# **3D Visual Grounding with Reasoning LLM**

Anonymous submission

Paper ID



Figure 1. **Overview.** We propose a data generation pipeline that can generate visual grounding queries and reasoning responses. We use our collected data to conduct post-training, resulting in an LLM that can conduct structured reasoning in 3D visual grounding tasks.

### Abstract

001 The recent development of Large Language Models (LLMs) with strong reasoning ability has driven research in various 002 domains such as mathematics, coding, and scientific rea-003 004 soning. Meanwhile, 3D visual grounding, as an important task in the 3D understanding area, still remains challeng-005 ing due to the limited reasoning ability of recent 3D visual 006 grounding models. Most of the current methods incorpo-007 rate a text encoder and visual feature encoder to generate 008 cross-modal fuse features and predict the referring object. 009 010 These models often require supervised training on extensive

3D annotation data. Meanwhile, recent zero-shot visual 011 grounding methods require access to proprietary LLM dur-012 ing test time, which leads to large inference costs. To over-013 come these limitations, we propose a 3D visual grounding 014 data generation pipeline capable of synthesizing 3D scenes 015 along with corresponding target queries and ground truth 016 answers for training visual grounding models. Additionally, 017 we leverage the generated visual grounding data and con-018 duct post-training on Llama-3.1-8B-Instruct, resulting in a 019 strong 3D visual grounding LLM that outperforms existing 020 SoTA zero-shot methods, demonstrating its effectiveness. 021

### **1. Introduction**

3D Visual Grounding (3DVG) is a crucial task in 3D un-023 derstanding, aiming to identify a target object within a 024 3D scene based on a given textual query. Recent mod-025 els [19, 26, 31, 32, 34, 41] typically integrate a text en-026 027 coder and a visual encoder to generate cross-modal features for target object prediction. While these methods 028 029 achieve strong performance, they rely on large-scale annotated 3D visual grounding datasets for supervised train-030 ing. To mitigate this, recent research has explored zero-031 shot approaches, leveraging proprietary LLMs like GPT-032 4 [3] as agents to generate grounding predictions using vi-033 034 sual grounding tools [36] or in-context examples and programs [40]. However, these approaches heavily depend on 035 proprietary LLMs, increasing inference costs. In this work, 036 we aim to explore using an 8B LLM to solve the challeng-037 ing visual grounding task under zero-shot setting, by lever-038 aging LLM's reasoning ability, we hope to provide a solu-039 040 tion with strong visual grounding performance with existing zero-shot methods while being more cost-effective. 041

Recently, LLMs with reasoning ability [12, 18, 27, 35] 042 have achieved great success in various domains including 043 mathematics, coding, and scientific reasoning. DeepSeek-044 R1 [12] leverages reinforcement learning to achieve strong 045 thinking ability and SoTA performance on various challeng-046 ing benchmarks. However, despite the success of these rea-047 048 soning models, they mainly focus on the text and image understanding domain, while solving the 3D visual grounding 049 or understanding tasks by leveraging LLM's strong reason-050 ing ability remains under-explored, with most of the current 051 SoTA 3D visual grounding methods [19, 31] incorporates a 052 053 text encoder to encode the given query and conduct direct prediction based on the cross-modal fused feature from var-054 ious input modalities. In the 3D visual grounding task, the 055 given query can contain sophisticated spatial relationships, 056 057 and leveraging LLM to comprehend the input query is another potential solution that is not fully explored. 058

059 In this work, we aim to solve the 3D visual grounding task using a reasoning approach, we present an overview 060 of our work in Fig. 1. Given the query of the target object 061 in a 3D scene, we developed an LLM that can conduct a 062 rigorous reasoning process before making the final predic-063 064 tion. To train such LLM, we design a 3D visual grounding reasoning data collection pipeline, that aims to address 065 the current methods' dependency on 3D visual grounding 066 data. We use our collected data and conduct post-training 067 on the base model Llama-3.1-8B-Instruct [8], resulting in 068 069 a 3D visual grounding LLM with strong reasoning abil-070 ity, which achieves higher performance than existing SoTA zero-shot method on multiple 3D visual grounding bench-071 072 marks. Moreover, we discover the spatial hallucination phenomenon of proprietary LLM on spatial understanding 073 074 tasks during our data collection process, we further design

a prompt that can boost the accuracy of the response from LLM by a large margin in our data collection pipeline. We summarize our contribution as follows:

- We propose a novel 3D visual grounding reasoning data collection pipeline, which serves as a post-training dataset for our reasoning LLM on zero-shot 3D visual grounding.
   080
- We discover the spatial hallucination of current SoTA LLM and design an effective prompt to reduce the spatial hallucination, enabling more accurate reasoning and response during data collection.
   081 082 083 084
- We conduct post-training on our collected data and introduce **3DVG-LLM**, an LLM for 3D visual grounding task that achieves SoTA performance on 3D visual grounding benchmarks including NR3D and SR3D.
- To our best knowledge, this is the first successful attempt using synthetic 3D data to train a 3D understanding LLM
   with reasoning ability, we hope this work serves as a cornerstone for future 3D-LLM development with stronger reasoning and spatial understanding ability.
   089
   090
   091
   092
   093

