

000  
001  
002  
003  
004  
005  
006  
007  
008

054

055

056

# 3D Concept Learning and Reasoning from Multi-View Images

057

058

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

079

080

081

082

083

084

085

086

087

088

089

090

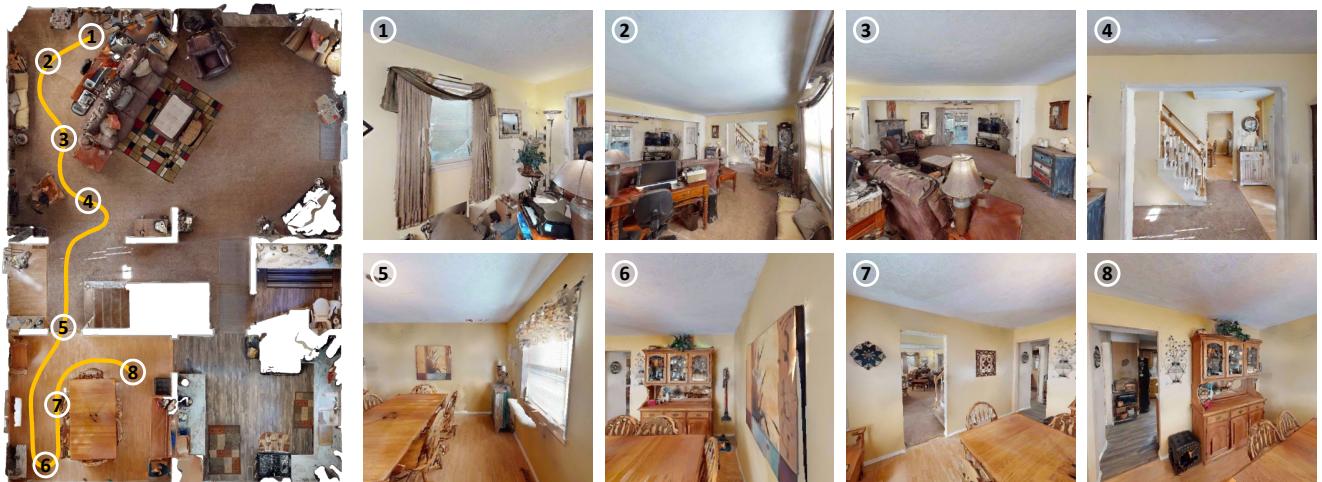
091

092

093

Anonymous CVPR submission

Paper ID 6248

**Concept:**

Q: Are there any **televisions**?  
A: Yes

**Counting:**

Q: How many **chairs** are **close** to the **table** in the room with **plant** on the **cabinet**?  
A: 6

**Relation:**

Q: Facing the **computer** from the **curtain**, is there a **lamp** on the **right**?  
A: Yes

**Comparison:**

Q: Are there **fewer pictures** in the **larger** room than the other room?  
A: No

Q: Is there a **sofa** in the room with a **printer**?  
A: Yes

Q: How many rooms have **sofas**?  
A: 1

Q: What's on the **cabinet** in the **smaller** room?  
A: Plant

Q: Is the **computer** **closer** to a **printer** or a **lamp**?  
A: Printer

Figure 1. An exemplar scene with multi-view images and question-answer pairs of our 3DMV-VQA dataset. 3DMV-VQA contains four question types: concept, counting, relation, comparison. Orange words denote semantic concepts; blue words denote the relations.

## Abstract

Humans are able to accurately reason in 3D by gathering multi-view observations of the surrounding world. Inspired by this insight, we introduce a new large-scale benchmark for 3D multi-view visual question answering (3DMV-VQA). This dataset is collected by an embodied agent actively moving and capturing RGB images in an environment using the Habitat simulator. In total, it consists of approximately 5k scenes, 600k images, paired with 50k questions. We evaluate various state-of-the-art models for visual reasoning on our benchmark and find that they all perform poorly. We suggest that a principled approach for 3D reasoning from multi-view

images should be to infer a compact 3D representation of the world from the multi-view images, which is further grounded on open-vocabulary semantic concepts, and then to execute reasoning on these 3D representations. As the first step towards this approach, we propose a novel 3D concept learning and reasoning (3D-CLR) framework that seamlessly combines these components via neural fields, 2D pre-trained vision-language models, and neural reasoning operators. Experimental results suggest that our framework outperforms baseline models by a large margin, but the challenge remains largely unsolved. We further perform an in-depth analysis of the challenges and highlight potential future directions.

108

## 1. Introduction

109

Visual reasoning, the ability to composite rules on internal representations to reason and answer questions about visual scenes, has been a long-standing challenge in the field of artificial intelligence and computer vision. Several datasets [23, 33, 69] have been proposed to tackle this challenge. However, they mainly focus on visual reasoning on 2D single-view images. Since 2D single-view images only cover a limited region of the whole space, such reasoning inevitably has several weaknesses, including occlusion, and failing to answer 3D-related questions about the entire scene that we are interested in. As shown in Fig. 1, it's difficult, even for humans, to count the number of chairs in a scene due to the object occlusion, and it's even harder to infer 3D relations like "closer" from a single-view 2D image.

124

On the other hand, there's strong psychological evidence that human beings conduct visual reasoning in the underlying 3D representations [55]. Recently, there have been several works focusing on 3D visual question answering [2, 16, 62, 64]. They mainly use traditional 3D representations (*e.g.*, point clouds) for visual reasoning. This is inconsistent with the way human beings perform 3D reasoning in real life. Instead of being given an entire 3D representation of the scene at once, humans will actively walk around and explore the whole environment, ingesting image observations from different views and converting them into a holistic 3D representation that assists them in understanding and reasoning about the environment. Such abilities are crucial for many embodied AI applications, such as building assistive robots.

138

To this end, we propose the novel task of 3D visual reasoning from multi-view images taken by active exploration of an embodied agent. Specifically, we generate a large-scale benchmark, 3DMV-VQA (3D multi-view visual question answering), that contains approximately 5k scenes and 50k question-answering pairs about these scenes. For each scene, we provide a collection of multi-view image observations. We generate this dataset by placing an embodied agent in the Habitat-Matterport environment [47], which actively explores the environment and takes pictures from different views. We also obtain scene graph annotations from the Habitat-Matterport 3D semantics dataset (HM3DSem) [61], including ground-truth locations, segmentations, semantic information of the objects, as well as relationships among the objects in the environments, for model diagnosis. To evaluate the models' 3D reasoning abilities on the entire environment, we design several 3D-related question types, including concept, counting, relation and comparison.

156

Given this new task, the key challenges we would like to investigate include: 1) how to efficiently obtain the compact visual representation to encode crucial properties (*e.g.*, semantics and relations) by integrating all incomplete observations of the environment in the process of active exploration for 3D visual reasoning? 2) How to ground the semantic con-

cepts on these 3D representations that could be leveraged for downstream tasks, such as visual reasoning? 3) How to infer the relations among the objects, and perform step-by-step reasoning?

As the first step to tackling these challenges, we propose a novel model, 3D-CLR (3D Concept Learning and Reasoning). First, to efficiently obtain a compact 3D representation from multi-view images, we use a neural-field model based on compact voxel grids [57] which is both fast to train and effective at storing scene properties in its voxel grids. As for concept learning, we observe that previous works on 3D scene understanding [1, 3] lack the diversity and scale with regard to semantic concepts due to the limited amount of paired 3D-and-language data. Although large-scale vision-language models (VLMs) have achieved impressive performances for zero-shot semantic grounding on 2D images, leveraging these pretrained models for effective open-vocabulary 3D grounding of semantic concepts remains a challenge. To address these challenges, we propose to encode the features of a pre-trained 2D vision-language model (VLM) into the compact 3D representation defined across voxel locations. Specifically, we use the CLIP-LSeg [37] model to obtain features on multi-view images, and propose an alignment loss to map the features in our 3D voxel grid to 2D pixels. By calculating the dot-product attention between the 3D per-point features and CLIP language embeddings, we can ground the semantic concepts in the 3D compact representation. Finally, to answer the questions, we introduce a set of neural reasoning operators, including FILTER, COUNT, RELATION operators and so on, which take the 3D representations of different objects as input and output the predictions.

