# S²-MLLM: Boosting Spatial Reasoning Capability of MLLMs for 3D Visual Grounding with Structural Guidance

Beining Xu[1*], Siting Zhu[1*], Zhao Jin[2], Junxian Li[1], Hesheng Wang[1†]

[1]Shanghai Jiao Tong University, [2]Nanyang Technological University

## Abstract

*3D Visual Grounding (3DVG) focuses on locating objects in 3D scenes based on natural language descriptions, serving as a fundamental task for embodied AI and robotics. Recent advances in Multi-modal Large Language Models (MLLMs) have motivated research into extending them to 3DVG. However, MLLMs primarily process 2D visual inputs and struggle with understanding 3D spatial structure of scenes solely from these limited perspectives. Existing methods mainly utilize viewpoint-dependent rendering of reconstructed point clouds to provide explicit structural guidance for MLLMs in 3DVG tasks, leading to inefficiency and limited spatial reasoning. To address this issue, we propose S²-MLLM, an efficient framework that enhances spatial reasoning in MLLMs through implicit spatial reasoning. We introduce a spatial guidance strategy that leverages the structure awareness of feed-forward 3D reconstruction. By acquiring 3D structural understanding during training, our model can implicitly reason about 3D scenes without relying on inefficient point cloud reconstruction. Moreover, we propose a structure-enhanced module (SE), which first employs intra-view and inter-view attention mechanisms to capture dependencies within views and correspondences across views. The module further integrates multi-level position encoding to associate visual representations with spatial positions and viewpoint information, enabling more accurate structural understanding. Extensive experiments demonstrate that S²-MLLM unifies superior performance, generalization, and efficiency, achieving significant performance over existing methods across the ScanRefer, Nr3D, and Sr3D datasets. Code will be available upon acceptance.*

## 1. Introduction

3DVG aims to locate referred objects in 3D scenes based on textual descriptions [2, 9, 28], serving as a key capability for
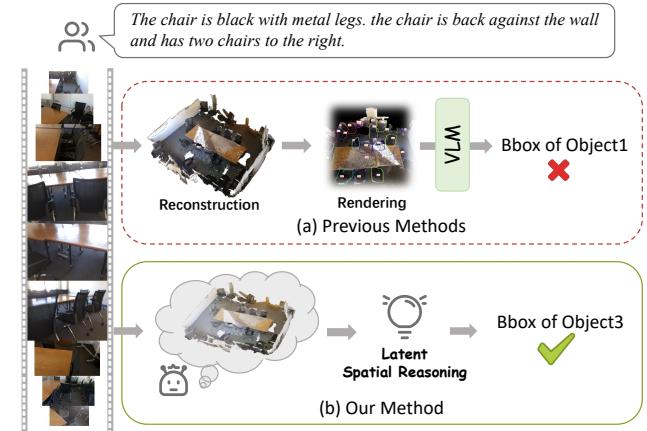


Figure 1. Comparison of previous methods and our method. (a) Previous methods typically reconstruct point clouds of 3D scenes explicitly and then render 2D images to obtain structure guidance. (b) Our method leverages spatial guidance to understand the 3D structure during training, allowing the model to perform implicit spatial reasoning in the latent space without requiring point-cloud reconstruction at inference.

embodied AI [39, 95, 97, 98] and augmented reality [35, 54, 71, 96]. Compared to 2D visual grounding, 3DVG demands a thorough understanding of spatial relationships and 3D scene structures [19, 26, 36, 65, 75], thereby posing greater challenges.

Traditional approaches [2, 6, 11, 16, 17, 23, 24, 47, 64, 70, 80, 89] primarily rely on training within a limited scope of datasets, resulting in limited generalization and scalability in real-world scenes [29, 34]. To overcome this limitation, leveraging MLLMs [1, 31–33, 60, 72] for 3DVG through their reasoning and generalization abilities has emerged as a promising direction. However, a significant gap remains between the 2D-centric training of MLLMs and spatial understanding capabilities demanded by 3DVG [25, 34, 84]. MLLMs are primarily trained on extensive image-text data without exposure to the 3D physical world [41, 50], which involves additional factors such as depth, viewpoint, layout, and structure [35, 36]. Therefore, MLLMs are incapable of understanding 3D scenes solely from 2D images [19, 91].

---

*These authors contributed equally.
†Corresponding Author

Existing approaches [29, 34, 50] attempt to compensate for this limitation by explicitly reconstructing point clouds of 3D scenes. These reconstructed point clouds are then rendered into 2D representations that preserve the layout and spatial relations, such as multi-view images [29, 34] or Bird's Eye View (BEV) images [50]. MLLMs can directly process these rendered images to select target objects. However, rendered images from specific viewpoints are unable to reflect the comprehensive structure of 3D scenes [38, 56] and are inevitably affected by viewpoint selection as well as occlusion. Moreover, these methods need to reconstruct point clouds explicitly during inference, leading to low efficiency.

To address these challenges, we propose a novel framework $S^2$-MLLM, which enhances the spatial reasoning of MLLMs for 3DVG. As shown in Fig. 1, our key insight is to encourage the model to implicitly internalize 3D structure awareness during training. This design enables our $S^2$-MLLM to reason implicitly about 3D scenes within the latent feature space, without requiring extra reconstruction or rendering at inference. Specifically, we propose to introduce spatial guidance by leveraging the capability of feed-forward 3D reconstruction [59, 61, 74]. These reconstruction techniques can directly derive 3D structures from multi-view inputs, demonstrating an inherent capability for spatial understanding [68, 87]. Building upon this property, we integrate the reconstruction objective into the training pipeline through end-to-end joint optimization, enabling the model to learn structure-aware visual representations and spatial reasoning capability.

Moreover, we propose a structure-enhanced (SE) module to further enhance the spatial understanding of MLLMs from the perspectives of position and viewpoint. We observe that (i) without explicit association between position cues (*e.g.*, 3D coordinates, camera rays) and visual appearance, MLLMs struggle to understand fine-grained spatial relations such as distance, direction, and relative relations; (ii) MLLMs are pretrained on independent image–text pairs, making it difficult to maintain semantic consistency across views. Therefore, we first integrate a multi-level position encoding to enhance the modeling of position and viewpoint information. Secondly, we employ inter-view attention mechanism to enforce semantic alignment during viewpoint transitions across views. Within each view, we conduct intra-view attention to capture dependencies between patches, which improves local and global context understanding. We conduct extensive experiments on both in-domain [2, 9] and out-of-domain benchmarks [4, 42]. $S^2$-MLLM achieves outstanding performance in terms of accuracy, efficiency, and generalization. This trade-off enables practical deployment in real-world applications and embodied robotics.

Overall, we provide the following contributions:

- We propose $S^2$-MLLM, an effective framework that enhances the spatial reasoning of MLLMs, thereby improving performance for 3DVG.
- We introduce a **spatial guidance** strategy that utilizes the capability of feed-forward 3D reconstruction to encourage our model to perform latent spatial reasoning.
- We design a **structure-enhanced module (SE)** module that enhances spatial understanding through modeling position and viewpoint, as well as employs intra-view and inter-view attention to strengthen contextual alignment and cross-view consistency.
- Extensive experiments on *ScanRefer*, *Nr3D*, and *Sr3D* datasets demonstrate that $S^2$-MLLM significantly outperforms baselines. Multiple out-of-domain evaluations indicate the generalization ability of $S^2$-MLLM.

## 2. Related Work

### 2.1. Supervised 3D Visual Grounding

3DVG is the task of 3D object localization in points or RGB-D scans via natural language descriptions [2, 6, 11, 16, 23, 24, 47, 64, 70, 80, 89]. Existing traditional methods typically adopt a fully supervised paradigm, which can be categorized into two-stage methods and one-stage methods depending on their network architectural designs. Two-stage methods [2, 8–11, 80, 85, 89] follow a proposal–matching pipeline: 3D object proposals are obtained from pretrained detector or segmentor [27, 49, 53] and then matched with the query. In recent years, more methods [3, 17, 23, 44, 76, 100] integrates multi-modal information, such as 2D images and multi-view contexts. To avoid relying on pre-trained proposal generators, one-stage methods [16, 24, 40, 70] directly regress the 3D box by densely aligning language with point-level features. Although these methods have achieved impressive accuracy on public benchmarks, they often suffer from limited generalization when applied to real-world 3D scenes.

### 2.2. Zero-shot 3D Visual Grounding

The reasoning and generalization abilities of LLMs have motivated several studies to explore zero-shot 3DVG. These approaches typically decompose the 3DVG task into a sequence of sub-tasks that can be processed by LLMs, reducing the reliance on large-scale and high-quality 3D annotations. However, as LLMs cannot directly process 3D information, the modality gap remains a fundamental challenge. Early approaches [30, 72, 73, 81] convert object attributes into textual descriptions, then use the reasoning ability of LLMs to choose the object that best matches the query. This design completely ignores the importance of scene-level context, which is essential for 3DVG. VLM-based (Vision Language Models) methods [29, 34, 82] mainly employ VLMs for query analysis, viewpoint selection, and

target object selection. Complex structures and spatial relations in 3D scenes cannot be captured solely from textual descriptions or images from specific viewpoints, which causes these methods to fail in situations that require spatial reasoning.

## 2.3. MLLMs for 3D Scene Understanding

MLLMs have also been recently extended to 3D scene understanding [7, 92], enabling the construction of a generalist model that is capable of handling multiple 3D tasks, including 3DVG [15, 19, 79, 99]. To enable LLMs to understand 3D scenes, early approaches [12, 18] integrate the point cloud encoder with LLMs, relying on large-scale training to enhance LLMs' ability to process point clouds. [13, 21, 67, 86] incorporate object-centric 3D representations that further improves performance. However, point-cloud–based approaches suffer from a large modality gap and limited 3D annotations [50, 72]. These limitations motivate a shift toward learning 3D scenes from 2D image sequences. Llava-3d [94] projects multi-view CLIP features into 3D voxels and aggregates them to recover coarse 3D scene structure. Recent methods [68, 90, 91] instead treat 3D scenes as dynamic video sequences. GPT4Scene [50] relies on the BEV image rendered from reconstructed point clouds to obtain global information, while ROSS3D [58] introduces cross-view and BEV generation tasks to encourage understanding of the layout. However, the former is time-consuming during inference due to reconstruction, and the latter is easily influenced by obstructions and noise, making both unreliable for learning consistent 3D structures.

## 3. Methods

We propose a task-specific model for 3D visual grounding based on a pre-trained MLLM [88], which effectively integrates information from texts and image sequences Our method represents 3D scenes as video sequences, which preserves abundant texture and semantic cues of 2D images. Given multi-view RGB-D images $\{(I_v, D_v)\}_{v=1}^V$, related camera parameters of 3D scene with candidate objects $\{o_i\}_{i=1}^N$, and a natural language description $Q$, our model predicts the 3D bounding box and category of target object $\hat{o}$ that best matches the description. The overview of our method is shown in Fig 2.

In the following, we describe our framework in detail. Sec. 3.1 introduces the spatial guidance strategy. Sec. 3.2 details the intra-view and inter-view attention mechanisms. Sec. 3.3 presents the multi-level position encoding. Finally, Sec. 3.4 describes the overall training objectives.

