P^*) = \{(f_v, f_s)\}_{P^*}$ in the same way to find the 3D points set with highest similarity. The localization process of both text query and image query is shown in Fig. 364 3. 365

4.3.2 **TOPOMETRIC MAP CONSTRUCTION**

342

343

344

357

366

367

374

375

376

368
 369
 370
 371
 (We describe the topometric map construction pipeline step-by-step with more details, which can be regarded as a mapping and updating strategy. Definition and attributes of vertices and edges are clarified, including the acquisition process.)

As defined in Section 3, topometric map G = (V, E) consists of vertices and edges. We define a vertice **v** and edge **e** as

 \mathbf{v} : { id, node_type, class, bounding_box, caption}, (4)

e : { id, edge_type, start_node, end_node, relationship, caption }. (5)

377 Mimicking the mental representation of cognitive maps, we construct the topometric map in a **mapping and updating** strategy based on the learned Topo-Field *F*.

379

380

381

382

384

386

387

388

389

390 391

392

393

394

395

396 397

399

400

401

403 404

425

426 427



Figure 5: **Topometric map construction example.** The topometric map is represented as a graph from a top-down view according to the position of nodes. Map structure shows number of nodes and edges. A planning path from a seen view to target is shown as an example employing topometric map, the path is highlighted in green showing the related nodes and edges. Visited nodes are listed on the right. The line with gradient colors represents the waypoints based on the planning results while different colors represent different predicted regions of waypoints.

Mapping. we first averagely sample k points $P_{1,...,k}$ in the environment (each grid of $0.5m \times 0.5m$ 398 with a point in our practice) and infer their related regions according to Section 4.3.1. Supposing there are m regions in total $r_{1,...,m}$, we calculate the extent of each region in the bounding-box format according to positions of points within the same region. The topo map region vertice set is then initialized as $\mathbf{v}_r = {\mathbf{v}_{r_1}, \mathbf{v}_{r_2}, \dots, \mathbf{v}_{r_m}}$. For each \mathbf{v} , {id} is set, {node_type} is {region}, 402 {class} and {caption} is set according to the inferred region label, and {bounding_box} is set to

$$\{[Min(x), Max(x)], [Min(y), Max(y)], [Min(z), Max(z)]\}, (x, y, z) \in P_{1,\dots,k}.$$
(6)

405 On the other hand, while employing Detic Zhou et al. (2022) to detect object instances as mentioned 406 in Section 4.1, instances with high confidence (more than 60% in our practice) are recorded as object vertices candidates. For each \mathbf{v} , {node_type} is {object}, {class} and {caption} is set according 407 to the prediction result, and {bounding_box} is set according to the back-projected masked pixels 408 similar to equation 6. With the mapped nodes, we leverage LLM to describe the layouts with con-409 nectivity, distances, and relationships of regions and objects in JSON format based on the vertices' 410 attributes and poses. During this process, edges are built among vertices. For object-object edge 411 \mathbf{e}_{o-o} , we follow Gu et al. (2024) which mainly consider bounding-box overlap. For object-region 412 edge e_{o-r} , we consider an object belongs to the region if the object b-box is in the region b-box 413 and filter the unreasonable relation noise powered by LLM (e.g., it's almost impossible that a bike 414 is in bedroom). For region relationships, the adjacency and position relationship of region b-box is 415 considered. Examples of LLM prompts to build relationships and JSONs are listed in appendix for 416 reference. Fig. 2 shows the pipeline of metric-topological map construction.

417 **Updating**. RGB-D image sequence for training F or a newly captured sequence can be used for 418 constructed topometric map fine-tuning. For object vertices, if an object is detected by more than 3 419 frames in sequence, the object b-box will be compared with the constructed vertices. A new vertice 420 will be added if no vertice corresponds to it with the above-mentioned process. For region vertices, 421 we calculate embeddings $F(p_I)$ of sampled back-projected pixels p_I in each image I. $F(p_I)$ will 422 be matched with the constructed region set $r_{1,...,m}$, and extent of a region r will be updated if $F(p_I)$ matches $\{C_r, S_r\}$ and p_I exceeds the {bounding_box} extent of vertice \mathbf{v}_r . LLM to update edges 423 will be called each 50 frames. 424

4.4 TRAINING

The pipeline of ground truth data generation is described in Section 4.1 to train F. To fit the im-428 plicit representation introduced in Section 4.2 to the target feature space, we design the loss function 429 through a contrastive approach. For the vision-language feature optimization, the tempered similar-430 ity matrix on point P is denoted as 431

$$\operatorname{Sim}_{v} = \tau\{f_{v}\}_{P}\{e_{v}\}_{P},\tag{7}$$

Methods	Scene1	Scene2	Scene3	Scene4	Scene5	Scene6	Scene7	Scene8	Scene9	Scene1(
CLIP-Field(2022)	0.242	0.165	0.130	0.142	0.127	0.138	0.227	0.200	0.102	0.060
VLMaps(2023)	0.177	0.194	0.127	0.098	0.148	0.187	0.199	0.221	0.092	0.087
LERF(2023)	0.268	0.189	0.165	0.153	0.136	0.169	0.216	0.252	0.110	0.091
RegionPLC(2023)	0.290	0.202	0.173	0.168	0.152	0.154	0.243	0.248	0.086	0.088
Topo-Field	0.886	0.900	0.884	0.894	0.872	0.858	0.901	0.897	0.821	0.839
Position Samples	169k	185k	111k	112k	106k	176k	130k	121k	205k	211k