We conduct experiments on our proposed 3DMV-VQA benchmark. Experimental results show that our proposed 3D-CLR outperforms all baseline models a lot. However, failure cases and model diagnosis show that challenges still exist concerning the grounding of small objects and the separation of close object instances. We provide an in-depth analysis of the challenges and discuss potential future directions. To sum up, we have the following contributions in this paper.

- We propose the novel task of 3D concept learning and reasoning from multi-view images.
- By having robots actively explore the embodied environments, we collect a large-scale benchmark on 3D multi-view visual question answering (3DMV-VQA).
- We devise a model that incorporates a neural radiance field, 2D pretrained vision and language model, and neural reasoning operators to ground the concepts and perform 3D reasoning on the multi-view images. We illustrate that our model outperforms all baseline models.
- We perform an in-depth analysis of the challenges of this new task and highlight potential future directions.

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216     **2. Related Work**  
217  
218     **Visual Reasoning** There have been numerous tasks focusing  
219     on learning visual concepts from natural language, including  
220     visually-grounded question answering [18, 19], text-image  
221     retrieval [59] and so on. Visual reasoning has drawn much  
222     attention recently as it requires human-like understanding of  
223     the visual scene. A wide variety of benchmarks have been  
224     created over the recent years [7, 8, 23, 27, 33, 69]. However,  
225     they mainly focus on visual reasoning from 2D single-view  
226     images, while there's strong psychological evidence that  
227     human beings perform visual reasoning on the underlying  
228     3D representations. In this paper, we propose the novel task  
229     of visual reasoning from multi-view images, and collect  
230     a large-scale benchmark for this task. In recent years, nu-  
231     merous visual reasoning models have also been proposed,  
232     ranging from attention-based methods [5, 30], graph-based  
233     methods [28], to models based on large pretrained vision-  
234     language model [9, 38]. These methods model the reasoning  
235     process implicitly with neural networks. Neural-symbolic  
236     methods [6, 40, 65] explicitly perform symbolic reasoning  
237     on the objects representations and language representations.  
238     They use perception models to extract 2D masks as a first  
239     step, and then execute operators and ground concepts on  
240     these pre-segmented masks, but are limited to a set of pre-  
241     defined concepts on simple scenes. [26] proposes to use the  
242     feature vectors from occupancy networks [42] to do visual  
243     reasoning in the 3D space. However, they also use a syn-  
244     thetic dataset, and learn a limited set of semantic concepts  
245     from scratch. We propose to learn 3D neural field features  
246     from 2D multi-view real-world images, and incorporate a  
247     2D VLM for open-vocabulary reasoning.

248     **3D Reasoning** Understanding and reasoning about 3D  
249     scenes has been a long-standing challenge. Recent works  
250     focus on leveraging language to explore 3D scenes, such  
251     as object captioning [3, 4] and object localization from lan-  
252     guage [1, 17, 29]. Our work is mostly related to 3D Visual  
253     Question Answering [2, 16, 62, 64] as we both focus on an-  
254     swering questions and reasoning about 3D scenes. However,  
255     these works use point clouds as 3D representations, which  
256     diverts from the way human beings perform 3D reasoning.  
257     Instead of being given an entire 3D representation all at once,  
258     human beings would actively move and explore the environ-  
259     ment, integrating multi-view information to get a compact  
260     3D representation. Therefore, we propose 3D reasoning from  
261     multi-view images. In addition, since 3D assets paired with  
262     natural language descriptions are hard to get in real-life sce-  
263     narios, previous works struggle to ground open-vocabulary  
264     concepts. In our work, we leverage 2D VLMs for zero-shot  
265     open-vocabulary concept grounding in the 3D space.

266     **Embodied Reasoning** Our work is also closely related to  
267     Embodied Question Answering (EQA) [11, 67] and Interac-  
268     tive Question Answering (IQA) [22, 35], which also involve  
269     an embodied agent exploring the environment and answering

270     the question. However, the reasoning mainly focuses on the  
271     outcome or the history of the navigation on 2D images and  
272     does not require a holistic 3D understanding of the environ-  
273     ment. There are also works [12, 20, 51, 54, 56, 68] targeting  
274     instruction following in embodied environments, in which an  
275     agent is asked to perform a series of tasks based on language  
276     instructions. Different from their settings, for our benchmark  
277     an embodied agent actively explores the environment and  
278     takes multi-view images for 3D-related reasoning.

279     **Neural Fields** Our approach utilizes neural fields to pa-  
280     rameterize an underlying 3D compact representations of  
281     scenes for reasoning. Neural field models (*e.g.*, [43]) have  
282     gained much popularity since they can reconstruct a vol-  
283     umetric 3D scene representation from a set of images. Recent  
284     works [21, 24, 57, 66] have pushed it further by using clas-  
285     sic voxel-grids to explicitly store the scene properties (*e.g.*,  
286     density, color and feature) for rendering, which allows for  
287     real-time rendering and is utilized by this paper. Neural fields  
288     have also been used to represent dynamic scenes [14, 44],  
289     appearance [43, 45, 49, 53, 63], physics [34], robotics [32, 52],  
290     acoustics [39] and more general multi-modal signals [13].  
291     There are also some works that integrate semantics or lan-  
292     guage in neural fields [31, 60]. However, they mainly fo-  
293     cus on using language for manipulation, editing or gen-  
294     eration. [26] leverages neural descriptor field [52] for 3D  
295     concept grounding. However, they require ground-truth oc-  
296     cupancy values to train the neural field, which can not be  
297     applied to real-world scenes. In this paper, we propose to  
298     leverage voxel-based neural radiance field [57] to get the  
299     compact representations for 3D visual reasoning.

## 3. Dataset Generation

### 3.1. Multi-View Images

301     Our dataset includes 5k 3D scenes from the Habitat-  
302     Matterport 3D Dataset (HM3D) dataset [47], and approx-  
303     imately 600k images rendered from the 3D scenes. The  
304     images are rendered via Habitat [50, 58].

305     **Scene Generation** We build our benchmark on top of the  
306     HM3DSem dataset [61], which is a large-scale dataset of  
307     3D real-world indoor scenes with densely annotated seman-  
308     tics. It consists of 142,646 object instance annotations across  
309     216 3D spaces and 3,100 rooms within those spaces. HM3D  
310     dataset uses texture information to annotate pixel-accurate  
311     object boundaries, which provides large-scale object anno-  
312     tations and ensures the scale, quality, and diversity of 3D  
313     visual reasoning questions of our benchmark.

314     To construct a benchmark that covers questions of differ-  
315     ent difficulty levels, it's crucial that we include 3D scenes of  
316     different scales in our benchmark. We start with single rooms  
317     in HM3D scenes, which has an appropriate amount of seman-  
318     tic concepts and relationships to base some simple questions  
319     on. To get the scale of single rooms, we calculate bounding  
320     322     323

324 boxes of rooms according to floor instance segmentations.  
 325 We then proceed to generate bounding boxes for scenes with  
 326 multiple adjacent rooms. For more complex holistic scene un-  
 327 derstanding, we also include whole-house scenes, which may  
 328 contain tens of rooms. Overall, the 3DMV-VQA benchmark  
 329 contains three levels of scenes (2000 single-room scenes,  
 330 2000 multi-room scenes and 100 whole-house scenes).