#### 2. Related Work

#### 2.1. 3D Visual Grounding

3D visual grounding is an important 3D understanding task 096 that aims to predict the target object in the 3D scene based 097 on a given language query of the target. Most of the re-098 cent 3D visual grounding methods [16, 19, 21, 34] adopt 099 separate encoders for each input modality such as text and 100 point cloud to generate cross-modal features and conduct 101 prediction. Some end-to-end methods like SAT [38], Lan-102 guageRefer [29], and UniT3D [6] utilized a unified multi-103 modal transformer to conduct 3D visual grounding. Re-104 cently, COT3DREF [1] and ReGround3D [42] proposed to 105 leverage chain-of-thoughts to retrieve objects related to a 106 given query, however, they did not leverage the language 107 models to conduct reasoning for better spatial understand-108 ing. Furthermore, these methods utilize supervised train-109 ing, which requires extensive annotations on the 3D data. 110 Besides the supervised method, recent works also focus on 111 the zero-shot approaches [36, 40]. These methods lever-112 age LLM to generate code and solve the 3D visual ground-113 ing task, which achieves some success in zero-shot 3D vi-114 sual grounding. However, these methods rely on in-context 115 examples and proprietary LLMs [3, 24] during test time, 116 which introduce more inference costs. In this work, we ex-117 plore using an 8B open-source LLM to solve the 3D vi-118 sual grounding task through inference time reasoning un-119 der zero-shot setting. Our method does not require pro-120 prietary LLM nor 3D data that requires extensive annota-121 tion, while achieving better performance compared with the 122 SoTA zero-shot methods, showcasing the potential of open-123 source LLM for zero-shot 3D visual grounding, and paving 124 the way for a more cost-effective solution in the field. 125

075

076

077

094



#

149

150

156

Figure 2. An illustration of our 3D scene layout generation pipeline, which contains 5 steps: 1) Set up an empty 3D scene with certain dimension. 2) Choose a spatial relationship and decide the anchor and distractor objects as well as their dimension and location. 3) Place them in the 3D scene. 4) Determine the target object from the distractors and parse query. 5) Generate more random objects to enrich scene diversity. Note that some spatial relationships do not have anchor object and only have distractors, such as Largest/Smallest. Note that the furniture in the figure are for visualization purpose only, our data pipeline does not include generating the appearance of the furniture. Please refer to our attached scene generation code for more technical detail on the implementations.

#### 126

#### 2.2. LLM with Reasoning Ability

127 Recent advancements in LLMs have significantly enhanced their reasoning capabilities, enabling them to tackle com-128 plex problems across various domains such as mathematics, 129 coding, and scientific reasoning. Models like DeepSeek-130 131 R1 [12] incorporate reinforcement learning during train-132 ing and achieve SoTA performance across multiple bench-133 marks. Similarly, LLaVA-CoT [35] and VisualCoT [30] introduce explicit reasoning processes to improve VLM on 134 visual question-answering task. Prism [27] proposed to de-135 couple the perception and reasoning stages to enhance VLM 136 performance. While these models have demonstrated re-137 138 markable success in text and image understanding, leveraging reasoning LLMs for zero-shot 3D understanding re-139 mains largely unexplored, with most existing 3D visual 140 grounding methods still require supervised training from 141 142 the extensive annotated 3D data. To address this gap, we present a pioneering work that integrates reasoning LLM 143 to tackle the challenging 3D visual grounding task. We 144 hope this work serves as a strong baseline for reasoning-145 enabled LLMs in the 3D scene understanding domain and 146 147 inspires further research at the intersection of reasoning-148 driven LLMs and 3D scene comprehension.

#### 3. 3D Visual Grounding Data Collection

3.1. Overview

In this work, we design a novel 3DVG data collection 151 framework for 3D reasoning visual grounding LLM train-152 ing. In this section, we introduce our program-based 3D 153 scene generation pipeline in Sec. 3.2, reasoning stages in 154 Sec. 3.3, and the reasoning data collection in Sec. 3.4. 155

#### 3.2. 3D Scene Generation

The collection of 3D visual grounding data requires a huge 157 amount of effort including extensive human annotations and 158 indoor 3D RGB-D scan collection. Unlike previous super-159 vised methods that leverage human-annotated 3D data, we 160 propose to generate synthetic 3D scene data for model train-161 ing and zero-shot 3D visual grounding on real-world 3D 162 scenes. There are many existing methods for 3D scene gen-163 eration [10, 22, 37], however, they mostly focus on gen-164 erating scene layout with realistic object arrangements or 165 following the specified room style from users' generation 166 prompt, while our goal is to collect object-centric 3D scene 167 that focuses on objects' spatial arrangements as well as the 168 target object's corresponding grounding query, which is an 169 important component of 3D visual grounding data. 170