440 Table 2: Comparison of position attributes inference results on the test set of different scenes from the Matterport3D dataset. The average region prediction accuracy of sampled 3D points is 441 used as metric. 442

Methods	Scene1	Scene2	Scene3	Scene4
CLIP-Field(2022)	2.541	2.748	2.922	2.651
VLMaps*(2023)	2.112	1.894	1.181	1.595
LERF(2023)	1.276	1.175	1.148	1.129
Topo-Field	0.742	0.830	0.374	0.327

Table 3: Quantitative comparison of image query localization results with other methods. The similarity weighted average distance (m) between the target view point cloud and the predicted point cloud is evaluated. VLMaps* is a self-implemented version with image localization ability.

where τ is the temperature term, $\{f_v\}_P$ and $\{e_v\}_P$ is the calculated implicit representation feature and target embedding according to P. Using cross-entropy loss, the vision-language loss can be calculated as

$$\mathcal{L}_{v} = -exp(-\operatorname{dist}_{P})(H(\operatorname{Sim}_{v}) + H(\operatorname{Sim}_{v}^{T})),$$
(8)

where dist_P is the distance from P to camera, and H is the cross-entropy function. For the semantic loss, similarity on points P can be calculated as

$$\operatorname{Sim}_{s} = \tau\{f_{s}\}_{P}\{e_{s}\}_{P}.$$
(9)

460 Similarly, semantic loss can be denoted as

$$\mathcal{L}_{s} = -\mathrm{conf}(H(\mathrm{Sim}_{s}) + H(\mathrm{Sim}_{s}^{T})), \tag{10}$$

463 where conf is the prediction confidence from the detection model. The total loss is computed by:

$$\mathcal{L} = \mathcal{L}_v + \mathcal{L}_s. \tag{11}$$

In our experiments, an NVIDIA RTX3090 GPU is utilized and the batch size is set to 12544 to 466 maximize the capability of our VRAM. As model instances, CLIP with SwinB is employed in Detic Zhou et al. (2022), CLIP Radford et al. (2021) encoder is ViT-B/32 and Sentence-BERT Reimers 468 & Gurevych (2019) encoder is all-mpnet-base-v2. The MHE has 18 levels of grids and the dimen-469 sion of each grid is 8, with loq_2 hash map size of 20 and only 1 hidden MLP layer of size 600. 470 We train the neural implicit network for 100 epochs with optimizer Adam, employing a decayed 471 learning rate of 1e - 4 and 3e - 3 decay rate. Each epoch contains 3e6 samples. Codes and scripts 472 are released in supplementary for reproducibility.

473 474

475

449

450

451 452

453

454

455 456

457

458

459

461

462

464 465

467

5 EXPERIMENTAL RESULTS

476 Our experiments are conducted on real-world datasets to validate the established layout-object-477 position association. The data environment is of single-floor residential buildings with multiple 478 rooms which is the common working scenario of household robots widely studied in this field. We 479 employed Matterport3D (Chang et al., 2017) as well as apartment environment (Zhu et al., 2022) 480 dataset to demonstrate that our approach can be generalized in diverse scenarios.

481

482 5.1 POSITION ATTRIBUTES INFERENCE 483

To demonstrate the built LOP association integrates positions with layout features, we designed 484 experiments that accept 3D positions as input to infer the region information. For quantitative eval-485 uation, we divided the RGB-D sequences into training and testing sets. The Topo-Field is trained



Figure 6: **Qualitative comparison of image query localization** results in heatmaps form among state-of-the-art methods and our method with image input. Our approach localizes the position of queried image in an exact smaller range.

Methods	Scene1	Scene2	Scene3	Scene4
CLIP-Field	0.242	0.165	0.130	0.142
Baseline1	0.852	0.891	0.863	0.874
Baseline2	0.865	0.887	0.872	0.879
Baseline3	0.872	0.891	0.875	0.886
Topo-Field	0.886	0.900	0.884	0.894

Table 4: **Ablation of target feature processing pipeline** of the neural field construction. The average region prediction accuracy of sampled points from different scenes on the Matterport3D dataset is used as the metric. Declarations of baselines are clarified in Fig. 7 and Section 4.