331 **Image Rendering** After we get the bounding box of each  
 332 scene, we load the scene into the Habitat simulator. We  
 333 also put a robot agent with an RGB sensor at a random  
 334 initial point in the bounding box. The data is collected via  
 335 exploration of the robot agent. Specifically, at each step of  
 336 the data collection process, we sample a navigable point and  
 337 make the agent move to the point along the shortest path.  
 338 When the agent has arrived at a point, we rotate the agent  
 339 30° along z-axis for 12 times so that the agent can observe  
 340 the 360° view of the scene at the position. It can also look up  
 341 and down, with a random mild angle from [-10°, 10°] along  
 342 the x-axis. A picture is taken each time the agent rotates to  
 343 a new orientation. In total 12 pictures are taken from each  
 344 point. While traveling between points, the robot agent further  
 345 takes pictures. We also exploit a policy such that when the  
 346 camera is too far from or too close to an object and thus the  
 347 agent cannot see anything, we discard the bad-view images.  
 348

### 349 3.2. Questions and Answers

350 We pair each scene with machine-generated questions  
 351 from pre-defined templates. All questions are open-ended  
 352 and can be answered with a single word (samples in Fig. 1).

353 **Concepts and Relationships** To generate questions and  
 354 answers, we utilize the semantic annotations of HM3DSem  
 355 [61] to get the semantic concepts and their bounding boxes,  
 356 as well as the bounding boxes of the rooms. We merge seman-  
 357 tic concepts with similar meanings (e.g., L-shaped sofa to sofa,  
 358 desk chair / computer chair e.g. to chair). We also define 11  
 359 relationships: inside, above, below, on the top of, close,  
 360 far, large, small, between, on the left, and on the right.  
 361 Before generating questions, we first generate a scene graph  
 362 for each scene containing all concepts and relationships.

363 **Question Types** We define four types of questions: concept,  
 364 counting, relation and comparison.

365 • **Concept.** Conceptual questions query whether there's an  
 366 object of a certain semantic concept in the scene, or whether  
 367 there's a room containing the objects of the semantic con-  
 368 cept.

369 • **Counting.** Counting-related questions ask about how many  
 370 instances of a semantic concept are in the scene, or how  
 371 many rooms contain objects of the semantic concept.

372 • **Relation.** Relational questions ask about the 11 rela-  
 373 tionships and their compositions. Based on the number of rela-  
 374 tions in a question, we have one-hop to three-hop questions  
 375 for the relation type.

376 • **Comparison.** The comparison question type focuses on the  
 377 comparison of two objects, two semantic concepts or two  
 378 rooms. It can be combined with the relational concepts to  
 379 compare two objects (e.g., larger, closer to, more left etc).  
 380 It also compares the number of instances of two semantic  
 381 concepts, or the number of objects of certain concepts in  
 382 different rooms.  
 383

384 **Bias Control.** Similar to previous visual reasoning bench-  
 385 marks [26, 33], we use machine-generated questions since  
 386 the generation process is fully controllable so that we can  
 387 avoid dataset bias. Questions are generated from pre-defined  
 388 templates, and transformed into natural language questions  
 389 with associated semantic concepts and relationships from  
 390 the scene. We manually define 41 templates for question  
 391 generation. We use depth-first search to generate questions.  
 392 We perform bias control based on three perspectives: tem-  
 393 plate counts, answer counts, and concept counts. For select-  
 394 ing templates, we sort the templates each time we generate  
 395 a question to ensure a balanced question distribution. We  
 396 force a flat answer distribution for each template by rejec-  
 397 tion sampling. Specifically, once we generate a question and  
 398 an answer, if the number of the questions having the same  
 399 answer and template is significantly larger than other an-  
 400 swers, we discard it and continue searching. Once we find  
 401 an answer that fits in the ideal answer distribution, we stop  
 402 the depth-first searching for this question. We also force a  
 403 flat concept distribution for each template using the same  
 404 method. In addition to controlling the number of concepts  
 405 mentioned in the templates, we also control the number of  
 406 relation tuples consisting of the same concept sets.  
 407

## 4. Method

408 Fig. 2 illustrates an overview of our framework. Specifi-  
 409 cally, our framework consists of three steps. First, we learn  
 410 a 3D compact representation from multi-view images using  
 411 neural field. And then we propose to leverage pre-trained 2D  
 412 vision-and-language model to ground concepts on 3D space.  
 413 This is achieved by 1) generating 2D pixel features using  
 414 CLIP-LSeg; 2) aligning the features of 3D voxel grid and 2D  
 415 pixel features from CLIP- LSeg [37]; 3) dot-product attention  
 416 between the 3D features and CLIP language features [37].  
 417 Finally, to perform visual reasoning, we propose neural rea-  
 418 soning operators, which execute the question step by step on  
 419 the 3D compact representation and outputs a final answer.  
 420 For example, we use FILTER operators to ground semantic  
 421 concepts on the 3D representation, GET\_INSTANCE to get  
 422 all instances of a semantic class, and COUNT\_RELATION to  
 423 count how many pairs of the two semantic classes have the  
 424 queried relation.

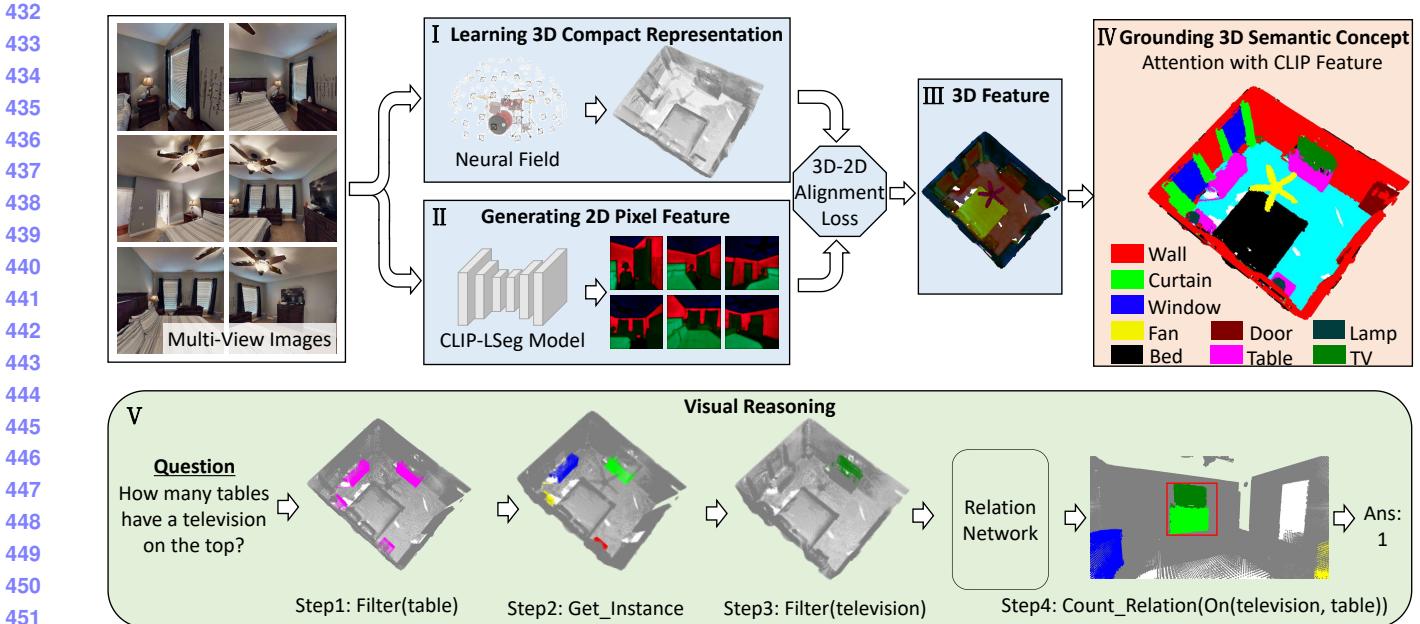


Figure 2. An overview of our 3D-CLR framework. First, we learn a 3D compact scene representation from multi-view images using neural fields (I). Second, we use CLIP-LSeg model to get per-pixel 2D features (II). We utilize a 3D-2D alignment loss to assign features to the 3D compact representation (III). By calculating the dot-product attention between the 3D per-point features and CLIP language embeddings, we could get the concept grounding in 3D (IV). Finally, the reasoning process is performed via a set of neural reasoning operators, such as FILTER, GET\_INSTANCE and COUNT\_RELATION (V). Relation operators are learned via relation networks.