171 In this work, we design a program-based 3D scene layout generation pipeline that can generate an object-centric 172 3D scene following several common spatial relationships 173 between objects. These common spatial relationships in-174 175 clude: closest/farthest, largest/smallest, left/right, and beside. These common spatial relationships and correspond-176 ing scene layouts can be generated by our code program. 177 Here, we provide a generation example of a 3D scene lay-178 179 out in Fig. 2, which contains the 3D scene layout following query "The Chair that is closest to the Table." with n dis-180 181 tractor chairs in the scene:

- 182 1. Define a 3D scene size with a certain dimension  $H \times W$ .
- 183 2. Randomly generate n chairs and one table in the 3D scene.
- 3. Calculate the Euclidean distance between each chair and table to find the closest chair to the table; this chair will be the target object of the 3D visual grounding query.
- 4. Generate more objects (but not chairs or tables) to furtherenrich the 3D scene diversity.

Following this procedure, we collect a wide range of 3D 190 scene layouts that follow our specified common spatial re-191 lationship and are not confined by the limited number of 192 3D scene layouts in the real-world available 3D scans. For 193 194 each spatial relationship, we design multiple templates from which we can parse the target's object class to generate the 195 final target object query. We prompt GPT-40 [17] to gen-196 erate some indoor scene objects as well as their common 197 dimension. When placing objects, we randomly select ob-198 199 ject classes and apply a random variation on their size to 200 enhance diversity. We attached our 3D scene layout generation code in the appendix for reference. 201

202 Note that because our base model LLM already has very 203 strong reasoning ability, we do not focus on designing complicated spatial relationships to strengthen its reasoning and 204 spatial understanding ability, but only use several common 205 relationships with reasoning data to encourage the LLM to 206 think step-by-step. Our experiment results demonstrate that 207 our LLM can still adapt to a 3D visual grounding query that 208 has an unseen spatial relationship after post-training on our 209 210 collected data. Furthermore, because our goal is to generate 3D scene data that follow our specified spatial relationship 211 between objects, during the random object placements, we 212 only refrain the object from being placed overlapping and 213 214 do not focus on placing the object in a way that is realistic in real-world scenarios. Given our goal is to guide LLM's 215 reasoning process, we believe the rationality of object ar-216 217 rangement will not affect LLM's reasoning process.

## **218 3.3. 3D Visual Grounding with Reasoning Stages**

Recent advanced LLMs with reasoning ability adopt special
tokens or tags [12, 35] to separate the reasoning and conclusion stage, which enables a more structured and organized
thinking process of the LLM. In the 3D visual grounding



Figure 3. An overview of the reasoning data collection pipeline, after the grounding data generator produce the 3D scene and query, we use GPT-4 $\circ$  to generate structured reasoning response. Only the response with correct answer and format will be stored in our 3D visual grounding and reasoning data for later training.

domain, we believe it is also important to think step-by-step223before generating the final prediction given the difficulty224of interpreting the spatial relationship of the target query.225Here, we proposed four different reasoning stages, which226encourage the LLM to generate a more structured thinking227process. The four reasoning stages include:228

- Related Object Selection. Given all the objects in the 3D scene, the LLM listed the objects that might be relevant to the given query. Each listed object will include information such as its object ID, class name, location, and dimension.
- Situation Estimation. If a 3D situation (viewer's location) is provided, the LLM will estimate its possible coordinate in the 3D scene. Otherwise, the LLM will just assume it is standing in the middle of the 3D scene (which is usually the case).
- **Reasoning.** This is the main stage for LLM to conduct reasoning based on the previously given information. The LLM will conduct necessary math calculations and logical reasoning to derive the final answer.
- **Conclusion.** In this stage, the model outputs the final prediction in a structured format, which enables easier parsing and understanding for the user.

234

235

236

237

238

239

240

241

242

243

244



Figure 4. A comparison of  $GPT-4\circ$  response with and without extra prompt. We highlight extra prompt in yellow. The correct reasoning from  $GPT-4\circ$  in green, and the error reasoning in red. The accuracy gain with extra prompt can be found in ablation study.