519

520

514

504

505 506 507

> according to Section 4.4 on the training set and tested in the test set. As the region inference task can be treated as a multi-class classification task for each input, the accuracy, precision, and F1-score are used as metrics. Tab. A.8 shows the region inference results on 10 real-world scenes in Matterport 3D (Chang et al., 2017) with different scales and layouts indicating the average accuracy exceeds 85%.

5.2 LOCALIZATION WITH PROMPT QUERIES

Localization with Text Queries: For objects of the same category in different regions, we input 526 the textual description of the target object in the form of "object in the region" and infer the specific 527 location of the target, comparing the results with the predictions from current state-of-the-art visual-528 language algorithms. Fig. 4 demonstrates the advancements of Topo-Field in object localization 529 tasks involving region information, which allows for the localization of specific target objects based 530 on the description and features of the region, while other methods confuse objects from different 531 regions. Tab. 1 shows the quantitative results on 4 scenes of different layouts compared to other methods with an average accuracy of more than 88% and less distance from targets. For the metrics, 532 the average distance (m) of predicted point cloud and ground truth point cloud is evaluated, together 533 with counting whether the center of predicted points is in the correct room. Ground truth comes 534 from the Matterport3D provided object instance labels. More results can be seen in the appendix. 535

Localization with Image Queries: To validate the help of region information in the image view
localization task. We localize the images from the test set in the trained Topo-Field. Selected views
include representative objects of the scene (e.g., TV in the living room) and views with similarlooking objects or context (e.g., bathroom washbasin and kitchen sink) which is challenging. The
localization results are shown in Fig. 6 in the form of heatmaps and Tab. 3 shows the quantitative



Figure 7: Ablation of our LOP information encoding and feature fusion strategy for target features. (CLIP-Field encoding strategy has been added for comparison.)

results which evaluates the weighted average distance of the target view and localized point cloud among all samples in a scene, using similarity as weight. VLMaps* is a self-implemented version, because origin VLMaps (Huang et al., 2023) does not implement the image localization task. To align with CLIP-Field (Shafiullah et al., 2022) and our work, the LSeg (Li et al., 2022) used in VLMap (Huang et al., 2023) is replaced by CLIP (Radford et al., 2021). The results show that Topo-Field constrains the localization results to a smaller range in the exact region. We sampled more than 40 images on each scene from Apartment (Zhu et al., 2022) and Matterport3D (Chang et al., 2017) dataset. By drawing the predicted camera view on the top-down view, we estimated the localization precision and found that most views can be ranged into a specific view on the target field of view, while other methods struggle to get precise results.

5.3 TOPOMETRIC MAP CONSTRUCTION

Fig. 4.3.2 shows an example of the built topometric map. Layout region nodes, object nodes with
bounding boxes, and entrance nodes connecting regions are shown with edges representing relationships. A planned navigable path is shown in the graph from an observed view in family room
to the TV room sofa in green. The path planning A* algorithm is employed to explore the topological structure to generate waypoints between nodes, and the waypoints are generated with the
planning API in Habitat Simulator (Savva et al., 2019) and shown in a line with gradient colors,
while different colors indicate different predicted regions of the waypoints.

5.4 ABLATION STUDY

577 Fig. 7 and Tab 4. show the ablation of our neural field LOP encoding strategy and feature fusion 578 where 1) CLIP-Field Shafiullah et al. (2022) means the origin feature encoding strategy that doesn't 579 explicitly consider the layout features. 2) Baseline1 is our first crude approach that directly super-580 vises the learned embedding from the encoded objects with region semantics. 3) Baseline2 encodes 581 the region description to the target vision-language and semantic feature space for supervision. 4) 582 Baseline3 takes the background pixels into account with the region labels. 5) Topo-Field further considers the context of the layout when supervising the object label semantics. These four main 583 versions of our numerous iterations of trying are listed as examples to show our work on the neural 584 field encoding of LOP association. 585

586 587

554

555 556

557

558

559

560

561

562

563

565 566

567

575

576

6 CONCLUSION AND LIMITATIONS

588

Inspired by postrhinal cortex (POR) neurons that prioritize spatial layouts over scene content for cognitive mapping, we propose Topo-Field, which integrates Layout-Object-Position (LOP) associations into a neural field and constructs a topometric map from the learned field for hierarchical robotic scene understanding. However, there are some limitations: 1) While we present a pipeline for topometric map construction, querying and path planning are currently implemented using traditional methods (e.g. A*). Future work will explore using large language models to integrate seman-

tic information for more advanced path planning. 2) Real-world deployment on mobile robots for
 long-term navigation is needed. 3) Future research will focus on updating and editing the topometric
 map to accommodate environmental changes.