### 3.1. Spatial Guidance Strategy

MLLMs are trained on massive collections of images, text, and videos, enabling them to integrate visual and textual information effectively [69, 78]. However, the 2D visual priors encoded in MLLMs are insufficient for inferring 3D structure from RGB images [91], including global layout, geometric properties, cross-view correspondences, and fine-grained spatial relations. The lack of 3D structure understanding restricts MLLMs from understanding the 3D physical world [41, 83]. Existing methods rely on reconstructing point clouds to produce BEV inputs [50] or rendered images [29, 34] to provide structure guidance for MLLMs. However, rendered images are easily influenced by the selection of viewpoints and occlusions [38, 56]. For example, relative spatial relations may change when observed from different viewpoints. Consequently, relying solely on view-specific observations often leads to incomplete or inconsistent scene representations. Moreover, these methods are time-consuming during inference due to reconstruction.

To bridge this gap, our goal is to equip the MLLM with the implicit understanding of 3D structure, thereby enhancing its spatial reasoning ability. Considering that feed-forward 3D reconstruction [59, 61, 74] directly predicts 3D dense structure from multi-view RGB images, these methods inherently capture the structural understanding of 3D scenes. To leverage this capability, we incorporate spatial guidance into MLLM by jointly optimizing the reconstruction objective.

**Overall Architecture** We build the reconstruction branch on top of Fast3R [74], which consists of a ViT encoder, a fusion transformer, and decoding heads. Our reconstruction branch contains an encoder $\mathcal{E}_v$ from MLLM, a projection layer $\mathcal{P}$, and a reconstruction decoder $\mathcal{D}$. $\mathcal{D}$ is composed of the fusion transformer and decoding heads of Fast3R. Since Fast3R [74] is optimized only for reconstruction and provides limited semantic alignment, we replace its ViT with the visual encoder from the MLLM to ensure both reconstruction and 3DVG rely on a unified representation space. Considering that 3D reconstruction relies on dense and structure-aware features rather than semantic features, we introduce a projection layer $\mathcal{P}$ to align and normalize the representations. Meanwhile, we retain only the fusion transformer and decoding heads of Fast3R [74], since these blocks are responsible for multi-view geometric aggregation and pointmap prediction.

**Training Strategy** Our reconstruction branch predicts the local pointmap $X_L$ and global pointmap $X_G$.

$$X_L, X_G = \mathcal{D}\big(\mathcal{P}(\mathcal{E}_v(I))\big). \tag{1}$$

For reconstruction, we adopt the confidence-weighted pointmap regression loss introduced in Fast3R [74]. During training, we optimize the reconstruction loss jointly with the other objectives required by our model. This training strategy encourages our model to internalize 3D structure within its latent space, which ultimately enhances spatial reasoning and improves performance in 3DVG.
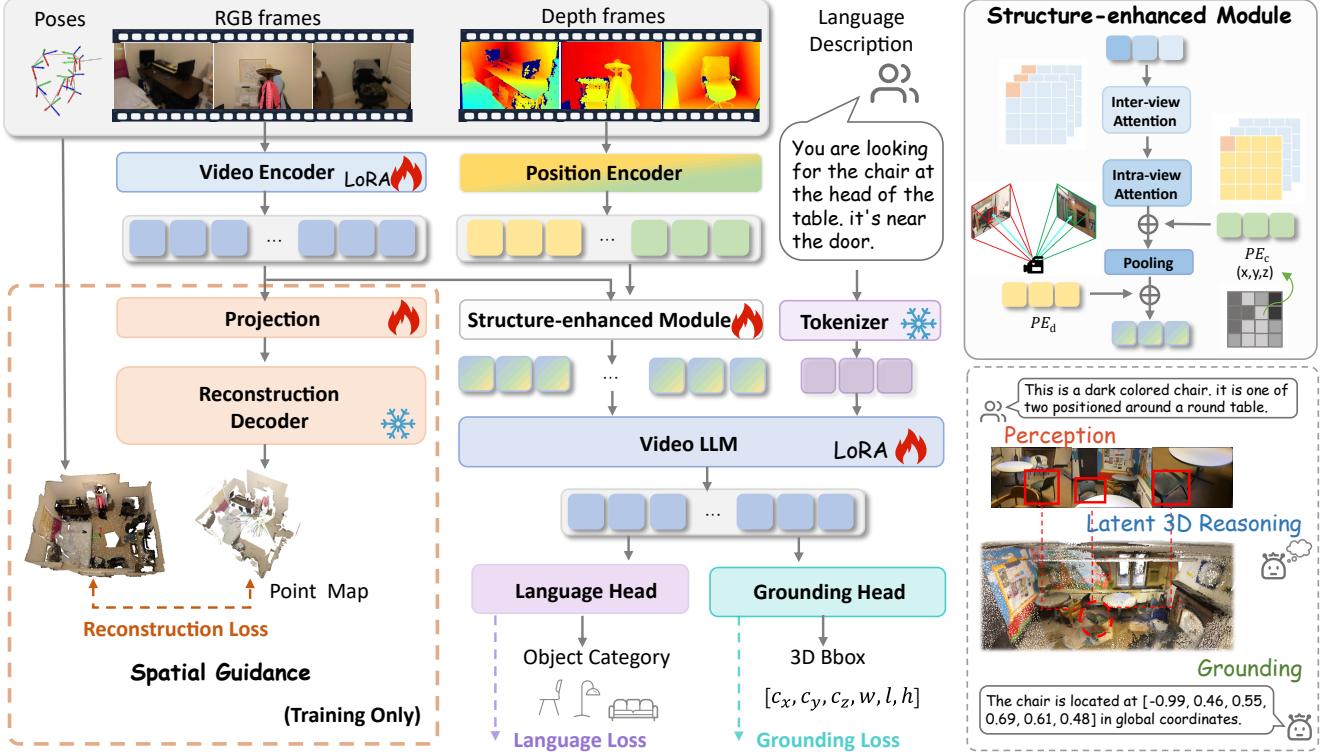
Figure 2. The Framework of S²-MLLM. Our model takes sampled multi-view RGB-D frames, related camera parameters, language descriptions, and the bounding boxes of object proposals as inputs. The shared video encoder and position encoder extract visual and geometric features from RGB-D frames. The structure-enhanced module (SE) integrates visual features and position information to form the visual input for the LLM. The Video LLM jointly processes the visual input and the tokenized query, enabling cross-modal understanding and reasoning. The grounding head predicts the 3D bounding box (Bbox) of the target object. The language head generates the category of the target object. while the reconstruction decoder predicts the point map to provide reconstruction supervision. Notice that the reconstruction is used only at training time.

In practice, the reconstruction loss converges early and provides stable structure supervision, encouraging our model to acquire structure-aware features in the early stage. It enables our model to benefit from structure understanding without compromising its multimodal alignment capability.

### 3.2. Intra-view and Inter-view Attention

MLLMs cannot maintain semantic consistency across views [77] and establish structure correspondences within views. For instance, the model is unable to determine whether chairs observed from different viewpoints correspond to the same physical object in 3D space. Multi-view features naturally form two independent types of interactions: within-view and across views. It is similar to temporal and spatial factorization in video modeling. Inspired by [5], we adopt a divided-attention design to separately capture spatial relations within each view and semantic correspondences across views.

Given multi-view image features $f \in \mathbb{R}^{B \times (V \cdot H \cdot W) \times dim}$, where $V$ is the number of views, $N = H \cdot W$ is the number of patches per view, and $dim$ is the feature dimension.

For each patch index $s \in \{1, \ldots, N\}$, we gather patches across views $f_s^{inter} \in \mathbb{R}^{B \times V \times dim}$ to compute the queries, keys, and values after LayerNorm. Inter-view attention operates across views for each patch index. The intra-attention within each view is formulated in the same way, where patches are grouped by view index $v$ rather than by patch index $s$.

### 3.3. Multi-level Position Encoding

Although spatial guidance encourages our MLLM to infer the 3D structure, it still lacks explicit associate visual representations with 3D positions, thereby limiting the understanding of fine-grained spatial relations. For example, while MLLMs can recognize objects in an image, they cannot accurately reason about their relative spatial relationships. To address this issue, we enhance visual representations with multi-level position embeddings to incorporate explicit spatial information.

Each pixel in an RGB-D frame can be projected into a 3D coordinate. Inspired by NeRF [43], each pixel also lies on a specific camera ray. To obtain geometry-enhanced vi-

4

sual representations, we encode both 3D coordinates and the camera ray's viewing direction together with the visual embeddings. This representation allows the model to explicitly associate visual representations with positions and viewing directions in 3D space.

Given an RGB frame, depth image $D$, camera intrinsic matrix $K \in \mathbb{R}^{3\times3}$ and extrinsic matrix $T \in \mathbb{R}^{4\times4}$, we calculate the 3D coordinates $p_{\text{world}} = (x, y, z)$ and camera ray's viewing direction $\boldsymbol{\omega}$ of each pixel $(u, v)$ of the image in the global coordinate system.

Formally, for a pixel coordinate with depth value $d = D(u, v)$, the corresponding 3D point in the world coordinate system is computed as:

$$p_{\text{world}} = T \begin{bmatrix} d\,K^{-1}(u, v, 1)^{\top} \\ 1 \end{bmatrix}. \qquad (2)$$

Each pixel corresponds to a camera ray that originates from the camera center and passes through the 3D point projected from the camera coordinate system. We denote a ray by $\boldsymbol{\omega} = (o, p, r)$, where $o$ is the ray origin, $p$ is the 3D point associated with the pixel, and $r$ is the viewing direction. Formally, the ray origin in world coordinates is given by the translation matrix $t$ of the extrinsic matrix: $o_{\text{world}} = t$. The ray termination point is equal to the 3D global coordinate $p_{\text{world}}$ of the pixel. The normalized camera ray's viewing direction $r$ is given by:

$$r = \frac{p_{\text{world}} - o_{\text{world}}}{\|p_{\text{world}} - o_{\text{world}}\|_2}. \qquad (3)$$

After obtaining visual features $f$ from the visual encoder $\mathcal{E}_v$ with a small patch size, we adopt sinusoidal position encoding $\phi(\cdot)$ to encode the average 3D coordinate $p^i_{\text{world}} = (x_i, y_i, z_i)$ of each patch $i$ in the global coordinate system following [91]. The 3D coordinate embeddings are added to the visual feature $f_i$ of each patch as $f^p_i$.