#### **3.4. Reasoning Data Collection**

247 To collect detailed reasoning data aligned with our four reasoning stages, we leverage the SoTA proprietary LLM 248 GPT-4o [17] to generate structured reasoning responses, 249 250 as illustrated in Fig. 3. We design a specialized prompt to 251 guide GPT-40 in adhering to our four-stage reasoning pro-252 cess. The full prompt is provided in the appendix. Approx-253 imately 10% of responses that deviate from the specified format or contain incorrect predictions are automatically 254 filtered. The remaining detailed visual grounding reason-255 256 ing responses are incorporated into our final training data, which consists of 3.2K samples. Each sample includes a 257 unique 3D scene layout with over 50 objects, a target ob-258 259 ject's language query, its corresponding object ID, and the reasoning response from GPT-40. Note that since the pri-260 mary goal of post-training is to enhance the LLM's reason-261 262 ing capabilities, we only utilize a small-scale dataset for data efficient post-training, and leave the scale-up training 263 using our data collection framework as a future direction. 264

#### **265 3.5. Reducing The Spatial Hallucination**

266 During the reasoning data collection, we found that despite 267 the strong reasoning ability of the SoTA language model like GPT-40, it still suffers from spatial hallucination and 268 generates wrong predictions sometimes, especially when 269 the given query is related to the left/right spatial relation-270 ship. This is likely caused by LLM's weak spatial under-271 standing ability, leading to errors in recognizing left/right 272 direction. For example, we found that GPT-40 often erro-273 neously thinks "when facing the negative Y direction, the 274 *left-hand side is negative X.*" This leads to a lower accuracy 275 when predicting target object related to left/right spatial re-276 lationship. To address this issue and ensure a more effi-277 cient data collection process, we modify the input prompt 278 and encourage GPT-40 to determine the left/right object 279 relationship by calculating the cross product between the 280 location coordinate vectors toward two objects, in other 281 words, we prompt GPT-40 to calculate the cross product 282 between viewer's coordinate (which is estimated in the Situ-283 ation Estimation Stage) and two objects' coordinate. Based 284 on the cross product's value, LLM can determine two ob-285 jects' left/right spatial relationship. We make a comparison 286 of the GPT-40 response with and without our extra prompt 287 in Fig. 4, and also conduct analysis in our ablation study and 288 found out the extra prompt largely boosted the response ac-289 curacy, suggesting the effectiveness of our designed prompt. 290

#

328

329

330

331

332

333

334

335

336

				Unique		Multiple		Overall	
Method	Source	Supervision	LLM	Acc@0.25	Acc@0.5	Acc@0.25	Acc@0.5	Acc@0.25	Acc@0.5
ScanRefer [5]	ECCV'20	Fully	-	67.6	46.2	32.1	21.3	39.0	26.1
InstanceRefer [39]	ICCV'21	Fully	-	77.5	66.8	31.3	24.8	40.2	32.9
3DVG-Transformer [41]	ICCV'21	Fully	-	77.2	58.5	38.4	28.7	45.9	34.5
BUTD-DETR [19]	ECCV'22	Fully	-	84.2	66.3	46.6	35.1	52.2	39.8
EDA [34]	CVPR'23	Fully	-	85.8	68.6	49.1	37.6	54.6	42.3
3D-VisTA [43]	ICCV'23	Fully	-	81.6	75.1	43.7	39.1	50.6	45.8
G3-LQ [32]	CVPR'24	Fully	-	88.6	73.3	50.2	39.7	56.0	44.7
MCLN [26]	ECCV'24	Fully	-	86.9	72.7	52.0	40.8	57.2	45.7
ConcreteNet [31]	ECCV'24	Fully	-	86.4	82.1	42.4	38.4	50.6	46.5
WS-3DVG [33]	ICCV'23	Weakly	-	-	-	-	-	27.4	22.0
LERF [20]	ICCV'23	Zero-Shot	CLIP [28]	-	-	-	-	4.8	0.9
OpenScene [25]	CVPR'23	Zero-Shot	CLIP [28]	20.1	13.1	11.1	4.4	13.2	6.5
LLM-Grounder [36]	ICRA'24	Zero-Shot	GPT-3.5 [24]	-	-	-	-	14.3	4.7
LLM-Grounder [36]	ICRA'24	Zero-Shot	GPT-4 [23]	-	-	-	-	17.1	5.3
ZSVG3D [40]	CVPR'24	Zero-Shot	GPT-4 [23]	63.8	58.4	27.7	24.6	36.4	32.7
3DVG-LLM	Ours	Zero-Shot	Llama-3.1-8B [8]	76.6	69.5	31.1	27.3	38.7	34.4

Table 1. Performance comparison with existing methods on ScanRefer validation set. We follow previous work [40] and report the grounding accuracy under different IoU threshold, and also report the accuracy of "*Unique*" (scenes with a single target object) and "*Multiple*" (scenes with distractors of the same class) subsets, along with overall performance.

## **4.** Experiments

#### **4.1. Implementation Details**

We use Llama-3.1-8B-Instruct [8] as our base model. We conduct post-training on Llama-3.1-8B-Instruct using our collected 3D visual grounding reasoning dataset, which consists of 3.2K 3D scenes, target queries, and structured reasoning processes. More details on training and implementation setup can be found in the appendix.

### **4.2. Benchmarks**

We adopt two large-scale 3D visual grounding benchmarks in our experiments, including ScanRefer [5] and
NR3D [4]. Both benchmarks leveraged 3D scenes collected
from ScanNet [7]. We follow the previous zero-shot method
ZSVG3D [40] and utilize the same object proposal generation procedure to obtain the 3D bounding boxes and their
class label for consistent evaluation and fair comparison.