598 599 REFERENCES

- Hernán Badino, Daniel Huber, and Takeo Kanade. Real-time topometric localization. In 2012 IEEE
 International conference on robotics and automation, pp. 1635–1642. IEEE, 2012.
- Fabian Blochliger, Marius Fehr, Marcin Dymczyk, Thomas Schneider, and Rol Siegwart. Topomap: Topological mapping and navigation based on visual slam maps. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 3818–3825. IEEE, 2018.
- Cesar Cadena, Luca Carlone, Henry Carrillo, Yasir Latif, Davide Scaramuzza, José Neira, Ian Reid,
 and John J Leonard. Past, present, and future of simultaneous localization and mapping: Toward
 the robust-perception age. *IEEE Transactions on robotics*, 32(6):1309–1332, 2016.
- Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva,
 Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor
 environments. *International Conference on 3D Vision (3DV)*, 2017.
- Angela Dai, Matthias Nießner, Michael Zollhöfer, Shahram Izadi, and Christian Theobalt. Bundle fusion: Real-time globally consistent 3d reconstruction using on-the-fly surface reintegration.
 ACM Transactions on Graphics (ToG), 36(4):1, 2017.
- Zhiwen Fan, Peihao Wang, Yifan Jiang, Xinyu Gong, Dejia Xu, and Zhangyang Wang. Nerf-sos: Any-view self-supervised object segmentation on complex scenes, 2022. URL https: //arxiv.org/abs/2209.08776.
- Christian Forster, Matia Pizzoli, and Davide Scaramuzza. Svo: Fast semi-direct monocular visual odometry. In 2014 IEEE international conference on robotics and automation (ICRA), pp. 15–22. IEEE, 2014.
- Stephen Friedman, Hanna Pasula, and Dieter Fox. Voronoi random fields: Extracting topological
 structure of indoor environments via place labeling. In *IJCAI*, volume 7, pp. 2109–2114, 2007.
- Luís Garrote, Cristiano Premebida, David Silva, and Urbano J Nunes. Hmaps-hybrid height-voxel maps for environment representation. In 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pp. 1197–1203. IEEE, 2018.
- Clara Gomez, Marius Fehr, Alex Millane, Alejandra C Hernandez, Juan Nieto, Ramon Barber, and
 Roland Siegwart. Hybrid topological and 3d dense mapping through autonomous exploration for
 large indoor environments. In 2020 IEEE International Conference on Robotics and Automation
 (ICRA), pp. 9673–9679. IEEE, 2020.
- Qiao Gu, Ali Kuwajerwala, Sacha Morin, Krishna Murthy Jatavallabhula, Bipasha Sen, Aditya
 Agarwal, Corban Rivera, William Paul, Kirsty Ellis, Rama Chellappa, et al. Conceptgraphs:
 Open-vocabulary 3d scene graphs for perception and planning. In 2024 IEEE International Con-*ference on Robotics and Automation (ICRA)*, pp. 5021–5028. IEEE, 2024.
- ⁶³⁷
 ⁶³⁸
 ⁶³⁸
 ⁶³⁹
 ⁶³⁹
 ⁶³⁹
 ⁶³⁹
 ⁶³⁹
 ⁶³⁰
 ⁶³⁰
 ⁶³¹
 ⁶³¹
 ⁶³²
 ⁶³³
 ⁶³³
 ⁶³⁴
 ⁶³⁵
 ⁶³⁵
 ⁶³⁶
 ⁶³⁶
 ⁶³⁷
 ⁶³⁷
 ⁶³⁷
 ⁶³⁸
 ⁶³⁹
 ⁶³⁹
 ⁶³⁹
 ⁶³⁹
 ⁶³⁹
 ⁶³⁹
 ⁶³⁰
 ⁶³⁰
 ⁶³¹
 ⁶³²
 ⁶³²
 ⁶³³
 ⁶³⁵
 ⁶³⁵
 ⁶³⁶
 ⁶³⁶
 ⁶³⁷
 ⁶³⁷
 ⁶³⁸
 ⁶³⁹
 ⁶³⁹
 ⁶³⁹
 ⁶³⁹
 ⁶³⁰
 ⁶³¹
 ⁶³²
 ⁶³²
 ⁶³⁵
 ⁶³⁵
 ⁶³⁵
 ⁶³⁵
 ⁶³⁶
 ⁶³⁷
 ⁶³⁷
 ⁶³⁸
 ⁶³⁹
 ⁶³⁹
- Chenguang Huang, Oier Mees, Andy Zeng, and Wolfram Burgard. Visual language maps for robot navigation. In 2023 IEEE International Conference on Robotics and Automation (ICRA), pp. 10608–10615. IEEE, 2023.
- Justin* Kerr, Chung Min* Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. Lerf: Language embedded radiance fields. In *International Conference on Computer Vision (ICCV)*, 2023.
- 647 Patrick A. LaChance, Travis P. Todd, and Jeffrey S. Taube. A sense of space in postrhinal cortex. *Science*, 365(6449):eaax4192, 2019.