Then, we perform feature aggregation via average pooling across neighboring patches to obtain a context-enriched feature $f^s_i$. Lastly, we introduce a learnable positional encoding $\psi(\cdot)$ to encode the camera ray direction $r_i$ of each patch center. $\psi(\cdot)$ is implemented as a multi-layer perceptron (MLP). The final position-aware visual representation of each patch is

$$f^{\text{vis}}_i = \text{AvgPool}\big(f_i + \phi(p^i_{\text{world}})\big) + \psi(r_i). \qquad (4)$$

### 3.4. Overall Loss Function

**Visual Grounding Loss** We follow prior work [19, 23, 67, 91, 93, 99] and formulate the 3DVG task as a classification for objects proposals. For each bounding box $b$ of an object in the scene, we obtain its feature $f_{\text{obj}}$ by averaging the patch features whose projected points lie inside $b$ with over $50\%$ coverage. $f_{\text{obj}}$ is further added with the 3D

position embedding of the object center. Given the hidden state $h$ of the $\langle\text{ground}\rangle$ token, we compute their similarity and optimize it using an InfoNCE [45] loss. Formally, the grounding loss is defined as

$$\mathcal{L}_{\text{ground}} = \text{InfoNCE}\big(f_{\text{obj}}, h\big), \qquad (5)$$

where $\text{InfoNCE}(\cdot)$ denotes the cross-entropy contrastive objective.

**Reconstruction Loss** Given the predicted pointmap $\hat{X}$ with confidence scores $\hat{\Sigma}$ and the ground-truth pointmap $X$, the regression loss is defined as:

$$\ell_{\text{regr}}(\hat{X}, X) = \left\| \tfrac{1}{\hat{z}}\hat{X} - \tfrac{1}{z}X \right\|_2, \quad z = \tfrac{1}{|X|}\sum_{x\in X}\|x\|_2, \quad (6)$$

and the final pointmap loss is formulated as:

$$\mathcal{L}_X(\hat{\Sigma}, \hat{X}, X) = \tfrac{1}{|X|}\sum \hat{\Sigma}_+ \cdot \ell_{\text{regr}}(\hat{X}, X) + \alpha\log(\hat{\Sigma}_+), \qquad (7)$$

where $\hat{\Sigma}_+ = 1 + \exp(\hat{\Sigma})$. Formally, the objective of reconstruction can be obtained by adding the loss of local and global point maps:

$$\mathcal{L}_{\text{recon}} = \mathcal{L}_{X_G} + \mathcal{L}_{X_L}. \qquad (8)$$

**Language Loss** Visual grounding errors also arise from misclassification of the target object category [34], which are common in complex scenes. For example, the model may focus on the correct spatial position but still classify a stool as a chair or confuse a nightstand with a cabinet. To enforce semantic consistency between the predicted object category and the type of target objects, we incorporate an additional language-guidance. We supervise text generation using a cross-entropy loss $\mathcal{L}_{\text{lang}}$ to encourage the model produce text responses like "*The [object category] is located at <ground> in the global coordinates*".

The final training objective combines all components:

$$\mathcal{L} = \lambda_{\text{g}}\mathcal{L}_{\text{ground}} + \lambda_{\text{r}}\mathcal{L}_{\text{recon}} + \lambda_{\text{l}}\mathcal{L}_{\text{lang}}, \qquad (9)$$

where $\lambda_{\text{g}}, \lambda_{\text{r}}, \lambda_{\text{l}}$ are balancing weights.

Concretely, we construct training queries in the form of a bounding box with the object category using a unified prompt template. Please refer to the supplementary for prompt details.

## 4. Experiments

### 4.1. Dataset and Metrics

We evaluate our method and baselines on ScanRefer [9] and ReferIt3D [2]. The ReferIt3D [2] benchmark contains two subsets, Nr3D and Sr3D, which provide natural and synthetic referring expressions, respectively. The ScanRefer [9] dataset includes 51,583 descriptions of 11,046 objects across ScanNet [14] scenes. we train our model on

Table 1. Accuracy comparison on Scanrefer [9] validation set at IoU thresholds of 0.25 and 0.5. We report results on the Unique subset (single-object scenes), the Multiple subset (scenes with same-class distractors), and the overall accuracy. * denotes results obtained by LoRA [20] fine-tuning with the same parameter size as ours, while other settings follow the original paper.

| Method | Venue | LLM | Unique | | Multiple | | Overall | |
|---|---|---|---|---|---|---|---|---|
| | | | Acc@0.25 | Acc@0.5 | Acc@0.25 | Acc@0.5 | Acc@0.25 | Acc@0.5 |
| ScanRefer [9] | ECCV'20 | - | 67.6 | 46.2 | 32.1 | 21.3 | 39.0 | 26.1 |
| InstanceRefer [80] | ICCV'21 | - | 77.5 | 66.8 | 31.3 | 24.8 | 40.2 | 32.9 |
| 3DVG-T [89] | ICCV'21 | - | 77.2 | 58.5 | 38.4 | 28.7 | 45.9 | 34.5 |
| BUTD-DETR [24] | ECCV'22 | - | 84.2 | 66.3 | 46.6 | 35.1 | 52.2 | 39.8 |
| EDA [70] | CVPR'23 | - | 85.8 | 68.6 | 49.1 | 37.6 | 54.6 | 42.3 |
| 3D-VisTA [99] | ICCV'23 | - | 81.6 | 75.1 | 43.7 | 39.1 | 50.6 | 45.8 |
| VPP-Net [55] | CVPR'24 | - | 86.1 | 67.1 | 50.3 | 39.0 | 55.7 | 43.3 |
| G3-LQ [63] | CVPR'24 | - | **88.6** | 73.3 | 50.2 | 39.7 | 56.0 | 44.7 |
| MCLN [51] | ECCV'24 | - | 86.9 | 72.7 | 52.0 | 40.8 | 57.2 | 45.7 |
| ConcreteNet [57] | ECCV'24 | - | 86.4 | **82.1** | 42.4 | 38.4 | 50.6 | 46.5 |
| ViewSRD [22] | ICCV'25 | - | 82.1 | 68.2 | 37.4 | 39.0 | 45.4 | 36.0 |
| BUTD-DETR [24]+AugRefer [62] | AAAI'25 | - | 85.2 | 69.0 | 47.7 | 37.2 | 53.9 | 42.4 |
| EDA [70]+AugRefer [62] | AAAI'25 | - | 86.2 | 70.8 | 50.0 | 39.1 | 55.7 | 44.0 |
| TSP3D [16] | CVPR'25 | - | 87.3 | 71.4 | 51.0 | 42.4 | 56.5 | 46.7 |
| WS-3DVG [66] | ICCV'23 | - | - | - | - | - | 27.4 | 22.0 |
| Video-3D-LLM [91]* | CVPR'25 | LLaVA-Video-7B [88] | 82.3 | 72.5 | 47.2 | 42.0 | 54.1 | 47.9 |
| LERF [30] | ICCV'23 | CLIP [52] | - | - | - | - | 4.8 | 0.9 |
| OpenScene [48] | CVPR'23 | CLIP [52] | 20.1 | 13.1 | 11.1 | 4.4 | 13.2 | 6.5 |
| ZSVG3D [81] | CVPR'24 | GPT-4 turbo [46] | 63.8 | 58.4 | 27.7 | 24.6 | 36.4 | 32.7 |
| SeeGround [34] | CVPR'25 | Qwen2-VL-72B [60] | 75.7 | 68.9 | 34.0 | 30.0 | 44.1 | 39.4 |
| **S²-MLLM** | Ours | LLaVA-Video-7B[88] | 87.4 | 77.8 | **52.4** | **46.6** | **59.2** | **52.7** |

Table 2. Accuracy comparison on Nr3D and Sr3D [2] validation set with both predicted at IoU thresholds of 0.25 and ground-truth bounding boxes as input.

| Method | Venue | Pred | | GT | |
|---|---|---|---|---|---|
| | | Sr3D | Nr3D | Sr3D | Nr3D |
| InstanceRefer [80] | ICCV'21 | 31.5 | 29.9 | 48.0 | 35.6 |
| LanguageRefer [11] | CoRL'22 | 39.5 | 28.6 | 56.0 | 43.9 |
| BUTD-DETR [24] | ECCV'22 | 52.1 | 43.3 | 67.0 | 54.6 |
| EDA [70] | CVPR'23 | 49.9 | 40.7 | 68.1 | 52.1 |
| MCLN [51] | ECCV'24 | **53.9** | 46.1 | **68.4** | **59.8** |
| ZSVG3D [81] | CVPR'24 | – | – | – | 39.0 |
| SeeGround [34] | CVPR'25 | – | – | – | 46.1 |
| **S²-MLLM** | Ours | **53.9** | **50.6** | 63.2 | **59.8** |

the combined ScanRefer [9], Nr3D, and Sr3D [2] datasets and evaluate it separately, following the setting of [91, 99]. This setup enables learning from both natural and template-based expressions while evaluating the model under varying levels of referential complexity. For ScanRefer, we report accuracy at IoU thresholds of 0.25 and 0.5. For ReferIt3D [2], we report results under two settings: (i) using ground-truth bounding boxes (GT) as object proposals, which is a commonly used setting in previous methods [24, 51], and (ii) using predicted bounding boxes (Pred), which simulates realistic inference conditions with noisy detections.

## 4.2. Comparison Methods

For ScanRefer [9], we compare S²-MLLM with traditional full-supervised methods [9, 70, 80], weakly-supervised method [66], and LLM-based methods including zero-shot methods [34] and Video-3D-LLM [91] (for fair comparison, we also apply LoRA [20] training on it). For ReferIt3D [2], we evaluate a series of recent open-source and reproducible methods [11, 24, 34, 51, 70, 80, 81]. The results of [11, 24, 51, 70, 80] taking predicted bounding boxes as inputs are directly taken from [16].