ScanRefer. ScanRefer provides 51,500 natural language
descriptions from 800 different 3D scenes. We follow previous works [19, 36, 40] and adopt the validation set of
ScanRefer for evaluation, which contains 9,508 language
queries. We report ScanRefer's standard evaluation metrics: Accuracy@0.25 and Accuracy@0.5, where 0.25 and
0.5 are different IoU thresholds of 3D bounding boxes.

NR3D. NR3D contains 41.5K natural, free-form language captions collected from humans through a referring game, featuring multiple spatial relationships. Unlike ScanRefer, NR3D provides the GT 3D bounding boxes in the scene along with their corresponding object IDs. Given a language query, the visual grounding model is required to

predict the target's object ID that the given query refers 320 to. We adhere to the original NR3D benchmark's setup, 321 categorizing text queries into "Easy" (scenarios with only 322 one same-class distractor) and "Hard" (scenarios with mul-323 tiple same-class distractors). Additionally, queries are also 324 classified as "View-Dependent" or "View-Independent," de-325 pending on whether specific viewpoints are necessary for 326 correctly identifying the target object. 327

#### 4.3. Performance

We compare the performance of our **3DVG-LLM** with various recent and SoTA visual grounding methods on ScanRefer in Table. 1. And the performance comparison on NR3D in Table. 2. We categorize the existing methods into fully supervised, weakly supervised, and zero-shot depending on their training strategies. Note that we categorize our method as zero-shot because our model is not trained on ScanRefer, NR3D, or any other real-world 3D visual grounding data.

On ScanRefer, our model outperforms the SoTA zero-337 shot method ZSVG3D [40], while using a smaller open-338 source LLM Llama-3.1-8B. Notably, unlike ZSVG3D, 339 which requires multiple in-context examples during infer-340 ence, our model conducts reasoning and prediction with-341 out any in-context example at test time. However, a per-342 formance gap remains between our model and SoTA fully 343 supervised methods. We attribute this gap to the advantages 344 of SoTA supervised methods, which either leverage RGB 345 visual features [19] or benefit from supervised training to 346 generate more accurate object proposals [26, 34]-both of 347 which are critical for optimal grounding performance. 348

On NR3D, our model also outperforms the zero-shot 349 SoTA method ZSVG3D. Moreover, we also tested the per-350

Method	Source	Supervision	LLM	Easy	Hard	Dep.	Indep.	Overall
ReferIt3DNet [4]	ECCV'20	Fully	-	43.6	27.9	32.5	37.1	35.6
TGNN [15]	AAAI'21	Fully	-	44.2	30.6	35.8	38.0	37.3
3DRefTransformer [2]	WACV'22	Fully	-	46.4	32.0	34.7	41.2	39.0
InstanceRefer [39]	ICCV'21	Fully	-	46.0	31.8	34.5	41.9	38.8
FFL-3DOG [9]	ICCV'21	Fully	-	48.2	35.0	37.1	44.7	41.7
LanguageRefer [29]	CoRL'22	Fully	-	51.0	36.6	41.7	45.0	43.9
3DVG-Transformer [41]	ICCV'21	Fully	-	48.5	34.8	34.8	43.7	40.8
TransRefer3D [13]	MM'21	Fully	-	48.5	36.0	36.5	44.9	42.1

GPT-4 [3]

Llama-3.1-8B [8]

Llama-3.1-8B [8]

Table 2. Performance on NR3D benchmark. Queries are labeled as "*Easy*" (one distractor) or "*Hard*" (multiple distractors), and as "*View-Dependent*" or "*View-Independent*" based on viewpoint requirements for grounding. † denotes the use of the oracle class label. We utilize oracle class label to estimate the upper bound performance of our model.

Table 3. Ablation study on the performance after post-training. We report the accuracy on NR3D of (1) base model Llama-3.1-8B, (2) base model after post-training on our collected data.

ECCV'22

ICCV'23

CVPR'24

Ours

Ours

Fully

Weakly

Zero-Shot

Zero-Shot

Zero-Shot

Post-training	Easy	Hard	Dep.	Indep.	Overall
1	$44.6 \\ 63.0$	$22.7 \\ 36.7$	$31.1 \\ 40.1$	$34.5 \\ 54.5$	$33.3 \\ 49.3$

formance of our model using the oracle class label to estimate the upper bound performance. The increased accuracy after incorporating the oracle class label suggests that
our performance is mainly limited by the quality of object
proposal, which can be improved by using a better detector.