691

648	Boyi Li, Kilian O Weinberger, Serge Belongie, Vladlen Koltun, and Rene Ranftl. Language-driven
649	semantic segmentation. In International Conference on Learning Representations, 2022. URL
650	https://openreview.net/forum?id=RriDjddCLN.
651	

- Dominic Maggio, Yun Chang, Nathan Hughes, Matthew Trang, Dan Griffith, Carlyn Dougherty, 652 Eric Cristofalo, Lukas Schmid, and Luca Carlone. Clio: Real-time task-driven open-set 3d scene 653 graphs. IEEE Robotics and Automation Letters, 9(10):8921-8928, 2024. doi: 10.1109/LRA. 654 2024.3451395. 655
- 656 Ricardo Martin-Brualla, Noha Radwan, Mehdi S. M. Sajjadi, Jonathan T. Barron, Alexey Dosovit-657 skiy, and Daniel Duckworth. NeRF in the Wild: Neural Radiance Fields for Unconstrained Photo Collections. In CVPR, 2021. 658
- 659 Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and 660 Ren Ng. NeRF: Representing scenes as neural radiance fields for view synthesis. In The European 661 Conference on Computer Vision (ECCV), 2020. 662
- Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics prim-663 itives with a multiresolution hash encoding. ACM transactions on graphics (TOG), 41(4):1-15, 664 2022. 665
- 666 John O'Keefe and Jonathan Dostrovsky. The hippocampus as a spatial map: preliminary evidence 667 from unit activity in the freely-moving rat. *Brain research*, 1971. 668
- Helen Oleynikova, Zachary Taylor, Roland Siegwart, and Juan Nieto. Sparse 3d topological graphs 669 for micro-aerial vehicle planning. In 2018 IEEE/RSJ International Conference on Intelligent 670 Robots and Systems (IROS), pp. 1-9. IEEE, 2018. 671
- 672 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, 673 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual 674 models from natural language supervision. In *International conference on machine learning*, pp. 675 8748-8763. PMLR, 2021.
- Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-677 networks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Lan-678 guage Processing and the 9th International Joint Conference on Natural Language Processing 679 (EMNLP-IJCNLP), pp. 3982–3992, Hong Kong, China, November 2019. Association for Com-680 putational Linguistics. doi: 10.18653/v1/D19-1410. URL https://aclanthology.org/ 681 D19-1410. 682
- Edmund T Rolls, Alessandro Treves, Robert G Robertson, Pierre Georges-François, and Stefano 683 Panzeri. Information about spatial view in an ensemble of primate hippocampal cells. Journal of 684 Neurophysiology, 79(4):1797–1813, 1998. 685
- 686 Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, 687 Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, Devi Parikh, and Dhruv Batra. Habitat: A 688 Platform for Embodied AI Research. In Proceedings of the IEEE/CVF International Conference 689 on Computer Vision (ICCV), 2019. 690
- Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In Conference on Computer Vision and Pattern Recognition (CVPR), 2016. 692
- 693 Nur Muhammad Mahi Shafiullah, Chris Paxton, Lerrel Pinto, Soumith Chintala, and Arthur 694 Szlam. Clip-fields: Weakly supervised semantic fields for robotic memory. arXiv preprint arXiv:2210.05663, 2022.
- Chengzhou Tang and Ping Tan. Ba-net: Dense bundle adjustment network. arXiv preprint 697 arXiv:1806.04807, 2018. 698
- 699 Edward C Tolman. Cognitive maps in rats and men. *Psychological review*, 55(4):189, 1948. 700
- Shimon Ullman. The interpretation of structure from motion. Proceedings of the Royal Society of 701 London. Series B. Biological Sciences, 203(1153):405-426, 1979.

702 Abdelrhman Werby, Chenguang Huang, Martin Büchner, Abhinav Valada, and Wolfram Burgard. 703 Hierarchical open-vocabulary 3d scene graphs for language-grounded robot navigation. Robotics: 704 Science and Systems, 2024. 705 Christopher Xie, Keunhong Park, Ricardo Martin-Brualla, and Matthew Brown. Fig-nerf: Figure-706 ground neural radiance fields for 3d object category modelling. In International Conference on 3D Vision (3DV), 2021. 708 709 Jihan Yang, Runyu Ding, Zhe Wang, and Xiaojuan Qi. Regionplc: Regional point-language con-710 trastive learning for open-world 3d scene understanding. arXiv preprint arXiv:2304.00962, 2023. 711 712 Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. PlenOctrees for 713 real-time rendering of neural radiance fields. In ICCV, 2021. 714 Taiping Zeng, Bailu Si, and Jianfeng Feng. A theory of geometry representations for spatial naviga-715 tion. Progress in Neurobiology, 211:102228, 2022. 716 717 Qiwen Zhang. Autonomous indoor exploration and mapping using hybrid metric/topological maps. 718 McGill University (Canada), 2015. 719 720 Qiwen Zhang, Ioannis Rekleitis, and Gregory Dudek. Uncertainty reduction via heuristic search 721 planning on hybrid metric/topological map. In 2015 12th Conference on Computer and Robot 722 Vision, pp. 222–229. IEEE, 2015. 723 Shuaifeng Zhi, Tristan Laidlow, Stefan Leutenegger, and Andrew Davison. In-place scene labelling 724 and understanding with implicit scene representation. In Proceedings of the International Con-725 ference on Computer Vision (ICCV), 2021. 726 727 Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and Ishan Misra. Detecting 728 twenty-thousand classes using image-level supervision. In ECCV, 2022. 729 730 Zihan Zhu, Songyou Peng, Viktor Larsson, Weiwei Xu, Hujun Bao, Zhaopeng Cui, Martin R. Os-731 wald, and Marc Pollefeys. Nice-slam: Neural implicit scalable encoding for slam. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), June 2022. 732 733 734 А APPENDIX 735 736 SCENE PARTATION EXAMPLE A.1 738 The scene can be partitioned into different regions using walls as dividers and lines can be aligned 739