#### 4.1. Learning 3D Compact Scene Representations

Neural radiance fields [43] are capable of learning a 3D representation that can reconstruct a volumetric 3D scene representation from a set of images. Voxel-based methods [21, 24, 57, 66] speed up the learning process by explicitly storing the scene properties (*e.g.*, density, color and feature) in its voxel grids. We leverage Direct Voxel Grid Optimization (DVGO) [57] as our backbone for 3D compact representation for its fast speed. DVGO stores the learned density and color properties in its grid cells. The rendering of multi-view images is by interpolating through the voxel grids to get the density and color for each sampled point along each sampled ray, and integrating the colors based on the rendering alpha weights calculated from densities according to quadrature rule [41]. The model is trained by minimizing the L2 loss between the rendered multi-view images and the ground-truth multi-view images. By extracting the density voxel grid, we can get the 3D compact representation (*e.g.*, By visualizing points with density greater than 0.5, we can get the 3D representation as shown in Fig. 2 I.)

#### 4.2. 3D Semantic Concept Grounding

Once we extract the 3D compact representation of the scene, we need to ground the semantic concepts for reasoning from language. Recent work from [26] has proposed to ground concepts from paired 3D assets and question-answers. Though promising results have been achieved on

synthetic data, it is not feasible for open-vocabulary 3D reasoning in real-world data, since it is hard to collect large-scale 3D vision-and-language paired data. To address this challenge, our idea is to leverage pre-trained 2D vision and language model [46, 48] for 3D concept grounding in real-world scenes. But how can we map 2D concepts into 3D neural field representations? Note that 3D compact representations can be learned from 2D multi-view images and that each 2D pixel actually corresponds to several 3D points along the ray. Therefore, it's possible to get 3D features from 2D per-pixel features. Inspired by this, we first add a feature voxel grid representation to DVGO, in addition to density and color, to represent 3D features. We then apply CLIP-LSeg [37] to learn per-pixel 2D features, which can be attended to by CLIP concept embeddings. We use an alignment loss to align 3D features with 2D features so that we can perform concept grounding on the 3D representations.

**2D Feature Extraction.** To get per-pixel features that can be attended by concept embeddings, we use the features from language-driven semantic segmentation (CLIP-LSeg) [37], which learns 2D per-pixel features from a pre-trained vision-language model (*i.e.*, [46]). Specifically, it uses the text encoder from CLIP, trains an image encoder to produce an embedding vector for each pixel, and calculates the scores of word-pixel correlation by dot-product. By outputting the semantic class with the maximum score of each pixel, CLIP-LSeg is able to perform zero-shot 2D semantic segmentation.

**3D-2D Alignment.** In addition to density and color, we also

540 store a 512-dim feature in each grid cell in the compact  
 541 representation. To align the 3D per-point features with 2D  
 542 per-pixel features, we calculate an L1 loss between each  
 543 pixel and each 3D point sampled on the ray of the pixel.  
 544 The overall L1 loss along a ray is the weighted sum of all  
 545 the pixel-point alignment losses, with weights same as the  
 546 rendering weights:  $\mathcal{L}_{\text{feature}} = \sum_{i=1}^K w_i (\|f_i - F(\mathbf{r})\|)$ , where  
 547  $\mathbf{r}$  is a ray corresponding to a 2D pixel,  $F(\mathbf{r})$  is the 2D feature  
 548 from CLIP-LSeg,  $K$  is the total number of sampled points  
 549 along the ray and  $f_i$  is the feature of point  $i$  by interpolating  
 550 through the feature voxel grid,  $w_i$  is the rendering weight.  
 551

552 **Concept Grounding through Attention.** Since our feature  
 553 voxel grid representation is learnt from CLIP-LSeg, by  
 554 calculating the dot-product attention  $\langle f, v \rangle$  between per-  
 555 point 3D feature  $f$  and the CLIP concept embeddings  $v$ ,  
 556 we can get zero-shot view-independent concept grounding  
 557 and semantic segmentations in the 3D representation, as is  
 558 presented in Fig. 2 IV.

### 559 4.3. Neural Reasoning Operators

560 Finally, we use the grounded semantic concepts for 3D  
 561 reasoning from language. We first transform questions into a  
 562 sequence of operators that can be executed on the 3D repre-  
 563 sentation for reasoning. We adopt a LSTM-based semantic  
 564 parser [65] for that. As [26, 40], we further devise a set of  
 565 operators which can be executed on the 3D representation.  
 566 Please refer to **Appendix** for a full list of operators.

567 **Filter Operators.** We filter all the grid cells with a certain  
 568 semantic concept.

569 **Get Instance Operators.** We implement this by utilizing  
 570 DBSCAN [15], an unsupervised algorithm which assigns  
 571 clusters to a set of points. Specifically, given a set of points  
 572 in the 3D space, it can group together the points that are  
 573 closely packed together for instance segmentation.

574 **Relation Operators.** We cannot directly execute the relation  
 575 on the 3D representation as we have not grounded rela-  
 576 tions. Thus, we represent each relation using a distinct  
 577 neural module (which is practical as the vocabulary of  
 578 relations is limited [36]). We first concatenate the voxel grid  
 579 representations of all the referred objects and feed them into  
 580 the relation network. The relation network consists of three  
 581 3D convolutional layers and then three 3D deconvolutional  
 582 layers. A score is output by the relation network indicating  
 583 whether the objects have the relationship or not. Since vanilla  
 584 3D CNNs are very slow, we use Sparse Convolution [10] in-  
 585 stead. Based on the relations asked in the questions, different  
 586 relation modules are chosen.

## 587 5. Experiments

### 588 5.1. Experimental Setup

589 **Evaluation Metric.** We report the visual question answering  
 590 accuracy on the proposed 3DMV-VQA dataset w.r.t the four

591 types of questions. The train/val/test split is 7:1:2.

592 **Implementation Details** For 3D compact representations,  
 593 we adopt the same architectures as DVGO, except skipping  
 594 the coarse reconstruction phase and directly training the fine  
 595 reconstruction phase. After that, we freeze the density voxel  
 596 grid and color voxel grid, for the optimization of the feature  
 597 voxel grid only. The feature grid has a world size of 100  
 598 and feature dim of 512. We train the compact representa-  
 599 tions for 100,000 iterations and the 3D features for another  
 600 20,000 iterations. For LSeg, we use the official demo model,  
 601 which has the ViT-L/16 image encoder and CLIP’s ViT-B/32  
 602 text encoder. We follow the official script for inference and  
 603 use multi-scale inference. For DBSCAN, we use an epsilon  
 604 value of 1.5, minimum samples of 2, and we use L1 as the  
 605 clustering method. For the relation networks, each relation  
 606 is encoded into a three-layer sparse 3D convolution network  
 607 with hidden size 64. The output is then fed into a one-layer  
 608 linear network to produce a score, which is normalized by  
 609 sigmoid function. We use cross-entropy loss to train the rela-  
 610 tion networks, and we use the one-hop relational questions  
 611 with “yes/no” answers to train the relation networks.