## 4.3. Implementation Details

We fine-tune our model based on LLaVA-Video [88] 7B with LoRA [20]. Training is performed across the combined dataset for one epoch with a batch size of 8 and a warmup ratio of 0.05. The learning rate is scheduled to peak for the LLM and for the video encoder during warmup. All experiments are conducted only on a single A100 GPU (80GB). The temperature $\tau$ for InfoNCE loss is 0.07. During training, ground-truth bounding boxes are provided as object proposals. We employ object proposals generated by Mask3D [53] as predicted bounding boxes at inference. During training, we freeze the reconstruction decoder and fine-tune the projection layer, the visual encoder, and the language model. During inference, the structural guidance branch is disabled. Please refer to the supplementary for more implementation details.
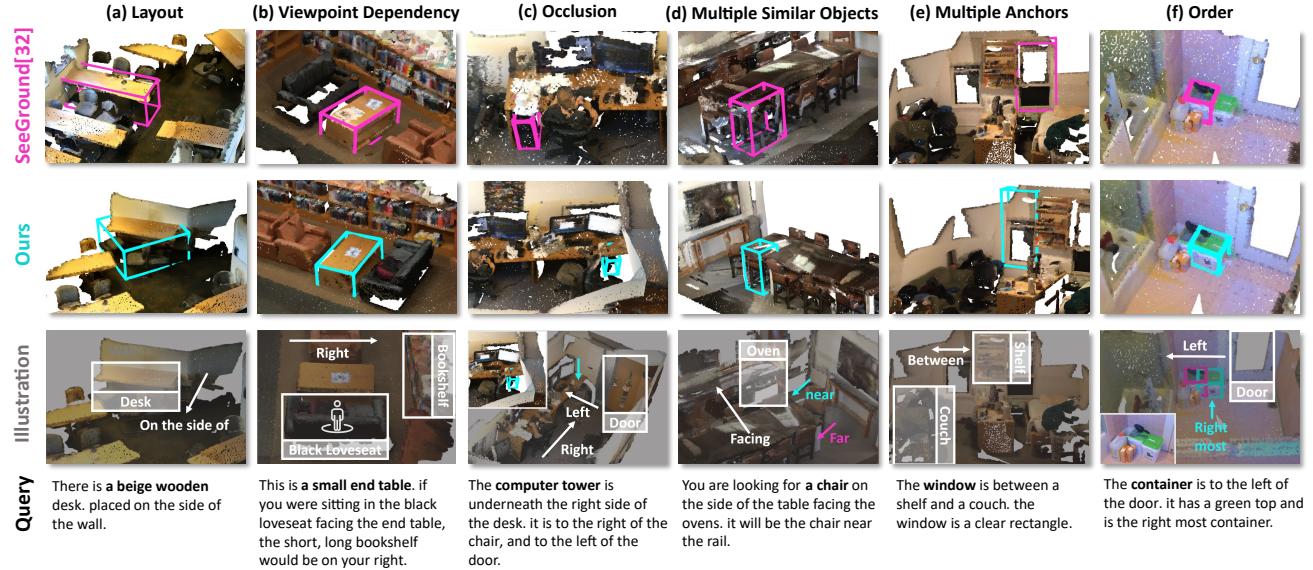
Figure 3. Qualitative comparison of 3DVG results in challenging spatial understanding cases in 3DVG. Incorrect predictions are highlighted in magenta, correct ones in cyan, and key spatial relations are underlined. Specifically, $S^2$-MLLM (a) reasons with scene layout priors. (b) understands the specified viewpoint. (c) predicts through structure cues despite partial occlusion. (d) distinguishes similar objects by jointly reasoning over multiple spatial relationships. (e) accurately handles spatial relations involving multiple reference objects. (f) understands relative positioning. Overall, $S^2$-MLLM demonstrates more reliable spatial understanding than the previous method.

## 4.4. Experiment Results

**ScanRefer** Tab. 1 presents the performance of our method and prior approaches on ScanRefer [9]. Overall, our method achieves the best performance across all evaluation settings, reaching Acc@0.25 of 59.2% and Acc@0.5 of 52.7% in overall accuracy. Notably, our method achieves a 10.0% improvement over the previous SOTA method in scenes containing multiple similar objects on Acc@0.5, which requires the model to jointly understand complex spatial relationships and accurately identify object attributes mentioned in the query (*e.g.*, the white chair next to the desk under the window). Meanwhile, our method improves by 4.8% on Acc@0.5 compared to previous methods. Additionally, compared with Video-3D-LLM [91] applying LoRA-finetuning, our method improves by over 5.1% on each metric. This improvement is attributed to our spatial guidance, which enables our model to reason about 3D layout and spatial relations, especially in complex environments.

**ReferIt3D** As shown in Tab. 2, our method achieves competitive results under the Pred setting, obtaining Acc@25 of 53.9% on Sr3D and 50.6% on Nr3D. These gains mainly come from our spatial guidance and SE module, which produce structure-aware visual features and improve robustness to bounding-box misalignment and proposal noise. Under the GT setting, $S^2$-MLLM achieves 59.8% on Nr3D. Although several methods perform better with ground-truth boxes, this idealized setting is rarely achievable in real applications. Therefore, experimental results under the

Table 3. **Efficiency comparison.** We report the training cost (in GPU hours), trainable parameters (in MB), and the inference latency (in seconds). $t_0$ represents the additional inference time of reconstructing point clouds.

| Method | GPU Hours $\downarrow$ | Trainable Parameters (MB)$\downarrow$ | Latency (s)$\downarrow$ |
|---|---|---|---|
| Video-3D-LLM [91] | 256 | 8078.79 | 1.04 |
| SeeGround [34] | - | - | $3.97 + t_0$ |
| $S^2$-MLLM(w/o SG) | 65 | 1763.53 | 1.16 |
| $S^2$-MLLM(Full) | 72 | 1767.50 | 1.16 |

predicted-box setting reflect real-world performance more accurately. The relatively lower performance on Sr3D under the GT setting is mainly due to the template-based queries. Traditional fully supervised methods can easily exploit these fixed patterns to match them with object features. In contrast, our performance on Nr3D highlights the advantage of $S^2$-MLLM in understanding natural human queries and aligning them with 3D spatial information.

**Qualitative Results** To further demonstrate the effectiveness of our approach, we visualize the visual grounding results of $S^2$-MLLM and SeeGround [34]. As shown in Fig. 3, $S^2$-MLLM demonstrates clear advantages in challenging spatial reasoning scenes, including layout understanding, viewpoint specification, occlusion, and grounding among multiple similar objects. The visualization results indicate that $S^2$-MLLM can capture 3D structure and understand complex spatial relations, indicating the effectiveness of spatial guidance.

Table 4. **Out-of-Distribution (OOD) Evaluation** on Multiscan [42] and ArkiScenes [4].

| Method | MultiScan | | ArtiScenes | |
|---|---|---|---|---|
| | Acc@0.25 | Acc@0.5 | Acc@0.25 | Acc@0.5 |
| MCLN [51] | 12.91 | 6.00 | 17.21 | 6.35 |
| SeeGround [34] | 53.41 | 53.41 | 38.82 | 38.43 |
| $S^2$-MLLM | **59.13** | **53.62** | **43.26** | **39.84** |

### 4.5. Efficiency and Generalization

**Efficiency** Tab. 3 compares the training efficiency and inference latency of our method and [34, 91]. All experiments are conducted on a single A100 GPU (80G). Due to limited computational resources, we utilize Qwen2-VL-7B [60] to measure the inference latency of See-Ground [34]. Our $S^2$-MLLM requires only $25\%$ of the trainable parameters and GPU hours required by the full-parameter fine-tuning method [91], while achieving higher accuracy. Although zero-shot methods are training-free, they typically perform multiple API calls within a single inference, leading to nearly $4\times$ higher inference time of See-Ground [34]. Notice that SeeGround [34] needs extra inference time $t_0$ for point clouds reconstruction, which further increases the actual latency. We further analyze the overhead introduced by our proposed spatial guidance. SG adds only around $10\%$ training time, with negligible additional trainable parameters. During inference, $S^2$-MLLM avoids extra point-cloud reconstruction and multi-view rendering, meaning that SG introduces no additional inference latency. Combining the performance in Tab. 1 and Tab. 2 with the efficiency comparison, our method achieves the best trade-off between performance and efficiency, demonstrating strong potential for real-world applications.

**Out-Of-Distribution Dimensions** We evaluate out-of-distribution (OOD) performance on MultiScan [42] and ArtiScenes [4], comparing $S^2$-MLLM with a traditional supervised method MCLN [51] and a zero-shot method See-Ground [34]. We train our model on the union datasets of ScanRefer [9] and ReferIt3D [2], which are built on Scan-Net [14]. Then we directly evaluate the OOD performance of our model without any fine-tuning. Compared with Scan-Net [14], MultiScan [42] and ArtiScenes [4] introduce substantial distribution shifts in scene layouts, object compositions, and language descriptions. As shown in Table 4, our $S^2$-MLLM achieves the best performance under both OOD benchmarks with Acc@25 of 59.13 on Multiscan [42] and 43.26 on ArkiScenes [4]. The results show that our model acquires spatial reasoning ability through spatial guidance rather than learning the distribution of a specific dataset.

### 4.6. Ablation

**Ablation on Our Modules** We validate the effectiveness of each component in $S^2$-MLLM. As shown in Tab. 5, all

Table 5. **Ablation study** on the ScanRefer [9] dataset. We evaluate the contribution of each proposed component and the impact of the number of input frames. (SG) Spatial Guidance; (MPE) Multi-level Position Encoding; (Attn) Intra-view and Inter-view Attention; (LG) Language Guidance.

| Ablation | Num Frames | Overall@0.25 | Overall@0.5 |
|---|---|---|---|
| w/o SG | 16 | 54.40 | 48.45 |
| | 24 | 56.46 | 49.88 |
| w/o MPE | 16 | 44.13 | 38.49 |
| w/o Attn | 16 | 59.13 | 52.30 |
| w/o LG | 16 | 57.75 | 50.85 |
| Full ($S^2$-MLLM) | 16 | **59.18** | **52.67** |
| | 24 | **60.59** | **53.66** |

components bring clear performance gains. Removing spatial guidance (SG) leads to a clear drop of $4.78\%$ Acc@0.25 and $4.22\%$ Acc@0.5 under the 16 frames setting, demonstrating the importance of enforcing structure supervision. Removing multi-level positional encoding (MPE) causes the largest degradation of $15.05\%$ at Acc@0.25, indicating that explicit position information is crucial for LLM-based 3DVG.

**Ablation on Frames.** We further conduct ablation studies on the effect of the number of input frames. Increasing frames from 16 to 24 consistently improves performance across all settings. The gains are more pronounced in the setting without spatial guidance (SG), reaching $+2.06\%$ at Acc@0.25 and $+1.43\%$ at Acc@0.5. In contrast, the improvement becomes marginal with spatial guidance enabled ($+1.41\%$ at Acc@0.25 and $+0.99\%$ at Acc@0.5), while GPU memory usage and training time increase significantly. This comparison demonstrates that spatial guidance already enables the model to extract reliable structure cues and perform spatial reasoning even from sparse observations, reducing the dependence on dense multi-view inputs. This highlights the effectiveness of SG in enhancing spatial understanding while keeping computational overhead low.

## 5. Conclusion

In this work, we propose $S^2$-MLLM, a novel framework that equips MLLMs with implicit 3D reasoning capability for 3D visual grounding. By integrating feed-forward reconstruction to provide spatial guidance and introducing a structure-enhanced module, our model is enable to understand 3D scenes and perform latent spatial reasoning without requiring explicit point-cloud reconstruction at inference. Extensive experiments on both in-domain and out-of-domain benchmarks evaluate the performance of $S^2$-MLLM. The results verify that $S^2$-MLLM achieves outstanding accuracy, efficiency, and generalization, demonstrating its potential for real-world embodied applications.