The scene can be partitioned into different regions using walls as dividers and lines can be aligned to these walls. This is similar in most scenarios, making the annotation of scene regions a straightforward task as shown in Fig. A1.

740

741 742

743 744 745

746

747 748

749

750

751 752

753



Figure A1: Using walls as dividers to associate lines with them, the scene can be divided into variousregions and 3D points can be labeled with related regions easily.

A.2 VISION-LANGUAGE EMBEDDINGS SIMILARITY OF REGION AND OBJECTS

758 To demonstrate that the relationship of the vision-language and semantic embeddings for different 759 regions is related to our intuition, we compare the similarity in region-region and object-region form 760 and show the results in Fig. A2. It can be seen that based on general knowledge, cognitively related 761 regions (e.g., the dining room and kitchen) and object-region pairs (e.g., sink and kitchen) are also 762 more correlated in the vision-language and semantic feature spaces.



Figure A2: The similarity of a set of region embeddings (as shown in a) and object-region embeddings (as shown in b). The left graph shows the vision-language embedding similarity and the right one shows the semantic embedding similarity.

798 799 800

796

797

801 802

A.3 ABLATION STUDY

To explicitly encode the region information, we apply the LVM to process the background pixels out of the object bounding box and LLM to encode the region label text. What's more, for object pixels, object label text is combined with the region text in the form of 'object in the region' before being encoded by LLM. To ablate the contribution of vision-language embeddings from CLIP and semantic embeddings from Sentence-BERT in encoding region features, we compare different weight settings between the v-s embeddings when inferring the regions with 3D position inputs. Results are shown in Fig. A3. It can be seen that both vision-language embeddings and semantic embeddings are indispensable, and weight settings with the greatest results are used for Topo-Field.



Figure A3: Ablation results on the accuracy of region prediction on Matterport3DChang et al. (2017) with 3D positions input. The w/o BG stands for not encoding background pixels to get region embeddings, and v-s weight ablates the weight of vision-language and semantic embeddings in the embeddings similarity contribution. Error bars show the results among samples from different scenes in Matterport3DChang et al. (2017).

A.4 HIERARCHICAL APPROACH COMPARISON

Hierarchical scene representation is widely studied with numerous tasks, mainly employing scalable receptive fields and representations to fine-tune results of scalable objects and local relations. As Fig.A.8 shows, VoxFusion introduced octree map with various voxel sizes, LERF employed feature pyramids. As far as we know, few of them explicitly consider the layout level information and the association with objects and positions. This idea comes from recent neuroscience findings, and similar theory has not yet been introduced in scene representations.

A.5 TOPOMETRIC SEARCH FOR PLANNING

We employ a simple A* approach for planning. Given a topometric graph G, the start point p, and the target destination object text t. First, the belonged region r of p is inferred according to the main paper. The existing objects nodes embeddings are compared with the encoded visual-language and semantic embeddings of t to find the target object node o. At the same time, if the region of destination object r_d is declared, the search process would be more simple by directly search among region nodes. Here lists the pseudocode of the employed A*.

A.6 TOPOMETRIC MAP NODES EXAMPLES

We list the attributes of nodes and edges in the topometric map as example here in Listing 1 - 4, including the object nodes, region nodes, and edges.

```
1
            "id": 0,
858
     2
            "node_type": region,
859
            "bbox_extent": [
860
                 4.163309999999999
861
                 4.207343,
     6
862
     7
                 2.53566175
863
     8
            1.
            "bbox_center": [
     9
```