#### **356 4.4. Ablation Studies**

BUTD-DETR [19]

WS-3DVG [33]

ZSVG3D [40]

**3DVG-LLM** 

3DVG-LLM<sup>†</sup>

Performance under different training settings. To ver-357 ify the effectiveness of post-training on our collected data, 358 we evaluate our model and base model's performance on 359 NR3D in Table. 3. After post-training on our collected 360 data, the model's performance increased by a notable mar-361 gin across all categories in NR3D, showcasing that conduct-362 ing post-training on our data can effectively enhance the 363 model's performance on 3D visual grounding task. 364

Importance of Reasoning. We also train the base model 365 on our collected data without the reasoning process. In this 366 367 training setting, the model directly generates the final pre-368 diction without the reasoning process. The resulting model is similar to the existing 3D-LLMs [11, 14], which leans 369 to generate short and concise answers due to the character-370 istics of the existing 3D understanding training data. Ta-371 372 ble. 4 showcases that involving the reasoning process dur-

Table 4. Performance comparison on NR3D when training the base model with and without reasoning process.

60.7

27.3

46.5

50.4

63.0

48.4

18.0

31.7

31.0

36.7

46.0

21.6

36.8

37.3

40.1

58.0

22.9

40.0

42.0

54.5

54.6

22.5

39.0

40.4

49.3

373

374

385

Reasoning	Easy	Hard	Dep.	Indep.	Overall
1	42.2 63.0	$25.4 \\ 36.7$	$30.2 \\ 40.1$	$35.4 \\ 54.5$	$33.5 \\ 49.3$
•	00.0	00.1	10.1	01.0	10.0

Table 5. Experiments on extra prompt's effectiveness on reducing the spatial hallucination of GPT-4 $\circ$ . We report the accuracy on 50 left relationship queries and 50 right relationship queries.

Extra Prompt	Left Query Acc	Right Query Acc
	40.0	52.0
1	82.0	90.0

ing training can largely improves the performance, which demonstrates the importance of our collected data.

Data collection prompt. To verify the effectiveness of 375 our proposed extra prompt on reducing the spatial hallu-376 cination of  $GPT-4\circ$  during the data collection process, we 377 conduct ablation study on a subset of our collected data, 378 which consists of 100 visual grounding queries sampled 379 from left and right spatial relationship (50 queries each). In 380 Table. 5, GPT-40 achieves a large accuracy boost on both 381 types of relationship queries with our extra prompt, show-382 casing the effectiveness of our prompt design in reducing 383 spatial hallucination in our data collection framework. 384

#### **4.5. Qualitative Results**

We provide some qualitative results in Fig. 5. We include an in-domain query with similar spatial relationships in our 387



Figure 5. Qualitative results from NR3D. Green box and red box represents the prediction from our model and Llama-3.1-8B, respectively. Boxes are re-plotted for better visualization and some sentences are skipped due to space limits. Best viewed when zoomed in.

training data, and out-of-domain queries that have unseen
spatial relationships in training data. The out-of-domain
examples demonstrate the effectiveness of our post-training
strategy to enable LLM reasoning.

#### **392 5.** Conclusion

In this work, we introduce a novel 3D visual grounding and
reasoning data collection pipeline, which can generate posttraining data to enhance LLM for 3D visual grounding task.

In the data collection process, we also identify and miti-396 gate spatial hallucination in SoTA proprietary LLMs by de-397 signing an effective prompt strategy, leading to a more re-398 liable and efficient data collection process. Leveraging our 399 collected training data, we conduct post-training and intro-400 duce **3DVG-LLM**, a powerful reasoning LLM specifically 401 tailored for 3D visual grounding tasks. In our experiments, 402 3DVG-LLM achieves SoTA performance among all zero-403 shot methods on ScanRefer and NR3D, demonstrating its 404 effectiveness in reasoning for 3D visual grounding task. 405