# References

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 1

[2] Panos Achlioptas, Ahmed Abdelreheem, Fei Xia, Mohamed Elhoseiny, and Leonidas Guibas. Referit3d: Neural listeners for fine-grained 3d object identification in real-world scenes. In *European conference on computer vision*, pages 422–440. Springer, 2020. 1, 2, 5, 6, 8, 3

[3] Eslam Bakr, Yasmeen Alsaedy, and Mohamed Elhoseiny. Look around and refer: 2d synthetic semantics knowledge distillation for 3d visual grounding. *Advances in neural information processing systems*, 35:37146–37158, 2022. 2

[4] Gilad Baruch, Zhuoyuan Chen, Afshin Dehghan, Tal Dimry, Yuri Feigin, Peter Fu, Thomas Gebauer, Brandon Joffe, Daniel Kurz, Arik Schwartz, et al. Arkitscenes: A diverse real-world dataset for 3d indoor scene understanding using mobile rgb-d data. *arXiv preprint arXiv:2111.08897*, 2021. 2, 8, 1, 5

[5] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *Icml*, page 4, 2021. 4

[6] Daigang Cai, Lichen Zhao, Jing Zhang, Lu Sheng, and Dong Xu. 3djcg: A unified framework for joint dense captioning and visual grounding on 3d point clouds. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16464–16473, 2022. 1, 2

[7] Shivam Chandhok. Scenegpt: A language model for 3d scene understanding. *arXiv preprint arXiv:2408.06926*, 2024. 3

[8] Chun-Peng Chang, Shaoxiang Wang, Alain Pagani, and Didier Stricker. Mikasa: Multi-key-anchor & scene-aware transformer for 3d visual grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14131–14140, 2024. 2

[9] Dave Zhenyu Chen, Angel X Chang, and Matthias Nießner. Scanrefer: 3d object localization in rgb-d scans using natural language. In *European conference on computer vision*, pages 202–221. Springer, 2020. 1, 2, 5, 6, 7, 8, 3, 4

[10] Dave Zhenyu Chen, Qirui Wu, Matthias Nießner, and Angel X Chang. D 3 net: A unified speaker-listener architecture for 3d dense captioning and visual grounding. In *European Conference on Computer Vision*, pages 487–505. Springer, 2022.

[11] Shizhe Chen, Pierre-Louis Guhur, Makarand Tapaswi, Cordelia Schmid, and Ivan Laptev. Language conditioned spatial relation reasoning for 3d object grounding. *Advances in neural information processing systems*, 35:20522–20535, 2022. 1, 2, 6

[12] Sijin Chen, Xin Chen, Chi Zhang, Mingsheng Li, Gang Yu, Hao Fei, Hongyuan Zhu, Jiayuan Fan, and Tao Chen. Ll3da: Visual interactive instruction tuning for omni-3d understanding reasoning and planning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 26428–26438, 2024. 3

[13] Yilun Chen, Shuai Yang, Haifeng Huang, Tai Wang, Runsen Xu, Ruiyuan Lyu, Dahua Lin, and Jiangmiao Pang. Grounded 3d-llm with referent tokens. *arXiv preprint arXiv:2405.10370*, 2024. 3

[14] 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. 5, 8, 2

[15] Rao Fu, Jingyu Liu, Xilun Chen, Yixin Nie, and Wenhan Xiong. Scene-llm: Extending language model for 3d visual understanding and reasoning. *arXiv preprint arXiv:2403.11401*, 2024. 3

[16] Wenxuan Guo, Xiuwei Xu, Ziwei Wang, Jianjiang Feng, Jie Zhou, and Jiwen Lu. Tsp3d: Text-guided sparse voxel pruning for efficient 3d visual grounding. *arXiv preprint arXiv:2502.10392*, 2025. 1, 2, 6

[17] Zoey Guo, Yiwen Tang, Ray Zhang, Dong Wang, Zhigang Wang, Bin Zhao, and Xuelong Li. Viewrefer: Grasp the multi-view knowledge for 3d visual grounding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15372–15383, 2023. 1, 2

[18] Ziyu Guo, Renrui Zhang, Xiangyang Zhu, Yiwen Tang, Xianzheng Ma, Jiaming Han, Kexin Chen, Peng Gao, Xianzhi Li, Hongsheng Li, et al. Point-bind & point-llm: Aligning point cloud with multi-modality for 3d understanding, generation, and instruction following. *arXiv preprint arXiv:2309.00615*, 2023. 3

[19] 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. 1, 3, 5

[20] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022. 6, 1

[21] Jiangyong Huang, Silong Yong, Xiaojian Ma, Xiongkun Linghu, Puhao Li, Yan Wang, Qing Li, Song-Chun Zhu, Baoxiong Jia, and Siyuan Huang. An embodied generalist agent in 3d world. *arXiv preprint arXiv:2311.12871*, 2023. 3

[22] Ronggang Huang, Haoxin Yang, Yan Cai, Xuemiao Xu, Huaidong Zhang, and Shengfeng He. Viewsrd: 3d visual grounding via structured multi-view decomposition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9726–9736, 2025. 6

[23] Shijia Huang, Yilun Chen, Jiaya Jia, and Liwei Wang. Multi-view transformer for 3d visual grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15524–15533, 2022. 1, 2, 5

[24] 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. 1, 2, 6

[25] Ayush Jain, Alexander Swerdlow, Yuzhou Wang, Sergio Arnaud, Ada Martin, Alexander Sax, Franziska Meier, and Katerina Fragkiadaki. Unifying 2d and 3d vision-language understanding. *arXiv preprint arXiv:2503.10745*, 2025. 1

[26] Baoxiong Jia, Yixin Chen, Huangyue Yu, Yan Wang, Xuesong Niu, Tengyu Liu, Qing Li, and Siyuan Huang. Sceneverse: Scaling 3d vision-language learning for grounded scene understanding. In *European Conference on Computer Vision*, pages 289–310. Springer, 2024. 1

[27] Li Jiang, Hengshuang Zhao, Shaoshuai Shi, Shu Liu, Chi-Wing Fu, and Jiaya Jia. Pointgroup: Dual-set point grouping for 3d instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and Pattern recognition*, pages 4867–4876, 2020. 2

[28] Zhao Jin, Munawar Hayat, Yuwei Yang, Yulan Guo, and Yinjie Lei. Context-aware alignment and mutual masking for 3d-language pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10984–10994, 2023. 1

[29] Zhao Jin, Rong-Cheng Tu, Jingyi Liao, Wenhao Sun, Xiao Luo, Shunyu Liu, and Dacheng Tao. Spazer: Spatial-semantic progressive reasoning agent for zero-shot 3d visual grounding. *arXiv preprint arXiv:2506.21924*, 2025. 1, 2, 3

[30] 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. 2, 6

[31] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024. 1

[32] Junxian Li, Xinyue Xu, Sai Ma, and Sichao Li. Faithact: Faithfulness planning and acting in mllms. *arXiv preprint arXiv:2511.08409*, 2025.

[33] Junxian Li, Di Zhang, Xunzhi Wang, Zeying Hao, Jingdi Lei, Qian Tan, Cai Zhou, Wei Liu, Yaotian Yang, Xinrui Xiong, et al. Chemvlm: Exploring the power of multimodal large language models in chemistry area. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 415–423, 2025. 1

[34] Rong Li, Shijie Li, Lingdong Kong, Xulei Yang, and Jun-wei Liang. Seeground: See and ground for zero-shot open-vocabulary 3d visual grounding. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 3707–3717, 2025. 1, 2, 3, 5, 6, 7, 8

[35] Daizong Liu, Yang Liu, Wencan Huang, and Wei Hu. A survey on text-guided 3d visual grounding: Elements, recent advances, and future directions. *arXiv preprint arXiv:2406.05785*, 2024. 1

[36] Daizong Liu, Yang Liu, Wencan Huang, and Wei Hu. A survey on text-guided 3-d visual grounding: Elements, recent advances, and future directions. *IEEE Transactions on Neural Networks and Learning Systems*, 2025. 1

[37] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916, 2023. 1

[38] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *Advances in Neural Information Processing Systems (NeurIPS) 33*, 2020. 2, 3

[39] Ziyang Lu, Yunqiang Pei, Guoqing Wang, Peiwei Li, Yang Yang, Yinjie Lei, and Heng Tao Shen. Scaneru: Interactive 3d visual grounding based on embodied reference understanding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 3936–3944, 2024. 1

[40] Junyu Luo, Jiahui Fu, Xianghao Kong, Chen Gao, Haibing Ren, Hao Shen, Huaxia Xia, and Si Liu. 3d-sps: Single-stage 3d visual grounding via referred point progressive selection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16454–16463, 2022. 2

[41] Xianzheng Ma, Yash Bhalgat, Brandon Smart, Shuai Chen, Xinghui Li, Jian Ding, Jindong Gu, Dave Zhenyu Chen, Songyou Peng, Jia-Wang Bian, et al. When llms step into the 3d world: A survey and meta-analysis of 3d tasks via multi-modal large language models. *arXiv preprint arXiv:2405.10255*, 2024. 1, 3

[42] Yongsen Mao, Yiming Zhang, Hanxiao Jiang, Angel Chang, and Manolis Savva. Multiscan: Scalable rgbd scanning for 3d environments with articulated objects. *Advances in neural information processing systems*, 35:9058–9071, 2022. 2, 8, 1, 5

[43] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 4

[44] Taiki Miyanishi, Daichi Azuma, Shuhei Kurita, and Motoaki Kawanabe. Cross3dvg: Cross-dataset 3d visual grounding on different rgb-d scans. In *2024 International Conference on 3D Vision (3DV)*, pages 717–727. IEEE, 2024. 2

[45] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 5

[46] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022. 6

[47] Qihang Peng, Henry Zheng, and Gao Huang. Proxytransformation: Preshaping point cloud manifold with proxy attention for 3d visual grounding. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 24582–24592, 2025. 1, 2

[48] Songyou Peng, Kyle Genova, Chiyu Jiang, Andrea Tagliasacchi, Marc Pollefeys, Thomas Funkhouser, et al. Openscene: 3d scene understanding with open vocabularies. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 815–824, 2023. 6

[49] Charles R Qi, Or Litany, Kaiming He, and Leonidas J Guibas. Deep hough voting for 3d object detection in point

clouds. In *proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9277–9286, 2019. 2

[50] Zhangyang Qi, Zhixiong Zhang, Ye Fang, Jiaqi Wang, and Hengshuang Zhao. Gpt4scene: Understand 3d scenes from videos with vision-language models. *arXiv preprint arXiv:2501.01428*, 2025. 1, 2, 3

[51] 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, 2024. 6, 8

[52] 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

[53] Jonas Schult, Francis Engelmann, Alexander Hermans, Or Litany, Siyu Tang, and Bastian Leibe. Mask3d: Mask transformer for 3d semantic instance segmentation. *arXiv preprint arXiv:2210.03105*, 2022. 2, 6

[54] Yichen Sha, Siting Zhu, Hekui Guo, Zhong Wang, and Hesheng Wang. Towards autonomous indoor parking: A globally consistent semantic slam system and a semantic localization subsystem. In *2025 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 21151–21157, 2025. 1

[55] Xiangxi Shi, Zhonghua Wu, and Stefan Lee. Aware visual grounding in 3d scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14056–14065, 2024. 6

[56] Shubham Tulsiani, Richard Tucker, and Noah Snavely. Layer-structured 3d scene inference via view synthesis. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018. 2, 3

[57] 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, 2024. 6

[58] Haochen Wang, Yucheng Zhao, Tiancai Wang, Haoqiang Fan, Xiangyu Zhang, and Zhaoxiang Zhang. Ross3d: Reconstructive visual instruction tuning with 3d-awareness. *arXiv preprint arXiv:2504.01901*, 2025. 3

[59] Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea Vedaldi, Christian Rupprecht, and David Novotny. Vggt: Visual geometry grounded transformer. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 5294–5306, 2025. 2, 3

[60] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. 1, 6, 8

[61] Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris Chidlovskii, and Jerome Revaud. Dust3r: Geometric 3d vision made easy. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20697–20709, 2024. 2, 3