Alg	gorithm 1 AStar(G, r, o)	
1:	$openSet \leftarrow \{r\}$	▷ Set of nodes to be evaluate
2:	$cameFrom \leftarrow \{\}$	▷ Mapping of nodes to their parent node
3:	$gScore[r] \leftarrow 0$	▷ Cost from start along best known pat
4:	$fScore[r] \leftarrow h(r, o)$	▷ Estimated total cost from start to goa
5:	while openSet is not empty do	C
6:	$current \leftarrow$ node in $openSet$ with lowest $fScor$	e value
7.	if $current = o$ then	
8.	return ReconstructPath(<i>cameFrom</i> , o)	
9. 9.	end if	
10.	remove current from open Set	
11.	for each neighbor n of current do	
11.	tentative $CScore \leftarrow aScore[current] + d($	current n)
12.	if $tantativeCScore < gScore[carrent] + a($	
13:	In term $uiveGScore < gScore[n]$ then aam a Enom[n] < aam on t	
14:	$camerrom[n] \leftarrow carrent$	
15:	$gScore[n] \leftarrow tentativeGScore$	
16:	$fScore[n] \leftarrow gScore[n] + h(n, o)$	
1/:	If n not in openSet then	
18:	add n to openSet	
19:	end if	
20:	end if	
21:	end for	
22:	end while	
23:	return "No path found"	
24:	function RECONSTRUCTPATH(cameFrom, current	nt)
25:	$path \leftarrow [current]$	
26:	while current is in cameFrom do	
27:	$current \leftarrow cameFrom[current]$	
28:	insert <i>current</i> at the beginning of <i>path</i>	
29:	end while	
30:	return <i>path</i>	
31:	end function	

-8.821845, 2.6915385, 1.259409125], "class": "bedroom", "caption": "A bedroom at the northwest of the house with warm lighting. Main objects include a bed in the center, a large closet, and a dresser at the corner." 16 }, Listing 1: Region node 1 { "id": 1, "node_type": object, "bbox_extent": [0.3569, 0.2297, 0.101.8], "bbox_center": [0.3222, -1.1108, -0.5062], "class": "picture", "caption": "A white framed picture hanging on the wall." 16 }, Listing 2: Object node 1 { "id": 0, "node_type": Entrance, "bbox_extent": [0.5, 1.6, 2.8,], "bbox_center": [-3.244, -0.276, 0.487], "class": "Entrance", "caption": "Entrance connecting bedroom and living room." 16 }, Listing 3: Entrance node 1 { "id": 2, "edge_type": region_entrance, "start_node": { "id": 0, "node_type": region, "bbox_extent": [4.163309999999999, 4.207343, 2.53566175], "bbox_center": [-8.821845, 2.6915385, 1.259409125

```
972
     16
                 ],
973
                 "region_tag": "bedroom"
     17
974
    18
            },
975
    19
            "end_node": {
                "id": 0,
     20
976
                "node_type": Entrance,
     21
977
                 "bbox_extent": [
     22
978
                     0.5,
     23
979
     24
                     1.6,
980
     25
                     2.8,
981
     26
                 ],
                 "bbox_center": [
982
                     -3.244,
     28
983
                     -0.276,
     29
984
     30
                     0.487
985
     31
                ],
                 "class": "Entrance",
986
     32
                 "caption": "Entrance connecting bedroom and living room."
     33
987
     34
            },
988
            "relationship": connected,
     35
989
            "position_relation": "b to the southeast of a",
     36
            "position_reason": "The x-coordinate of the center of bbox of
990
     37
            end_node (-3.244) is larger than that of start_node (-8.821845), and
991
            the y-coordinates of the center of bbox of end_node (-0.276) is less
992
            than that of start_node (4.207343). Therefore, b is to the southeast
993
           of a."
994
            "caption": "The pathway from bedroom to living room."
     38
995
     39 },
996
                                    Listing 4: Region entrance edge
997
998
999
     1 {
            "id": 2,
1000
     2
            "node_type": object_region,
1001
     3
            "start_node": {
     4
1002
                "id": 7,
     5
1003
                "node_type": object,
     6
1004
                 "bbox_extent": [
     7
                     2.155,
1005
     8
                     2.052,
     9
1006
                     0.883
     10
1007
     11
                 ],
1008
     12
                 "bbox_center": [
1009
     13
                     5.598,
1010 14
                     2.566,
                     0.136
1011 <sup>15</sup>
1012 <sup>16</sup>
                ],
                 "class": "bed",
     17
1013
                "caption": "a bed with a white comforter and a pillow"
     18
1014 19
            },
1015 20
            "end_node": {
                 "id": 0,
1016 <sup>21</sup>
1017<sup>22</sup>
                 "node_type": region,
                 "bbox_extent": [
     23
1018
                     4.163309999999999,
     24
1019 25
                     4.207343,
                     2.53566175
1020 26
1021 <sup>27</sup>
                 ],
    28
                 "bbox_center": [
1022
                     -8.821845,
     29
1023
     30
                     2.6915385,
1024 31
                     1.259409125
1025 32
                 ],
                 "class": "bedroom"
     33
```

```
1026
                "caption": "A bedroom at the northwest of the house with warm
    34
1027
          lighting. Main objects include a bed in the center, a large closet,
1028
           and a dresser at the corner."
1029 35
           },
           "relationship": belong,
1030 <sup>36</sup>
           "position_relation": "a in the center of b",
    37
1031
           "caption": "According to the bbox center position and extent, the bed
    38
1032
            is in the center of bedroom."
1033
    39 },
1034
                                   Listing 5: Object region edge
1035
1036
1037
       A.7 PROMPT EXAMPLE FOR REGION NODE CONNECTIVITY DESCRIPTION
1038
1039
       With topometric mapped nodes, we leverage LLM to describe the connectivity of nodes according
       to the general knowledge and bounding box 3D position. In listing 5, here we provide a prompt
1040
      example to describe the connectivity relationship between content objects and regions and set up the
1041
      edge.
1042
1043
     1 {
1044
    2 DEFAULT_PROMPT_POST = """
1045 3 You are an excellent graph managing agent. Given a graph nodes set of an
          environment,
1046
1047 <sup>4</sup> you can explore the relationships of nodes with their attributes and
        build edges among
1048
     5 them.
1049 6
1050 7 The input is a list of JSONS describing two types of nodes, including the
           object and
1051
    8 region. You need to produce a JSON string (and nothing else) and set up
1052
          edges between them with keys: "relationship", "position_relation" and
1053
            "caption".
1054 9
1055 10 Each of the JSON fields will have the following fields:
1056 11 1. id: a unique number
1057 <sup>12</sup> 2. node_type: type of this node
    13 3. bbox_extent: the 3D bounding box extents
1058 14 4. bbox_center: the 3D bounding box center
1059 15 5. class: an extremely brief description
1060 16 6. caption: a sentence describing node attributes in detail
1061 <sup>17</sup>
1062 <sup>18</sup> Produce a "relationship" field that best describes the relationship of
          the object node and region node. Set "false" if the object is not
1063
          related to the area or is not reasonable, the relationship is refused
1064
           . Produce a
1065 19 "position_relation" field describing the position relationship between
          object and region according to their
1066
1067 20 bounding box information in the 3D space. Before producing the "
          position_relation" field, produce a "caption" field that explains why
1068
            the "position_relation" field is reasonable.
1069 <sub>21</sub>
1070 22 The built edges should include following fields:
1071 23 1. id: a unique number of each edge in order
1072 24 2. node_type: according to the connected node type in the form "
          start_node\_end_node"
1073 25 3. start_node: keep JSON values of the object node unchanged
1074 26 4. end_node: keep JSON values of the region node unchanged
1075 27 5. relationship
1076 28 6. position_relation
1077 <sup>29</sup> 7. caption
    30 """
1078
                         Listing 6: Prompt example to set up edge with nodes.
1079
```