### 612 5.2. Baselines

613 Our baselines range from vanilla neural networks,  
 614 attention-based methods, fine-tuned from large-scale VLM,  
 615 and graph-based methods, to neural-symbolic methods.

- **LSTM.** The question is transferred to word embeddings  
 616 which are input into a word-level LSTM [25]. The last  
 617 LSTM hidden state is fed into a multi-layer perceptron  
 618 (MLP) that outputs a distribution over answers. This  
 619 method is able to model question-conditional bias since it  
 620 uses no image information.
- **CNN+LSTM.** The question is encoded by the final hidden  
 621 states from LSTM. We use a resnet-50 to extract frame-  
 622 level features of images and average them over the time  
 623 dimension. The features are fed to an MLP to predict the  
 624 final answer. This is a simple baseline that examines how  
 625 vanilla neural networks perform on 3DMV-VQA.
- **3D-Feature+LSTM.** We use the 3D features we get from  
 626 3D-2D alignment and downsample the voxel grids using  
 627 3D-CNN as input, concatenated with language features  
 628 from LSTM and fed to an MLP.
- **MAC [30].** MAC utilizes a Memory, Attention and Com-  
 629 position cell to perform iterative reasoning process. Like  
 630 CNN+LSTM, we use the average pooling over multi-view  
 631 images as the feature map.
- **MAC(V).** We treat the multi-view images along a trajectory  
 632 as a video. We modify the MAC model by applying a  
 633 temporal attention unit across the video frames to generate  
 634 a latent encoding for the video.

	Methods	Concept	Counting	Relation	Comparison	Overall	
648	Q-type (rand.)	49.4	10.7	21.6	49.2	26.4	702
649	Q-type (freq.)	50.8	11.3	23.9	50.3	28.2	703
650	LSTM	53.4	15.3	24.0	55.2	29.8	704
651	CNN+LSTM	57.8	22.1	35.2	59.7	37.8	705
652	MAC	62.4	19.7	47.8	62.3	46.7	706
653	MAC(V)	60.0	24.6	51.6	65.9	50.0	707
654	NS-VQA	59.8	21.5	33.4	61.6	38.0	708
655	ALPRO	65.8	12.7	42.2	68.2	43.3	709
656	LGCN	56.2	19.5	35.5	66.7	39.1	710
657	3D-Feature+LSTM	61.2	22.4	49.9	61.3	48.2	711
658	3D-CLR (Ours)	<b>66.1</b>	<b>41.3</b>	<b>57.6</b>	<b>72.3</b>	<b>57.7</b>	712
659							713
660							714
661							715
662							716

Table 1. Question-answering accuracy of 3D visual reasoning baselines on different question types.

- **NS-VQA** [65]. This is a 2D version of our 3D-CLR model. We use CLIP-LSeg to ground 2D semantic concepts from multi-view images, and the relation network also takes the 2D features as input. We execute the operators on each image and max pool from the answers to get our final predictions.
- **ALPRO** [38]. ALPRO is a video-and-language pre-training framework. A transformer model is pretrained on large webly-source video-text pairs and can be used for downstream tasks like Video Question answering.
- **LGCN** [28]. LGCN represents the contents in the video as a location-aware graph by incorporating the location information of an object into the graph construction.

### 5.3. Experimental Results

**Result Analysis.** We summarize the performances for each question type of baseline models in Table 1. All models are trained on the training set until convergence, tuned on the validation set, and evaluated on the test set. We provide detailed analysis below.

First, for the examination of language-bias of the dataset, we find that the performance of LSTM is only slightly higher than random and frequency, and all other baselines outperform LSTM a lot. This suggests that there's little language bias in our dataset. Second, we observe that encoding temporal information in MAC (*i.e.*, MAC(V)) is better than average-pooling of the features, especially in counting and relation. This suggests that average-pooling of the features may cause the model to lose information from multi-view images, while attention on multi-view images helps boost the 3D reasoning performances. Third, we also find that fine-tuning on large-scale pretrained model (*i.e.*, ALPRO) has relatively high accuracies in concept-related questions, but for counting it's only slightly higher than the random baseline, suggesting that pretraining on large-scale video-language dataset may improve the model's perception ability, but does not provide the model with the ability to tackle with more difficult reasoning types such as counting. Next, we

find that LGCN has poor performances on the relational questions, indicating that building a location-aware graph over 2D objects still doesn't equip the model with 3D location reasoning abilities. Last but not least, we find that 3D-based baselines are better than their 2D counterparts. 3D-Feature+LSTM performs well on the 3D-related questions, such as counting and relation, than most of the image-based baselines. Compared with 3D-CLR, NS-VQA can perform well in the conceptual questions. However, it underperforms 3D-CLR a lot in counting and relation, suggesting that these two types of questions require the holistic 3D understanding of the entire 3D scenes. Our 3D-CLR outperforms other baselines by a large margin, but is still far from satisfying. From the accuracy of the conceptual question, we can see that it can only ground approximately 66% of the semantic concepts. This indicates that our 3DMV-VQA dataset is indeed very challenging.

**Qualitative Examples.** In Fig. 3, we show four qualitative examples. From the examples, we show that our 3D-CLR can infer an accurate 3D representation from multi-view images, as well as ground semantic concepts on the 3D representations to get the semantic segmentations of the entire scene. Our 3D-CLR can also learn 3D relationships such as “close”, “largest”, “on top of” and so on. However, 3D-CLR also fails on some questions. For the third scene in the qualitative examples, it fails to ground the concepts “mouse” and “printer”. Also, it cannot accurately count the instances sometimes. We give detailed discussions below.

### 5.4. Discussions

We perform an in-depth analysis to understand the challenge of this dataset. We leverage the modular design of our 3D-CLR, replacing individual components of the framework with ground-truth annotations for model diagnosis. The result is shown in Fig 4. 3D-CLR w/ Semantic denotes our model with ground-truth semantic concepts from HM3DSem annotations. 3D-CLR w/ Instance denotes that we have ground-truth instance segmentations of semantic