[62] Xinyi Wang, Na Zhao, Zhiyuan Han, Dan Guo, and Xun Yang. Augrefer: Advancing 3d visual grounding via cross-modal augmentation and spatial relation-based referring. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 8006–8014, 2025. 6

[63] Yuan Wang, Yali Li, and Shengjin Wang. G^ 3-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. 6

[64] Yuan Wang, Ya-Li Li, WU Eastman ZY, and Shengjin Wang. Liba: Language instructed multi-granularity bridge assistant for 3d visual grounding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 8114–8122, 2025. 1, 2

[65] Zehan Wang, Haifeng Huang, Yang Zhao, Linjun Li, Xize Cheng, Yichen Zhu, Aoxiong Yin, and Zhou Zhao. 3drp-net: 3d relative position-aware network for 3d visual grounding. *arXiv preprint arXiv:2307.13363*, 2023. 1

[66] 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

[67] Zehan Wang, Haifeng Huang, Yang Zhao, Ziang Zhang, and Zhou Zhao. Chat-3d: Data-efficiently tuning large language model for universal dialogue of 3d scenes. *arXiv preprint arXiv:2308.08769*, 2023. 3, 5

[68] Diankun Wu, Fangfu Liu, Yi-Hsin Hung, and Yueqi Duan. Spatial-mllm: Boosting mllm capabilities in visual-based spatial intelligence. *arXiv preprint arXiv:2505.23747*, 2025. 2, 3

[69] Jiayang Wu, Wensheng Gan, Zefeng Chen, Shicheng Wan, and Philip S Yu. Multimodal large language models: A survey. In *2023 IEEE International Conference on Big Data (BigData)*, pages 2247–2256. IEEE, 2023. 3

[70] 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. 1, 2, 6

[71] Beining Xu, Siting Zhu, and Hesheng Wang. Sgloc: Semantic localization system for camera pose estimation from 3d gaussian splatting representation. In *2025 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 8872–8878, 2025. 1

[72] Runsen Xu, Zhiwei Huang, Tai Wang, Yilun Chen, Jiangmiao Pang, and Dahua Lin. Vlm-grounder: A vlm agent for zero-shot 3d visual grounding. *arXiv preprint arXiv:2410.13860*, 2024. 1, 2, 3

[73] Jianing Yang, Xuweiyi Chen, Shengyi Qian, Nikhil Madaan, Madhavan Iyengar, David F Fouhey, and Joyce Chai. Llm-grounder: 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

[74] Jianing Yang, Alexander Sax, Kevin J Liang, Mikael Henaff, Hao Tang, Ang Cao, Joyce Chai, Franziska Meier, and Matt Feiszli. Fast3r: Towards 3d reconstruction of 1000+ images in one forward pass. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 21924–21935, 2025. 2, 3

[75] Li Yang, Ziqi Zhang, Zhongang Qi, Yan Xu, Wei Liu, Ying Shan, Bing Li, Weiping Yang, Peng Li, Yan Wang, et al. Exploiting contextual objects and relations for 3d visual grounding. *Advances in Neural Information Processing Systems*, 36:49542–49554, 2023. 1

[76] 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*, pages 1856–1866, 2021. 2

[77] Chun-Hsiao Yeh, Chenyu Wang, Shengbang Tong, Ta-Ying Cheng, Ruoyu Wang, Tianzhe Chu, Yuexiang Zhai, Yubei Chen, Shenghua Gao, and Yi Ma. Seeing from another perspective: Evaluating multi-view understanding in mllms. *arXiv preprint arXiv:2504.15280*, 2025. 4

[78] Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey on multimodal large language models. *National Science Review*, 11(12): nwae403, 2024. 3

[79] Hanxun Yu, Wentong Li, Song Wang, Junbo Chen, and Jianke Zhu. Inst3d-lmm: Instance-aware 3d scene understanding with multi-modal instruction tuning. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 14147–14157, 2025. 3

[80] 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. 1, 2, 6

[81] 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

[82] Nader Zantout, Haochen Zhang, Pujith Kachana, Jinkai Qiu, Guofei Chen, Ji Zhang, and Wenshan Wang. Sort3d: Spatial object-centric reasoning toolbox for zero-shot 3d grounding using large language models. *arXiv preprint arXiv:2504.18684*, 2025. 2

[83] Jirong Zha, Yuxuan Fan, Xiao Yang, Chen Gao, and Xinlei Chen. How to enable llm with 3d capacity? a survey of spatial reasoning in llm. *arXiv preprint arXiv:2504.05786*, 2025. 3

[84] Xiaoyu Zhan, Wenxuan Huang, Hao Sun, Xinyu Fu, Changfeng Ma, Shaosheng Cao, Bohan Jia, Shaohui Lin, Zhenfei Yin, Lei Bai, et al. Actial: Activate spatial reasoning ability of multimodal large language models. *arXiv preprint arXiv:2511.01618*, 2025. 1

[85] Haomeng Zhang, Chiao-An Yang, and Raymond A Yeh. Multi-object 3d grounding with dynamic modules and language-informed spatial attention. *Advances in Neural Information Processing Systems*, 37:123237–123260, 2024. 2

[86] Jiawei Zhang, Chejian Xu, and Bo Li. Chatscene: Knowledge-enabled safety-critical scenario generation for autonomous vehicles. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15459–15469, 2024. 3

[87] Jiahui Zhang, Yuelei Li, Anpei Chen, Muyu Xu, Kunhao Liu, Jianyuan Wang, Xiao-Xiao Long, Hanxue Liang, Zexiang Xu, Hao Su, et al. Advances in feed-forward 3d reconstruction and view synthesis: A survey. *arXiv preprint arXiv:2507.14501*, 2025. 2

[88] Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. Video instruction tuning with synthetic data. *arXiv preprint arXiv:2410.02713*, 2024. 3, 6, 1

[89] Lichen Zhao, Daigang Cai, Lu Sheng, and Dong Xu. 3dvg-transformer: Relation modeling for visual grounding on point clouds. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2928–2937, 2021. 1, 2, 6

[90] Duo Zheng, Shijia Huang, Yanyang Li, and Liwei Wang. Learning from videos for 3d world: Enhancing mllms with 3d vision geometry priors. *arXiv preprint arXiv:2505.24625*, 2025. 3

[91] Duo Zheng, Shijia Huang, and Liwei Wang. Video-3d llm: Learning position-aware video representation for 3d scene understanding. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 8995–9006, 2025. 1, 3, 5, 6, 7, 8

[92] Hongyan Zhi, Peihao Chen, Junyan Li, Shuailei Ma, Xinyu Sun, Tianhang Xiang, Yinjie Lei, Mingkui Tan, and Chuang Gan. Lscenellm: Enhancing large 3d scene understanding using adaptive visual preferences. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 3761–3771, 2025. 3

[93] 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. 5

[94] Chenming Zhu, Tai Wang, Wenwei Zhang, Jiangmiao Pang, and Xihui Liu. Llava-3d: A simple yet effective pathway to empowering lmms with 3d-awareness. *arXiv preprint arXiv:2409.18125*, 2024. 3

[95] Siting Zhu, Guangming Wang, Hermann Blum, Jiuming Liu, Liang Song, Marc Pollefeys, and Hesheng Wang. Sni-slam: Semantic neural implicit slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21167–21177, 2024. 1

[96] Siting Zhu, Guangming Wang, Xin Kong, Dezhi Kong, and Hesheng Wang. 3d gaussian splatting in robotics: A survey. *arXiv preprint arXiv:2410.12262*, 2024. 1

[97] Siting Zhu, Renjie Qin, Guangming Wang, Jiuming Liu, and Hesheng Wang. Semgauss-slam: Dense semantic gaussian splatting slam. In *2025 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 21174–21181, 2025. 1

[98] Siting Zhu, Guangming Wang, Hermann Blum, Zhong Wang, Ganlin Zhang, Daniel Cremers, Marc Pollefeys, and Hesheng Wang. Sni-slam++: Tightly-coupled semantic neural implicit slam. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025. 1

[99] 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. 3, 5, 6, 1

[100] Ziyu Zhu, Zhuofan Zhang, Xiaojian Ma, Xuesong Niu, Yixin Chen, Baoxiong Jia, Zhidong Deng, Siyuan Huang, and Qing Li. Unifying 3d vision-language understanding via promptable queries. In *European Conference on Computer Vision*, pages 188–206. Springer, 2024. 2

# S²-MLLM: Boosting Spatial Reasoning Capability of MLLMs for 3D Visual Grounding with Structural Guidance

## Supplementary Material

## A. Implementation Details

### A.1. Training Details

During training, we freeze the reconstruction decoder and fine-tune the projection layer, the visual encoder, and the language model. Since the LLM output resides in the text space, whereas 3D visual grounding requires regression to region-level features, using LoRA [20] alone is insufficient for learning cross-modal alignment. Therefore, we additionally fully finetune the last four layers of the LLM, the projection layer, as well as the inter-view and intra-view attention during training. The rest of the model is fine-tuned using LoRA [20], including the visual encoder.

### A.2. Dataset and Arguments

Following [91, 99], we train our model on the combined dataset of ScanRefer [9], Nr3D [2], and Sr3D [2], and evaluate it separately on each corresponding validation set. We provide detailed statistics about the data used for training and evaluation, including out-of-distribution (OOD) evaluation in Tab. 6. Following [91], we convert all data into the LLaVA [37] format and report statistics based on this unified format.

In addition, we show detailed hyperparameters in our experiments in Tab. 7

Table 6. Detailed statistics of datasets.

| Split | Dataset | Samples | Scenes | Query Length |
|-------|---------|---------|--------|--------------|
| Train | ScanRefer [9] | 36665 | 562 | 17.83 |
|       | Nr3D [2] | 32919 | 511 | 25.38 |
|       | Sr3D [2] | 65844 | 1018 | 23.68 |
| Test | ScanRefer [9] | 9508 | 141 | 17.92 |
|      | Nr3D [2] | 8584 | 130 | 25.12 |
|      | Sr3D [2] | 17726 | 255 | 23.72 |
| OOD | MultiScan [42] | 1490 | 53 | 20.44 |
|     | ArkitScenes [4] | 2693 | 275 | 20.88 |

Table 7. Detailed statistics of hyperparameters.

| Parameter Name | Value |
|----------------|-------|
| **Training** | |
| LoRA rank | 64 |
| LoRA $\alpha$ | 16 |
| LoRA dropout | 0.05 |
| LoRA bias | None |
| FP/BF Precision | bf16 |
| tf32 | False |
| weight Decay | 0.0 |
| tuning MLP or ViT | True |
| training steps | 16928 |
| batch size | 1 |
| warmup ratio | 0.05 |
| lr | 2e-5 |
| lr of ViT | 2e-4 |
| optimizer | AdamW |
| max token length | 32768 |
| gradient accumulation steps | 8 |
| $\lambda_g$ | 1.0 |
| $\lambda_r$ | 0.3 |
| $\lambda_l$ | 1.0 |
| **Inference** | |
| temperature | 1.0 |
| num beams | 1 |
| top_p, top_k | 1.0, 50 |

Table 8. **Ablation study** on the ScanRefer [9] dataset. We evaluate the contribution of inter-view and intra-view attention (Attn).

| Ablation | Overall@0.25 | Overall@0.5 |
|----------|--------------|-------------|
| **Base** | 5.31 | 5.07 |
| **Base+Attn** | **41.74** | **35.78** |

## B. Additional Ablation Studies

To further verify the contribution of the inter-view and intra-view attention (Attn), we provide an additional ablation experiment. Specifically, we report results for LLava-Video 7B [88] (Base) and Base + Attn, which adds only Attn. This comparison enables a clearer evaluation of the effect of Attn, isolated from the influence of other components.