1080 A.8 ADDITIONAL EXPERIMENT RESULTS

Additional experiments results of object localization using text query inputs and view localization
 using image query inputs. Also, a table is provided showing the metric on exactly each region class
 from 4 scenes in Matterport3D dataset.

1085													
1086	Dogions		Scene1		Scene2		Scene3			Scene4			
1087	Regions	Acc.	Pre.	F1	Acc.	Pre.	F1	Acc.	Pre.	F1	Acc.	Pre.	F1
1088	Living Room	0.948	0.970	0.959	0.870	0.881	0.875	0.778	0.810	0.793	0.902	0.949	0.925
1089	Bedroom	0.943	0.825	0.880	0.925	0.923	0.924	0.687	0.767	0.725	0.920	0.870	0.894
1090	Bathroom	0.466	0.680	0.554	0.903	0.898	0.901	0.875	0.463	0.605	0.797	0.831	0.814
1091	Dining Room	-	-	-	0.961	0.794	0.870	0.774	0.732	0.752	0.933	0.887	0.910
1092	Lobby	0.681	0.941	0.790	0.853	0.951	0.899	0.978	0.510	0.671	0.855	0.698	0.769
1093	Family Room	-	-	-	-	-	-	0.903	0.571	0.700	0.926	0.936	0.931
1094	Kitchen	0.994	0.654	0.789	0.788	0.836	0.811	0.833	0.833	0.833	0.758	0.854	0.803
1095	Office	-	-	-	0.969	0.848	0.905	-	-	-	0.953	0.883	0.917
1096	Toilet	-	-	-	-	-	-	0.900	0.711	0.795	-	-	-
1097	Avg. Acc./Samples	0.886 / 169k			0.900 / 185k			0.884 / 111k			0.894 / 112k		

Table 5: Region prediction results on the test set of different scenes from the Matterport3DChang et al. (2017) dataset. Accuracy, precision, and F1 score are used as metrics.







Figure A6: Text query localization on scene 17DRP5sb8fyChang et al. (2017).



Figure A8: Text query localization on scene HxpKQynjfinChang et al. (2017).



Figure A9: Image query localization on scene 2t7WUuJeko7Chang et al. (2017).



Figure A10: Image query localization on scene 17DRP5sb8fyChang et al. (2017).



Figure A11: Image query localization on scene HxpKQynjfinChang et al. (2017).

CLIP-Field (Shafiullah et al., 2022)

VLMaps* (Huang et al., 2023)

Topo-Field(Ours)

GT View

Input Image