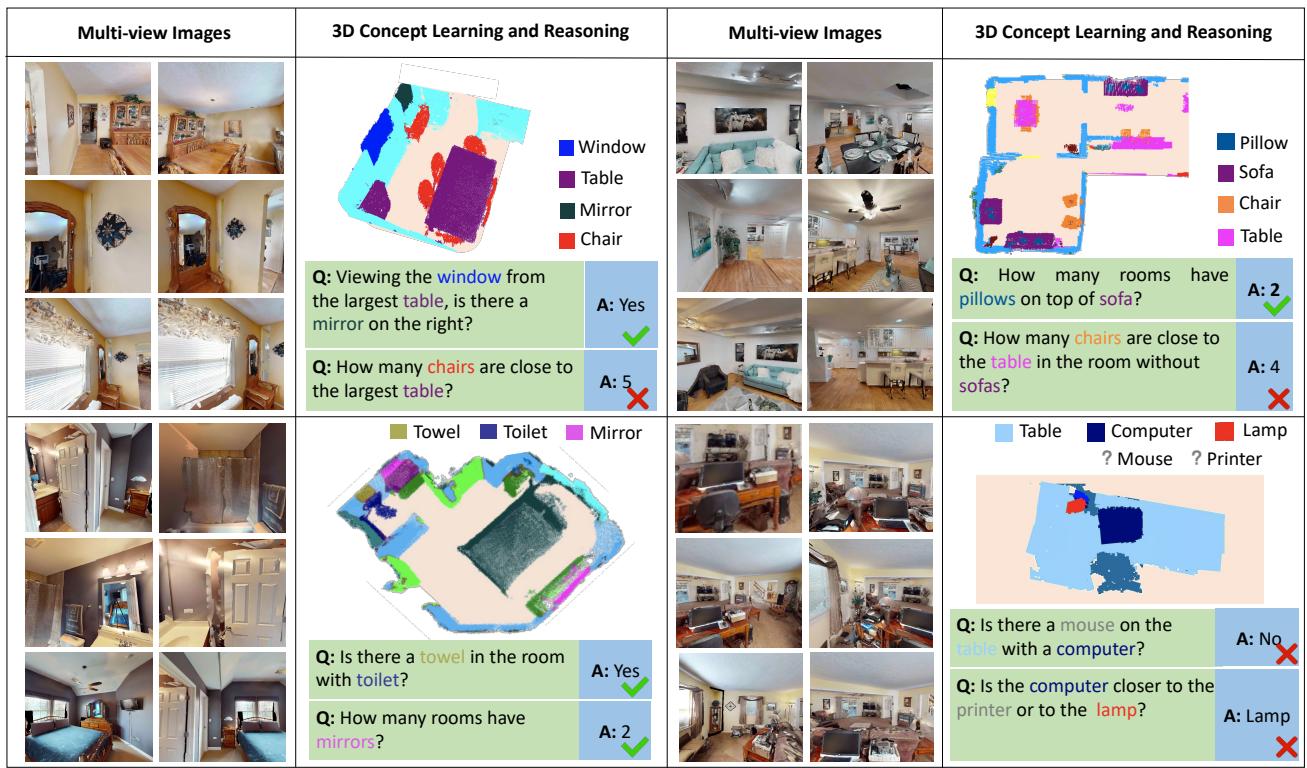


Figure 3. Qualitative examples of our 3D-CLR. We can see that 3D-CLR can ground most of the concepts and answer most questions correctly. However, it still fails sometimes, mainly because it cannot separate close object instances and ground small objects.

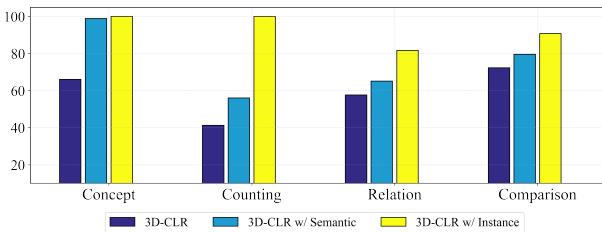


Figure 4. Model diagnosis of our 3D-CLR.

concepts. From Fig. 3 and Fig. 4, we summarize several key challenges of our benchmark:

**Very close object instances** From Fig. 4, we can see that even with ground-truth semantic labeling of the 3D points, 3D-CLR still has unsatisfying results on counting questions. This suggests that the instance segmentations provided by DBSCAN are not accurate enough. From the top two qualitative examples in Fig. 3, we can also see that if two chairs contact each other, DBSCAN will not tell them apart and thus have poor performance on counting. One crucial future direction is to improve unsupervised instance segmentations on very close object instances.

**Grounding small objects** Fig. 4 suggests that 3D-CLR fails to ground a large portion of the semantic concepts, which hinders the performance. From the last example in Fig. 3, we can see that 3D-CLR fails to ground small objects like “computer mouse”. Further examination indicates there are

two possible reasons: 1) CLIP-LSeg fails to assign the right features to objects with limited pixels; 2) The resolution of feature voxel grid is not high enough and therefore small objects cannot be represented in the compact representation. An interesting future direction would be learning exploration policies that enable the agents to get closer to uncertain objects that cannot be grounded.

**Ambiguity on 3D relations** Even with ground-truth semantic and instance segmentations, the performance of the relation network still needs to be improved. We find that most of the failure cases are correlated to the “inside” relation. From the segmentations in Fig. 3, we can see that 3D-CLR is unable to ground the objects in the cabinets. A potential solution can be joint depth and segmentation predictions.

## 6. Conclusion

In this paper, we introduce the novel task of 3D reasoning from multi-view images. By placing embodied robot that actively explores indoor environments, we collect a large-scale benchmark named 3DMV-VQA. We also propose a new 3D-CLR model that incorporates neural field, 2D VLM, as well as reasoning operators for this task and illustrate its effectiveness. Finally, we perform an in-depth analysis to understand the challenges of this dataset and also point out potential future directions. We hope that 3DMV-VQA can be used to push the frontiers of 3D reasoning.

864

## References

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

- [1] Panos Achlioptas, Ahmed Abdelreheem, Fei Xia, Mohamed Elhoseiny, and Leonidas J. Guibas. Referit3d: Neural listeners for fine-grained 3d object identification in real-world scenes. In *ECCV*, 2020. [2](#), [3](#)
- [2] Daich Azuma, Taiki Miyanishi, Shuhei Kurita, and Motoki Kawanabe. Scanqa: 3d question answering for spatial scene understanding. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 19107–19117, 2022. [2](#), [3](#)
- [3] Dave Zhenyu Chen, Angel X. Chang, and Matthias Nießner. Scanrefer: 3d object localization in rgb-d scans using natural language. In *ECCV*, 2020. [2](#), [3](#)
- [4] Dave Zhenyu Chen, Ali Gholami, Matthias Nießner, and Angel X. Chang. Scan2cap: Context-aware dense captioning in rgb-d scans. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3192–3202, 2021. [3](#)
- [5] Z Chen, L Ma, W Luo, and KKY Wong. Weakly-supervised spatio-temporally grounding natural sentence in video. In *ACL*, 2019. [3](#)
- [6] Zhenfang Chen, Jiayuan Mao, Jiajun Wu, Kwan-Yee Kenneth Wong, Joshua B Tenenbaum, and Chuang Gan. Grounding physical concepts of objects and events through dynamic visual reasoning. *ICLR*, 2021. [3](#)
- [7] Zhenfang Chen, Peng Wang, Lin Ma, Kwan-Yee K Wong, and Qi Wu. Cops-ref: A new dataset and task on compositional referring expression comprehension. In *CVPR*, 2020. [3](#)
- [8] Zhenfang Chen, Kexin Yi, Yunzhu Li, Mingyu Ding, Antonio Torralba, Joshua B Tenenbaum, and Chuang Gan. Comphy: Compositional physical reasoning of objects and events from videos. In *ICLR*, 2022. [3](#)
- [9] Zhenfang Chen, Qinhong Zhou, Yikang Shen, Yining Hong, Hao Zhang, and Chuang Gan. See, think, confirm: Interactive prompting between vision and language models for knowledge-based visual reasoning. *arXiv preprint arXiv:2301.05226*, 2023. [3](#)
- [10] Spconv Contributors. Spconv: Spatially sparse convolution library. <https://github.com/traveller59/spconv>, 2022. [6](#)
- [11] Abhishek Das, Samyak Datta, Georgia Gkioxari, Stefan Lee, Devi Parikh, and Dhruv Batra. Embodied question answering. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 2135–213509, 2018. [3](#)
- [12] Mingyu Ding, Yan Xu, Zhenfang Chen, David Daniel Cox, Ping Luo, Joshua B Tenenbaum, and Chuang Gan. Embodied concept learner: Self-supervised learning of concepts and mapping through instruction following. In *CoRL*. [3](#)
- [13] Yilun Du, M. Katherine Collins, B. Joshua Tenenbaum, and Vincent Sitzmann. Learning signal-agnostic manifolds of neural fields. In *Advances in Neural Information Processing Systems*, 2021. [3](#)
- [14] Yilun Du, Yinan Zhang, Hong-Xing Yu, Joshua B. Tenenbaum, and Jiajun Wu. Neural radiance flow for 4d view synthesis and video processing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021. [3](#)
- [15] Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. In *KDD*, 1996. [6](#)
- [16] Yasaman Etesam, Leon Kochiev, and Angel X Chang. 3dvqa: Visual question answering for 3d environments. In *2022 19th Conference on Robots and Vision (CRV)*, pages 233–240. IEEE, 2022. [2](#), [3](#)
- [17] Mingtao Feng, Zhen Li, Qi Li, Liang Zhang, Xiangdong Zhang, Guangming Zhu, Hui Zhang, Yaonan Wang, and Ajmal S. Mian. Free-form description guided 3d visual graph network for object grounding in point cloud. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3702–3711, 2021. [3](#)
- [18] Chuang Gan, Yandong Li, Haoxiang Li, Chen Sun, and Boqing Gong. Vqs: Linking segmentations to questions and answers for supervised attention in vqa and question-focused semantic segmentation. In *ICCV*, pages 1811–1820, 2017. [3](#)
- [19] Siddha Ganju, Olga Russakovsky, and Abhinav Kumar Gupta. What's in a question: Using visual questions as a form of supervision. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6422–6431, 2017. [3](#)
- [20] Xiaofeng Gao, Qiaozi Gao, Ran Gong, Kaixiang Lin, Govind Thattai, and Gaurav S Sukhatme. Dialfred: Dialogue-enabled agents for embodied instruction following. *arXiv*, 2022. [3](#)
- [21] Stephan J. Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien P. C. Valentin. Fastnerf: High-fidelity neural rendering at 200fps. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 14326–14335, 2021. [3](#), [5](#)
- [22] Daniel Gordon, Aniruddha Kembhavi, Mohammad Rastegari, Joseph Redmon, Dieter Fox, and Ali Farhadi. Iqa: Visual question answering in interactive environments. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4089–4098, 2018. [3](#)
- [23] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6325–6334, 2017. [2](#), [3](#)
- [24] Peter Hedman, Pratul P. Srinivasan, Ben Mildenhall, Jonathan T. Barron, and Paul E. Debevec. Baking neural radiance fields for real-time view synthesis. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5855–5864, 2021. [3](#), [5](#)
- [25] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9:1735–1780, 1997. [6](#)
- [26] Yining Hong, Yilun Du, Chunru Lin, Joshua B Tenenbaum, and Chuang Gan. 3d concept grounding on neural fields. *arXiv preprint arXiv:2207.06403*, 2022. [3](#), [4](#), [5](#), [6](#)
- [27] Yining Hong, Li Yi, Joshua B. Tenenbaum, Antonio Torralba, and Chuang Gan. Ptr: A benchmark for part-based conceptual, relational, and physical reasoning. In *NeurIPS*, 2021. [3](#)
- [28] Deng Huang, Peihao Chen, Runhao Zeng, Qing Du, Mingkui Tan, and Chuang Gan. Location-aware graph convolutional networks for video question answering. In *AAAI*, 2020. [3](#), [7](#)
- [29] Pin-Hao Huang, Han-Hung Lee, Hwann-Tzong Chen, and Tyng-Luh Liu. Text-guided graph neural networks for referring 3d instance segmentation. In *AAAI*, 2021. [3](#)