As shown in Tab. 8, the inclusion of Attn significantly improves the performance of S²-MLLM in 3DVG. Since other modules we proposed are highly effective and contribute substantially to the overall performance, removing the Attn module in the ablation study mentioned in Tab. 5 of the main text does not result in a significant performance drop.
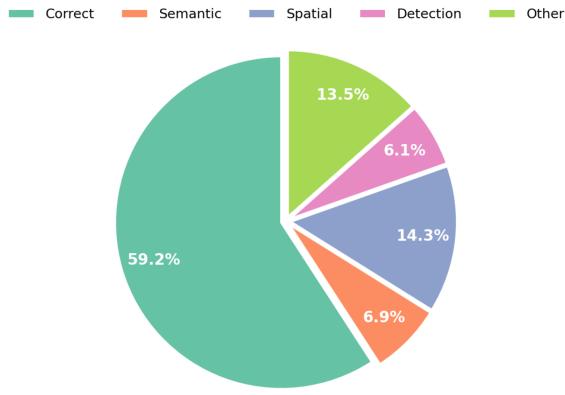
Figure 4. Error type analysis on ScanRefer [9] dataset.

## C. Additional Analysis

We analyzed the error types in the predictions of $S^2$-MLLM in Fig. 4 and visualize typical error cases in Fig. 5. We adopted the same definition of error types as [29], classifying errors into four types. Overall, $S^2$-MLLM achieves an accuracy of $59.2\%$ on ScanRefer [9], indicating that our model can understand language descriptions and 3D scenes in most cases.

Among the errors, we observed that: (1) Spatial: As shown in Fig. 5, these cases mainly involve complex relational descriptions with multiple anchor objects (*e.g.*, determining the target object based on multiple anchor objects or complex relational descriptions that require multi-step reasoning). (2) Semantic: These errors mainly occur due to misjudgment of fine-grained attributes of the target object and anchor objects (*e.g.*, the color of the keyboard, the pattern and usage of the curtain). (3) Detection: These errors arise because the bounding boxes provided by the detector are inaccurate. In these cases, detectors usually predict the correct object but predict an inaccurate 3D bounding region, often due to occlusion, partial visibility, or sparse viewpoints. This suggests that improving the robustness under limited view coverage could further enhance performance. (4) Other: This category primarily refers to inaccurate language descriptions, where the predicted object matches the description but differs from the ground truth. This indicates that current 3DVG datasets still have limitations and incompleteness, rather than being due to the model's performance.

## D. Additional Qualitative Results

### D.1. Visualizations on Nr3D

We additionally visualize the performance of $S^2$-MLLM and SeeGround [34] on Nr3D [2], further demonstrating the advantages of our approach. As shown in Fig. 6, $S^2$-MLLM is capable of accurately understanding the language descrip-

tions and the spatial relationships in the 3D scene, such as the relative size of two similar regions, the distance between two similar objects, and the window. This demonstrates that our model exhibits superior 3D spatial understanding and reasoning abilities.

### D.2. Ablation Visualizations

We provide additional visual examples on the ScanRefer [9] to supplement the ablation analysis, further illustrating the effectiveness of each component we proposed in $S^2$-MLLM. As shown in Fig. 7, each row corresponds to different ablation experiments. We observe that spatial guidance (SG) enables $S^2$-MLLM to accurately comprehend the spatial relationship (*e.g.*, the chair at the head of the table). The second row highlights the importance of LG: without language guidance (LG), $S^2$-MLLM fails to identify the brown couch correctly in the context of surrounding objects. With multi-level position encoding (MPE), encoding camera rays' viewing directions encourages $S^2$-MLLM to distinguish orientation and observation direction (*e.g.*, identifying the chair facing the desk). The fourth row shows the impact of inter-view and intra-view attention (Attn): $S^2$-MLLM without Attn struggles to locate the red and black office chair at the correct position relative to the table, while the full model successfully identifies it. In the presence of Attn, $S^2$-MLLM can stably identify the target object even when the query involves viewpoint transitions (*e.g.*, describing the target object from different perspectives or based on different anchor objects). These results demonstrate how each part we proposed contributes to improving the ability of $S^2$-MLLM to understand and reason in the 3D scenes.

### D.3. Out-of-Distribution Qualitative Examples

Fig. 8 and Fig. 9 provide additional qualitative examples on Multiscan [42] and Arkiscenes [4] compared with See-Groud [34]. As shown in Fig. 8 and Fig. 9, the scene layouts of MultiScan [42] and ArkiScenes [4] differ from those in ScanNet [14]. Despite being completely out of distribution, $S^2$-MLLM is still able to maintain strong spatial reasoning capabilities and achieves superior performance.

| **Other** | **Spatial** | **Semantic** | **Detection** |
|---|---|---|---|



Query: *There is **a square chair**. It is at the end of a long table.*

Query: *The **large door** to the room. The door is next to the trash can.*

Query: *This is a **white keyboard**. It is located below the computer screen on the right.*

Query: *There is **a square chair**. It is at the end of a long table.*

Query: ***A brown wooden chair.** It is under a wooden table.*

Query: *To the right of the double doors and on the floor is an orange pylon. Behind the orange pylon is a ladder and a vacuum cleaner. Behind the ladder and behind the vacuum cleaner is a stainless steel sink. **The stainless steel sink** is the item we are looking for*

Query: *The shower **curtain** is long and in an open position. It is golden in tone with a textured pattern.*

Query: ***This table** is not the only table in the room. This table is not the only table in the room it is the only cluttered table in the room. It is on the left side of the doors against the wall.*
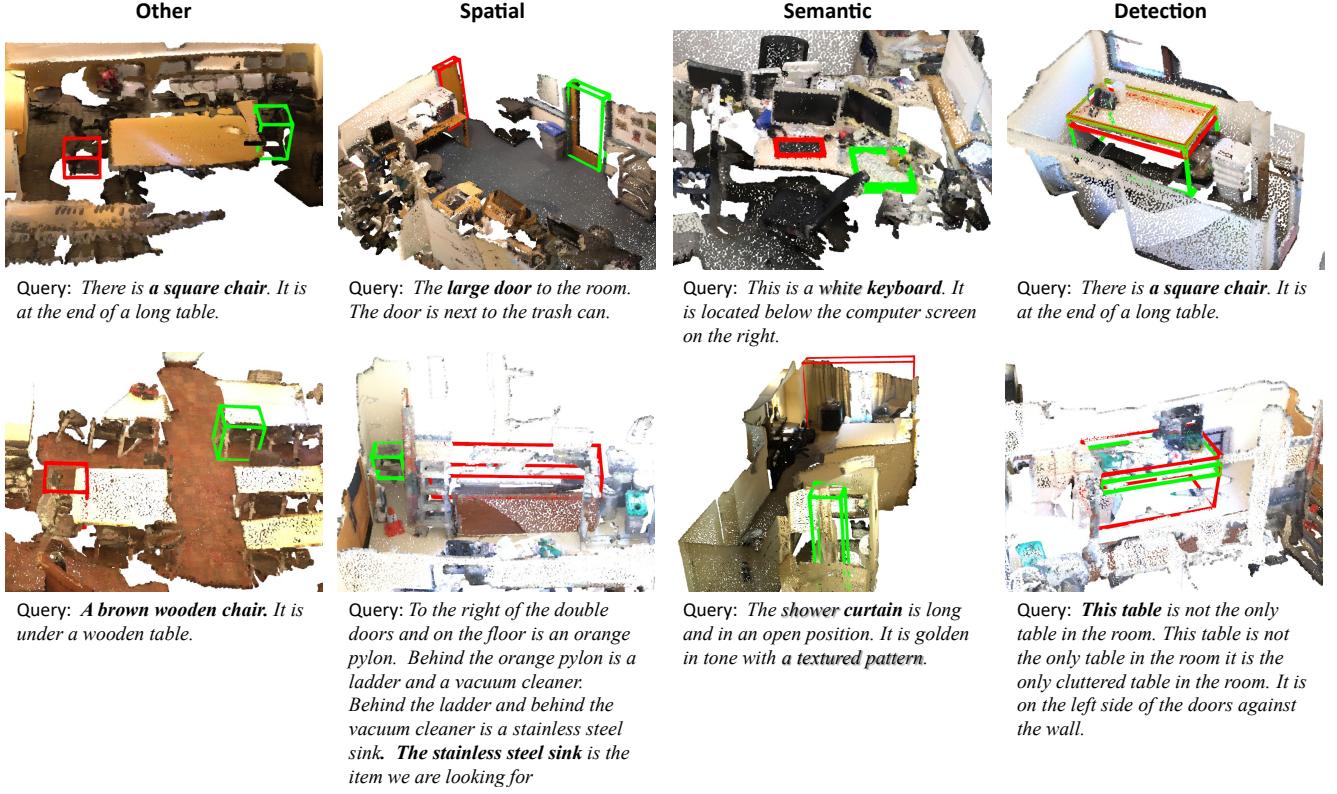
Figure 5. Typical types of failure cases in Scanrefer [9]. Ground Truth is highlighted in green, our predictions in red. Key semantic information is in shadow.

| **Comparison** | **Ground Truth** | **Comparison** | **Ground Truth** |
|---|---|---|---|



Query: *The rear, left hand side **pillow** on the bed that is closest to the cabinets.*