#### References 406

- 407 [1] Eslam Abdelrahman, Mohamed Ayman, Mahmoud Ahmed, 408 Habib Slim, and Mohamed Elhoseiny. Cot3dref: Chain-of-409 thoughts data-efficient 3d visual grounding. arXiv preprint 410 arXiv:2310.06214, 2023. 2
- 411 [2] Ahmed Abdelreheem, Ujjwal Upadhyay, Ivan Skorokhodov, 412 Rawan Al Yahya, Jun Chen, and Mohamed Elhoseiny. 413 3dreftransformer: Fine-grained object identification in real-414 world scenes using natural language. In Proceedings of the IEEE/CVF winter conference on applications of computer 415 416 vision, pages 3941-3950, 2022. 7
- 417 [3] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ah-418 mad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, 419 Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 420 Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 421 2023. 2, 7
- 422 [4] Panos Achlioptas, Ahmed Abdelreheem, Fei Xia, Mohamed 423 Elhoseiny, and Leonidas Guibas. Referit3d: Neural listeners 424 for fine-grained 3d object identification in real-world scenes. 425 In Computer Vision-ECCV 2020: 16th European Confer-426 ence, Glasgow, UK, August 23-28, 2020, Proceedings, Part 427 I 16, pages 422-440. Springer, 2020. 6, 7
- 428 [5] Dave Zhenyu Chen, Angel X Chang, and Matthias Nießner. 429 Scanrefer: 3d object localization in rgb-d scans using natural 430 language. In European conference on computer vision, pages 202-221. Springer, 2020. 6
- [6] Zhenyu Chen, Ronghang Hu, Xinlei Chen, Matthias 432 433 Nießner, and Angel X Chang. Unit3d: A unified transformer 434 for 3d dense captioning and visual grounding. In Proceed-435 ings of the IEEE/CVF international conference on computer 436 vision, pages 18109–18119, 2023. 2
  - [7] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 5828-5839, 2017. 6
- 442 [8] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Ab-443 hishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil 444 Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The 445 llama 3 herd of models. arXiv preprint arXiv:2407.21783, 446 2024. 2. 6. 7
- 447 [9] Mingtao Feng, Zhen Li, Qi Li, Liang Zhang, XiangDong 448 Zhang, Guangming Zhu, Hui Zhang, Yaonan Wang, and Aj-449 mal Mian. Free-form description guided 3d visual graph net-450 work for object grounding in point cloud. In Proceedings of 451 the IEEE/CVF international conference on computer vision, 452 pages 3722-3731, 2021. 7
- 453 [10] Weixi Feng, Wanrong Zhu, Tsu-jui Fu, Varun Jampani, Ar-454 jun Akula, Xuehai He, Sugato Basu, Xin Eric Wang, and 455 William Yang Wang. Layoutgpt: Compositional visual plan-456 ning and generation with large language models. Advances 457 in Neural Information Processing Systems, 36:18225-18250, 458 2023 3
- 459 [11] Rao Fu, Jingyu Liu, Xilun Chen, Yixin Nie, and Wenhan Xiong. Scene-llm: Extending language model for 460 461 3d visual understanding and reasoning. arXiv preprint 462 arXiv:2403.11401, 2024. 7

- [12] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, 463 Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi 464 Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning 465 capability in llms via reinforcement learning. arXiv preprint 466 arXiv:2501.12948, 2025. 2, 3, 4 467
- [13] Dailan He, Yusheng Zhao, Junyu Luo, Tianrui Hui, Shaofei Huang, Aixi Zhang, and Si Liu. Transrefer3d: Entity-andrelation aware transformer for fine-grained 3d visual grounding. In Proceedings of the 29th ACM international conference on multimedia, pages 2344-2352, 2021. 7
- [14] Yining Hong, Haoyu Zhen, Peihao Chen, Shuhong Zheng, Yilun Du, Zhenfang Chen, and Chuang Gan. 3d-llm: Injecting the 3d world into large language models. Advances in Neural Information Processing Systems, 36:20482-20494, 2023. 7
- [15] Pin-Hao Huang, Han-Hung Lee, Hwann-Tzong Chen, and Tyng-Luh Liu. Text-guided graph neural networks for referring 3d instance segmentation. In Proceedings of the AAAI Conference on Artificial Intelligence, pages 1610-1618, 2021. 7
- [16] Shijia Huang, Yilun Chen, Jiaya Jia, and Liwei Wang. Multiview transformer for 3d visual grounding. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 15524-15533, 2022. 2
- [17] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditva Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. arXiv preprint arXiv:2410.21276, 2024. 4, 5
- [18] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. arXiv preprint arXiv:2412.16720, 2024. 2
- [19] Ayush Jain, Nikolaos Gkanatsios, Ishita Mediratta, and Katerina Fragkiadaki. Bottom up top down detection transformers for language grounding in images and point clouds. In European Conference on Computer Vision, pages 417–433. Springer, 2022. 2, 6, 7
- [20] Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. Lerf: Language embedded radiance fields. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 19729–19739, 2023.6
- [21] Junyu Luo, Jiahui Fu, Xianghao Kong, Chen Gao, Haibing Ren, Hao Shen, Huaxia Xia, and Si Liu. 3d-sps: Singlestage 3d visual grounding via referred point progressive selection. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 16454–16463, 2022. 2
- [22] Basak Melis Öcal, Maxim Tatarchenko, Sezer Karaoğlu, and Theo Gevers. Sceneteller: Language-to-3d scene generation. In European Conference on Computer Vision, pages 362-378. Springer, 2024. 3
- [23] OpenAI. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023. 6
- [24] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Car-517 roll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini 518 Agarwal, Katarina Slama, Alex Ray, et al. Training language 519