- 972 [30] D. A. Hudson and Christopher D. Manning. Compositional  
973 attention networks for machine reasoning. *ICLR*, 2018. 3, 6  
974  
975 [31] Ajay Jain, Ben Mildenhall, Jonathan T. Barron, P. Abbeel,  
976 and Ben Poole. Zero-shot text-guided object generation with  
977 dream fields. *2022 IEEE/CVF Conference on Computer Vi-  
978 sion and Pattern Recognition (CVPR)*, pages 857–866, 2022.  
979 3  
980  
981 [32] Zhenyu Jiang, Yifeng Zhu, Maxwell Svetlik, Kuan Fang,  
982 and Yuke Zhu. Synergies between affordance and geometry:  
983 6-dof grasp detection via implicit representations. *ArXiv*,  
984 abs/2104.01542, 2021. 3  
985  
986 [33] J. Johnson, Bharath Hariharan, L. V. D. Maaten, Li Fei-Fei,  
987 C. L. Zitnick, and Ross B. Girshick. Clevr: A diagnostic  
988 dataset for compositional language and elementary visual  
989 reasoning. *2017 IEEE Conference on Computer Vision and  
990 Pattern Recognition (CVPR)*, pages 1988–1997, 2017. 2, 3, 4  
991  
992 [34] Stefan Kollmannsberger, Davide D’Angella, Moritz Jokeit,  
993 and Leon Alexander Herrmann. Physics-informed neural  
994 networks. *Deep Learning in Computational Mechanics*, 2021.  
995 3  
996  
997 [35] Natalia Konstantinova and Constantin Orasan. Interactive  
998 question answering. In *EMNLP*. IGI Global, 2013. 3  
999  
1000 [36] Barbara Landau and Ray Jackendoff. “what” and “where” in  
1001 spatial language and spatial cognition. *Behavioral and Brain  
1002 Sciences*, 16:217–238, 1993. 6  
1003  
1004 [37] Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen  
1005 Koltun, and René Ranftl. Language-driven semantic seg-  
1006 mentation. *ICLR*, 2022. 2, 4, 5  
1007  
1008 [38] Dongxu Li, Junnan Li, Hongdong Li, Juan Carlos Niebles,  
1009 and Steven C. H. Hoi. Align and prompt: Video-and-language  
1010 pre-training with entity prompts. *2022 IEEE/CVF Conference  
1011 on Computer Vision and Pattern Recognition (CVPR)*, pages  
1012 4943–4953, 2022. 3, 7  
1013  
1014 [39] Andrew Luo, Yilun Du, Michael J Tarr, Joshua B Tenenbaum,  
1015 Antonio Torralba, and Chuang Gan. Learning neural acoustic  
1016 fields. *arXiv preprint arXiv:2204.00628*, 2022. 3  
1017  
1018 [40] Jiayuan Mao, Chuang Gan, Pushmeet Kohli, Joshua B. Tenen-  
1019 baum, and Jiajun Wu. The neuro-symbolic concept learner:  
1020 Interpreting scenes words and sentences from natural supervi-  
1021 sion. *ArXiv*, abs/1904.12584, 2019. 3, 6  
1022  
1023 [41] Nelson L. Max. Optical models for direct volume rendering.  
1024 *IEEE Trans. Vis. Comput. Graph.*, 1:99–108, 1995. 5  
1025  
1026 [42] Lars M. Mescheder, Michael Oechsle, Michael Niemeyer,  
1027 Sebastian Nowozin, and Andreas Geiger. Occupancy net-  
1028 works: Learning 3d reconstruction in function space. *2019  
1029 IEEE/CVF Conference on Computer Vision and Pattern Recog-  
1030 nition (CVPR)*, pages 4455–4465, 2019. 3  
1031  
1032 [43] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik,  
1033 Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf:  
1034 Representing scenes as neural radiance fields for view syn-  
1035 thesis. In *Proc. ECCV*, 2020. 3, 5  
1036  
1037 [44] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and  
1038 Andreas Geiger. Occupancy flow: 4d reconstruction by learn-  
1039 ing particle dynamics. In *Proceedings of the IEEE Interna-  
1040 tional Conference on Computer Vision*, pages 5379–5389,  
1041 2019. 3  
1042  
1043 [45] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and  
1044 Andreas Geiger. Differentiable volumetric rendering: Learn-  
1045 ing implicit 3d representations without 3d supervision. In  
1046 *Proc. CVPR*, 2020. 3  
1047  
1048 [46] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya  
1049 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,  
1050 Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen  
1051 Krueger, and Ilya Sutskever. Learning transferable visual  
1052 models from natural language supervision. In *ICML*, 2021. 5  
1053  
1054 [47] Santhosh K. Ramakrishnan, Aaron Gokaslan, Erik Wijmans,  
1055 Oleksandr Maksymets, Alexander Clegg, John Turner, Eric  
1056 Undersander, Wojciech Galuba, Andrew Westbury, Angel X.  
1057 Chang, Manolis Savva, Yili Zhao, and Dhruv Batra. Habitat-  
1058 matterport 3d dataset (hm3d): 1000 large-scale 3d environ-  
1059 ments for embodied ai. *ArXiv*, abs/2109.08238, 2021. 2,  
1060 3  
1061  
1062 [48] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray,  
1063 Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever.  
1064 Zero-shot text-to-image generation. *ArXiv*, abs/2102.12092,  
1065 2021. 5  
1066  
1067 [49] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Mor-  
1068 ishimura, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned  
1069 implicit function for high-resolution clothed human digitiza-  
1070 tion. In *Proc. ICCV*, pages 2304–2314, 2019. 3  
1071  
1072 [50] Manolis Savva, Abhishek Kadian, Oleksandr Maksymets,  
1073 Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia  
1074 Liu, Vladlen Koltun, Jitendra Malik, Devi Parikh, and Dhruv  
1075 Batra. Habitat: A Platform for Embodied AI Research. In  
1076 *Proceedings of the IEEE/CVF International Conference on  
1077 Computer Vision (ICCV)*, 2019. 3  
1078  
1079 [51] Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan  
1080 Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer, and  
1081 Dieter Fox. Alfred: A benchmark for interpreting grounded  
1082 instructions for everyday tasks. In *CVPR*, 2020. 3  
1083  
1084 [52] Anthony Simeonov, Yilun Du, Andrea Tagliasacchi, Joshua B  
1085 Tenenbaum, Alberto Rodriguez, Pulkit Agrawal, and Vin-  
1086 cent Sitzmann. Neural descriptor fields: Se (3)-equivariant  
1087 object representations for manipulation. *arXiv preprint  
1088 arXiv:2112.05124*, 2021. 3  
1089  
1090 [53] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein.  
1091 Scene representation networks: Continuous 3d-structure-  
1092 aware neural scene representations. In *Proc. NeurIPS 2019*,  
1093 2019. 3  
1094  
1095 [54] Chan Hee Song, Jihyung Kil, Tai-Yu Pan, Brian M Sadler,  
1096 Wei-Lun Chao, and Yu Su. One step at a time: Long-horizon  
1097 vision-and-language navigation with milestones. *arXiv  
1098 preprint arXiv:2202.07028*, 2022. 3  
1099  
1100 [55] Elizabeth S Spelke, Karen Breinlinger, Kristen Jacobson, and  
1101 Ann Phillips. Gestalt Relations and Object Perception: A  
1102 Developmental Study. *Perception*, 22(12):1483–1501, 1993.  
1103 2  
1104  
1105 [56] Alessandro Suglia, Qiaozi Gao, Jesse Thomason, Govind  
1106 Thattai, and Gaurav Sukhatme. Embodied bert: A transformer  
1107 model for embodied, language-guided visual task completion.  
1108 *arXiv preprint arXiv:2108.04927*, 2021. 3  
1109  
1110 [57] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel  
1111 grid optimization: Super-fast convergence for radiance fields  
1112