Query: *The narrower of the two **stalls.***

Query: *The **lamp** next to the window.*

Query: *The small **blue couch** that is on the end.*

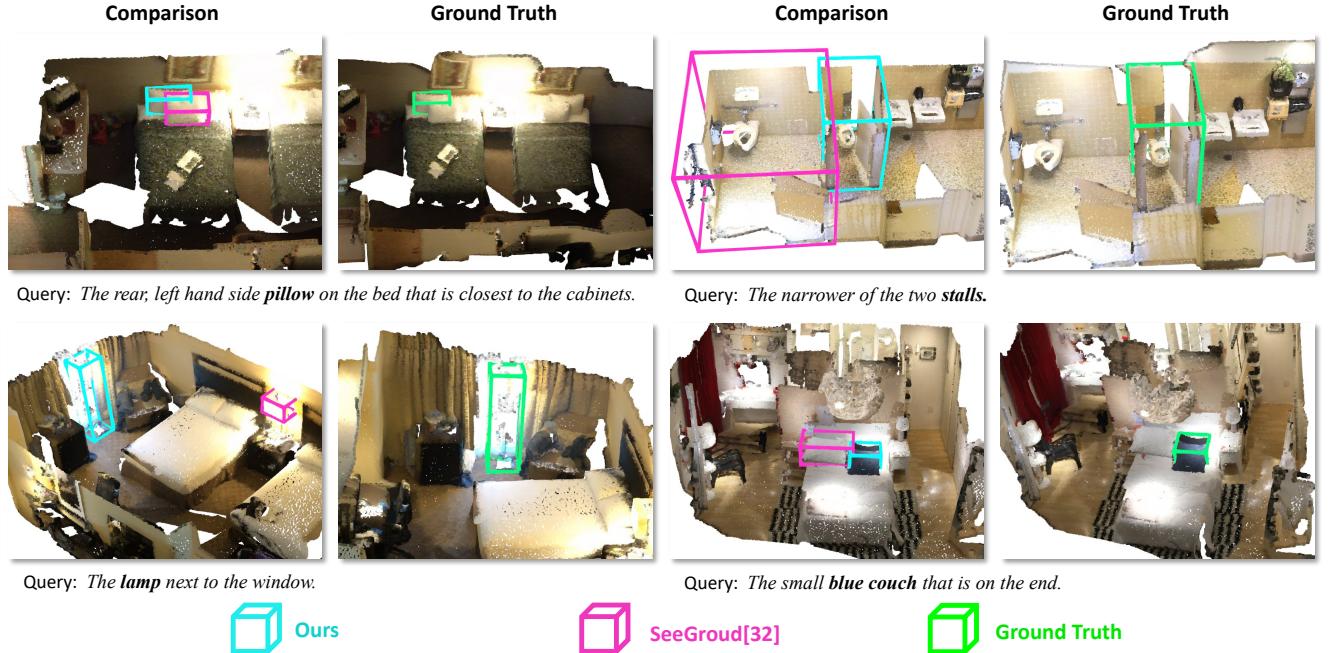Ours · · · SeeGroud[32] · · · Ground Truth

Figure 6. Qualitative comparison of 3DVG results in Nr3D [2]. Ground Truth is highlighted in green, our predictions in cyan, and predictions of SeeGround [34] in magenta.
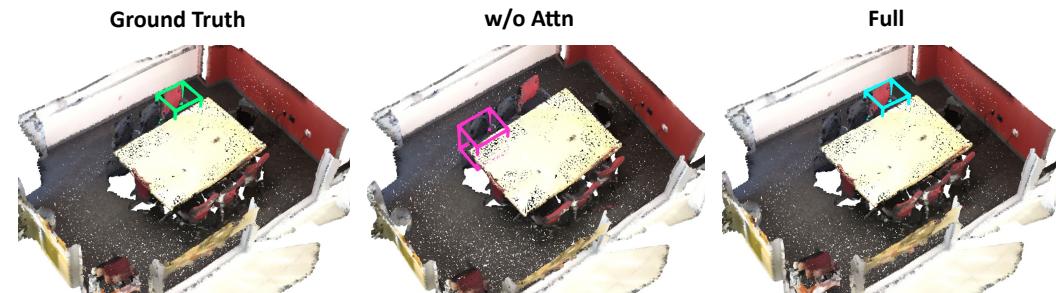
3

**Ground Truth**  **w/o SG**  **Full**

Query*: You are looking for **the chair** at the head of the table. it's near the door.*

**Ground Truth**  **w/o LG**  **Full**

Query*: **The brown couch** is in between the piano and the glass double doors. **The brown couch** is against the white wall.*

**Ground Truth**  **w/o MPE**  **Full**

Query: **The chair** is facing the left corner of the room, and is at the desk. the monitor is to the left of the chair, and there are shelves in front of the chair.

**Ground Truth**  **w/o Attn**  **Full**

Query: **A red and black office chair** is sitting on the right side of the table. the chair is on the end in the front with the wall behind it.

Figure 7. Qualitative ablation results in Scanrefer [9]. Ground Truth is highlighted in mygreen, predictions of our full model are in cyan, and predictions of the model without specific module are in magenta. (SG) Spatial Guidance; (MPE) Multi-level Position Encoding; (Attn) Intra-view and Inter-view Attention; (LG) Language Guidance.

Query: *Located above two adjacent sinks, **the white windowsill** is used for placing things.*

Query: *Facing the door, there is **a folded wooden chair** near the railing on the left side of the door.*

Query: *The **three white, connected wall cabinets** above the oven.*

Query: *On the elliptical table closest to the tv, there is a **black backpack**.*

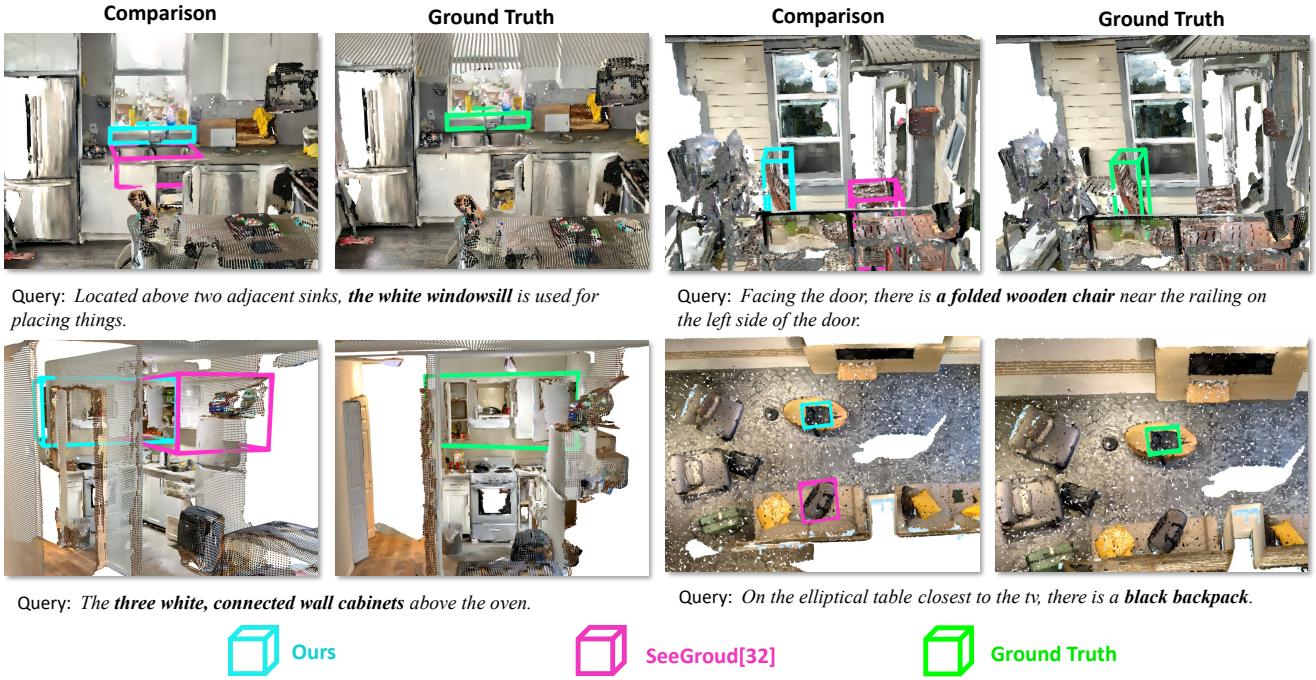**Ours**     **SeeGroud[32]**     **Ground Truth**

Figure 8. Qualitative comparison of 3DVG results in Multiscan [42]. Ground Truth is highlighted in green, our predictions in cyan, and predictions of SeeGround [34] in magenta.



Query: *In the bedroom, there is a **white drawer table** in the middle of the two white cabinets opposite the foot of the bed.*

Query: *Next to the entrance, facing the entrance, the upper half of the **brown cabinet** on the right side.*

Query: *In the kitchen, directly below the oven, there is **a rectangular cabinet support bar** used to support the oven.*

Query: ***The white cabinet** on the right side of the silver range hood on the black bar counter in the kitchen.*
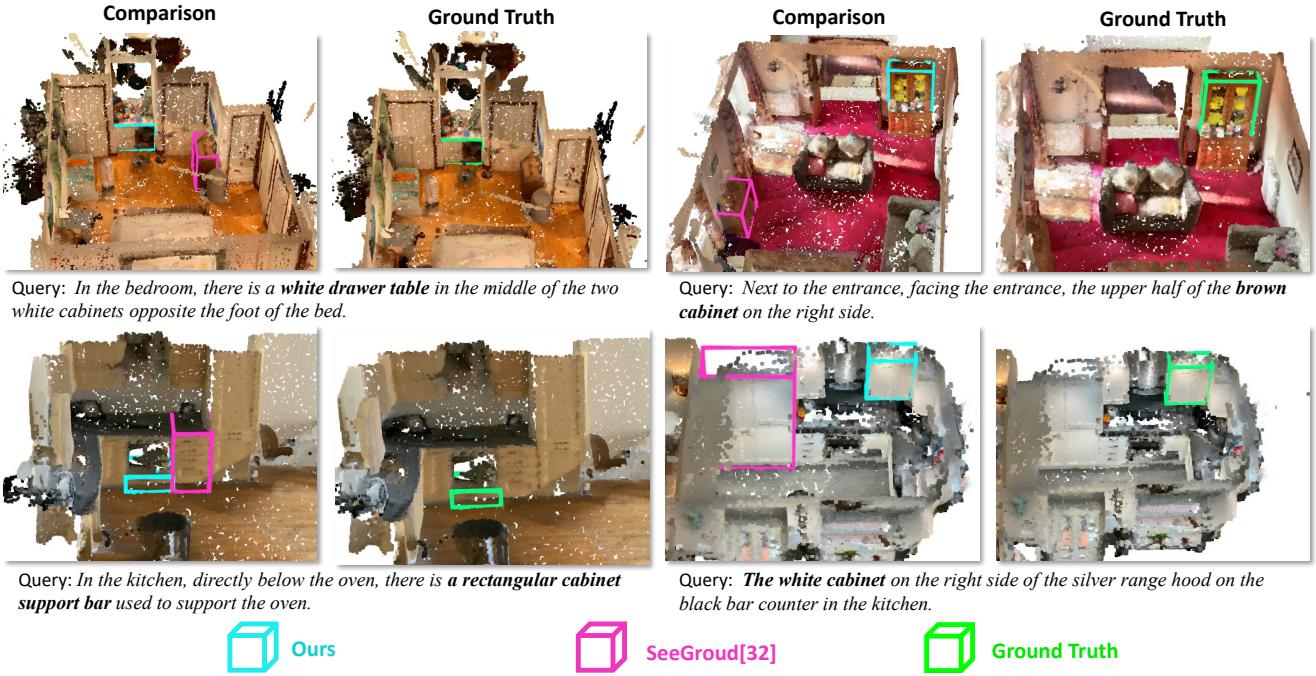
**Ours**     **SeeGroud[32]**     **Ground Truth**

Figure 9. Qualitative comparison of 3DVG results in ArkiScenes [4]. Ground Truth is highlighted in green, our predictions in cyan, and predictions of SeeGround [34] in magenta.