431

437

438

439

440

441

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

- 523 [25] Songyou Peng, Kyle Genova, Chiyu Jiang, Andrea
  524 Tagliasacchi, Marc Pollefeys, Thomas Funkhouser, et al.
  525 Openscene: 3d scene understanding with open vocabularies.
  526 In Proceedings of the IEEE/CVF Conference on Computer
  527 Vision and Pattern Recognition, pages 815–824, 2023. 6
- [26] Zhipeng Qian, Yiwei Ma, Zhekai Lin, Jiayi Ji, Xiawu Zheng, Xiaoshuai Sun, and Rongrong Ji. Multi-branch collaborative learning network for 3d visual grounding. In *European Conference on Computer Vision*, pages 381–398. Springer, 2025.
  2, 6
- [27] Yuxuan Qiao, Haodong Duan, Xinyu Fang, Junming Yang,
  Lin Chen, Songyang Zhang, Jiaqi Wang, Dahua Lin, and Kai
  Chen. Prism: A framework for decoupling and assessing
  the capabilities of vlms. *Advances in Neural Information Processing Systems*, 37:111863–111898, 2025. 2, 3
- [28] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya
  Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,
  Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning
  transferable visual models from natural language supervision. In *International conference on machine learning*, pages
  8748–8763. PMLR, 2021. 6
- 544 [29] Junha Roh, Karthik Desingh, Ali Farhadi, and Dieter
  545 Fox. Languagerefer: Spatial-language model for 3d visual
  546 grounding. In *Conference on Robot Learning*, pages 1046–
  547 1056. PMLR, 2022. 2, 7
- [30] Hao Shao, Shengju Qian, Han Xiao, Guanglu Song, Zhuofan Zong, Letian Wang, Yu Liu, and Hongsheng Li. Visual cot: Advancing multi-modal language models with a comprehensive dataset and benchmark for chain-of-thought reasoning. *Advances in Neural Information Processing Systems*, 37:8612–8642, 2025. 3
- [31] Ozan Unal, Christos Sakaridis, Suman Saha, and Luc
  Van Gool. Four ways to improve verbo-visual fusion for
  dense 3d visual grounding. In *European Conference on Computer Vision*, pages 196–213. Springer, 2025. 2, 6
- [32] Yuan Wang, Yali Li, and Shengjin Wang. G3-lq: Marrying hyperbolic alignment with explicit semantic-geometric modeling for 3d visual grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13917–13926, 2024. 2, 6
- [33] Zehan Wang, Haifeng Huang, Yang Zhao, Linjun Li, Xize
  Cheng, Yichen Zhu, Aoxiong Yin, and Zhou Zhao. Distilling coarse-to-fine semantic matching knowledge for weakly
  supervised 3d visual grounding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,
  pages 2662–2671, 2023. 6, 7
- [34] Yanmin Wu, Xinhua Cheng, Renrui Zhang, Zesen Cheng, and Jian Zhang. Eda: Explicit text-decoupling and dense alignment for 3d visual grounding. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 19231–19242, 2023. 2, 6
- 574 [35] Guowei Xu, Peng Jin, Hao Li, Yibing Song, Lichao Sun, and
  575 Li Yuan. Llava-cot: Let vision language models reason step576 by-step, 2025. 2, 3, 4

- [36] Jianing Yang, Xuweiyi Chen, Shengyi Qian, Nikhil Madaan, Madhavan Iyengar, David F Fouhey, and Joyce Chai. Llmgrounder: Open-vocabulary 3d visual grounding with large language model as an agent. In 2024 IEEE International Conference on Robotics and Automation (ICRA), pages 7694–7701. IEEE, 2024. 2, 6
  582
- [37] Yue Yang, Fan-Yun Sun, Luca Weihs, Eli VanderBilt, Alvaro Herrasti, Winson Han, Jiajun Wu, Nick Haber, Ranjay Krishna, Lingjie Liu, et al. Holodeck: Language guided generation of 3d embodied ai environments. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16227–16237, 2024. 3
- [38] Zhengyuan Yang, Songyang Zhang, Liwei Wang, and Jiebo Luo. Sat: 2d semantics assisted training for 3d visual grounding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1856–1866, 2021. 2
- [39] Zhihao Yuan, Xu Yan, Yinghong Liao, Ruimao Zhang, Sheng Wang, Zhen Li, and Shuguang Cui. Instancerefer: Cooperative holistic understanding for visual grounding on point clouds through instance multi-level contextual referring. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1791–1800, 2021. 6, 7
- [40] Zhihao Yuan, Jinke Ren, Chun-Mei Feng, Hengshuang Zhao, Shuguang Cui, and Zhen Li. Visual programming for zero-shot open-vocabulary 3d visual grounding. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 20623–20633, 2024. 2, 6, 7
- [41] Lichen Zhao, Daigang Cai, Lu Sheng, and Dong Xu. 3dvgtransformer: Relation modeling for visual grounding on point clouds. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2928–2937, 2021. 2, 6, 7
- [42] Chenming Zhu, Tai Wang, Wenwei Zhang, Kai Chen, and Xihui Liu. Scanreason: Empowering 3d visual grounding with reasoning capabilities. In *European Conference on Computer Vision*, pages 151–168. Springer, 2024. 2
- [43] Ziyu Zhu, Xiaojian Ma, Yixin Chen, Zhidong Deng, Siyuan Huang, and Qing Li. 3d-vista: Pre-trained transformer for 3d vision and text alignment. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2911– 2921, 2023. 6

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617