- 1080 reconstruction. 2022 IEEE/CVF Conference on Computer 1134  
1081 Vision and Pattern Recognition (CVPR), pages 5449–5459, 1135  
1082 2022. 2, 3, 5 1136
- 1083 [58] Andrew Szot, Alex Clegg, Eric Undersander, Erik Wijmans, 1137  
1084 Yili Zhao, John Turner, Noah Maestre, Mustafa Mukadam, 1138  
1085 Devendra Chaplot, Oleksandr Maksymets, Aaron Gokaslan, 1139  
1086 Vladimir Vondrus, Sameer Dharur, Franziska Meier, Wojciech 1140  
1087 Galuba, Angel Chang, Zsolt Kira, Vladlen Koltun, Jitendra 1141  
1088 Malik, Manolis Savva, and Dhruv Batra. Habitat 2.0: Training 1142  
1089 home assistants to rearrange their habitat. In *Advances in 1143  
1090 Neural Information Processing Systems (NeurIPS)*, 2021. 3 1144
- 1091 [59] Ivan Vendrov, Ryan Kiros, Sanja Fidler, and Raquel Urtasun. 1145  
1092 Order-embeddings of images and language. *CoRR*, 1146  
1093 abs/1511.06361, 2016. 3 1147
- 1094 [60] Can Wang, Menglei Chai, Mingming He, Dongdong Chen, 1148  
1095 and Jing Liao. Clip-nerf: Text-and-image driven manipulation 1149  
1096 of neural radiance fields. *ArXiv*, abs/2112.05139, 2021. 3 1150
- 1097 [61] Karmesh Yadav, Ram Ramrakhyta, Santhosh Kumar Ramakrishnan, 1151  
1098 Theo Gervet, John Turner, Aaron Gokaslan, Noah Maestre, 1152  
1099 Angel Xuan Chang, Dhruv Batra, Manolis Savva, et al. 1153  
1100 Habitat-matterport 3d semantics dataset. *arXiv preprint arXiv:2210.05633*, 2022. 2, 3, 4 1154
- 1101 [62] Xu Yan, Zhihao Yuan, Yuhao Du, Yinghong Liao, Yao Guo, 1155  
1102 Zhen Li, and Shuguang Cui. Clevr3d: Compositional language 1156  
1103 and elementary visual reasoning for question answering 1157  
1104 in 3d real-world scenes. *arXiv preprint arXiv:2112.11691*, 1158  
1105 2021. 2, 3 1159
- 1106 [63] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan 1160  
1107 Atzmon, Basri Ronen, and Yaron Lipman. Multiview neural 1161  
1108 surface reconstruction by disentangling geometry and appearance. 1162  
1109 *Proc. NeurIPS*, 2020. 3 1163
- 1110 [64] Shuquan Ye, Dongdong Chen, Songfang Han, and Jing Liao. 1164  
1111 3d question answering. *ArXiv*, abs/2112.08359, 2021. 2, 3 1165
- 1112 [65] Kexin Yi, Jiajun Wu, Chuang Gan, Antonio Torralba, Pushmeet 1166  
1113 Kohli, and Joshua B. Tenenbaum. Neural-symbolic 1167  
1114 vqa: Disentangling reasoning from vision and language understanding. 1168  
1115 In *NeurIPS*, 2018. 3, 6, 7 1169
- 1116 [66] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and 1170  
1117 Angjoo Kanazawa. Plenoctrees for real-time rendering of 1171  
1118 neural radiance fields. 2021 IEEE/CVF International Conference 1172  
1119 on Computer Vision (ICCV), pages 5732–5741, 2021. 3, 5 1173
- 1120 [67] Licheng Yu, Xinlei Chen, Georgia Gkioxari, Mohit Bansal, 1174  
1121 Tamara L Berg, and Dhruv Batra. Multi-target embodied 1175  
1122 question answering. In *CVPR*, pages 6309–6318, 2019. 3 1176
- 1123 [68] Kaizhi Zheng, Xiaotong Chen, Odest Chadwicke Jenkins, 1177  
1124 and Xin Eric Wang. Vlmbench: A compositional benchmark 1178  
1125 for vision-and-language manipulation. In *Proceedings of the 1179  
1126 Neural Information Processing Systems Track on Datasets and 1180  
1127 Benchmarks*, 2022. 3 1181
- 1128 [69] Yuke Zhu, O. Groth, Michael S. Bernstein, and Li Fei-Fei. 1182  
1129 Visual7w: Grounded question answering in images. 2016 1183  
1130 IEEE Conference on Computer Vision and Pattern Recognition 1184  
1131 (CVPR), pages 4995–5004, 2016. 2, 3 1185
- 1132 1186
- 1133 1187