

ViewSRD: 3D Visual Grounding via Structured Multi-View Decomposition

Ronggang Huang^{1*}, Haoxin Yang^{1*†}, Yan Cai¹,
Xuemiao Xu^{12345†}, Huaidong Zhang¹, Shengfeng He⁶

¹ South China University of Technology ² Guangdong Engineering Center for Large Model and GenAI Technology

³ State Key Laboratory of Subtropical Building and Urban Science

⁴ Ministry of Education Key Laboratory of Big Data and Intelligent Robot

⁵ Guangdong Provincial Key Lab of Computational Intelligence and Cyberspace Information

⁶ Singapore Management University

Abstract

3D visual grounding aims to identify and localize objects in a 3D space based on textual descriptions. However, existing methods struggle with disentangling targets from anchors in complex multi-anchor queries and resolving inconsistencies in spatial descriptions caused by perspective variations. To tackle these challenges, we propose ViewSRD, a framework that formulates 3D visual grounding as a structured multi-view decomposition process. First, the Simple Relation Decoupling (SRD) module restructures complex multi-anchor queries into a set of targeted single-anchor statements, generating a structured set of perspective-aware descriptions that clarify positional relationships. These decomposed representations serve as the foundation for the Multi-view Textual-Scene Interaction (Multi-TSI) module, which integrates textual and scene features across multiple viewpoints using shared, Cross-modal Consistent View Tokens (CCVTs) to preserve spatial correlations. Finally, a Textual-Scene Reasoning module synthesizes multi-view predictions into a unified and robust 3D visual grounding. Experiments on 3D visual grounding datasets show that ViewSRD significantly outperforms state-of-the-art methods, particularly in complex queries requiring precise spatial differentiation. Code is available at <https://github.com/visualjason/ViewSRD>.

1. Introduction

3D Visual Grounding (3DVG) aims to establish semantic correspondences between natural language descriptions and target objects in a 3D space [19, 44]. This task has gained significant attention in applications such as visual language navigation [25, 57], intelligent agents [4, 49], and autonomous vehicles [9, 13].

*The first two authors contributed equally.

†Corresponding authors: xuemx@scut.edu.cn, harxis@outlook.com.

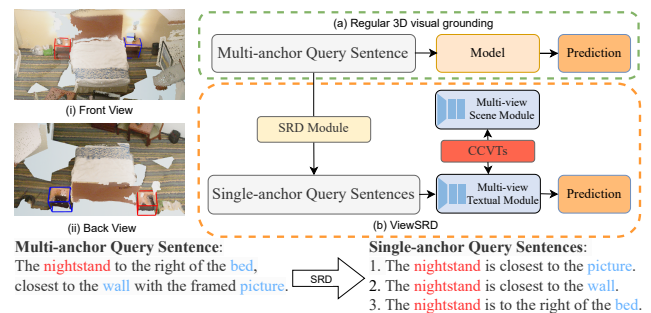


Figure 1. (a) Previous 3DVG methods struggle with ambiguities from complex multi-anchor queries and perspective shifts. (b) ViewSRD addresses this by using the SRD module to simplify queries and the CCVTs to capture viewpoint variations in both scene and textual modal, boosting cross-modal feature interaction and enhancing grounding accuracy.

Traditional single-view approaches rely on 2D sampled images to extract scene information [38] or construct scene graphs from textual descriptions [42]. However, these methods are inherently limited by their dependence on single-view cues, as language descriptions often presuppose specific viewpoints. To overcome this limitation, recent research has explored multi-view 3DVG, integrating multiple perspectives to enhance robustness [8, 31, 58]. Some methods process distinct descriptions for different viewpoints via manual annotation and learning [15, 37], while others incorporate spatial modules to encode relative spatial coordinates under specific perspectives [5]. However, they typically address only isolated aspects of the problem, limiting their effectiveness in handling complex multi-view scenarios.

Despite their potential, existing 3DVG models struggle to disentangle targets from anchors in multi-anchor textual descriptions [5, 14, 32]. Large language models (LLMs) often have difficulty interpreting such descriptions [17, 51], yet resolving these ambiguities is crucial for improving grounding accuracy [20]. Compounding this challenge, inconsisten-

cies between textual descriptions and spatial relationships arise when viewpoints change. As illustrated in Fig. 1, an object described as being to the right of another—such as “*The nightstand is to the right of the bed*”—from a front-facing view may appear on the left when observed from the opposite direction. These perspective-induced inconsistencies make it significantly harder for models to establish accurate correspondences between textual descriptions and visual information, further degrading performance. Ultimately, both the inherent complexity of multi-anchor queries and the challenges introduced by perspective shifts hinder the accurate interpretation of positional relationships in 3DVG, limiting the overall effectiveness of existing systems.

To tackle these challenges, we propose *ViewSRD*, a framework that formulates 3D visual grounding as a structured multi-view decomposition process. By leveraging the *Simple Relation Decoupling (SRD)* module, ViewSRD effectively disentangles target-anchor relationships in the complex multi-anchor queries, while the *Multi-view Textual-Scene Interaction (Multi-TSI)* module integrates multi-view information to enhance grounding accuracy. As illustrated in Fig. 1(b), ViewSRD first applies the SRD module to decompose complex multi-anchor queries into a set of simpler single-anchor queries, isolating interactions between the target and its anchors. This structured decomposition allows the model to more effectively learn positional relationships from textual descriptions. The Multi-TSI module then fuses textual and scene features across multiple viewpoints using *Cross-modal Consistent View Tokens (CCVTs)*, which explicitly encode viewpoint information as learnable cue for both textual and scene module. This mechanism ensures that the model accurately captures spatial interactions, even under perspective shifts. Finally, the *Textual-Scene Reasoning* module aggregates these multi-view features to accurately predict the final 3D VG results. Extensive experiments have validated the efficacy of our proposed ViewSRD across different 3DVG benchmarks, demonstrating its superior performance across diverse scenarios. In summary, our contributions are fourfold:

- We propose ViewSRD, a framework that formulates 3D visual grounding as a structured multi-view decomposition process, effectively handling complex multi-anchor queries and mitigating text-visual inconsistencies across different perspectives.
- We introduce the Simple Relation Decoupling (SRD) module, which restructures complex multi-anchor queries into simpler single-anchor statements, disentangling target-anchor relationships. This structured decomposition enables the model to extract more effective textual features for grounding.
- We develop the Multi-view Textual-Scene Interaction (Multi-TSI) module to explicitly encode viewpoint infor-

mation using cross-modal consistent view tokens. This mechanism ensures alignment between textual descriptions and visual features across different perspectives, reducing spatial ambiguities.

- We conduct extensive evaluations on 3D visual grounding datasets, where ViewSRD achieves state-of-the-art performance, yielding superior performance over prior work.

2. Related Work

3D Visual Grounding. 3D computer vision has made great progress in various fields [7, 26, 29, 33, 45, 55, 56], the 3D visual grounding (3DVG) task involves identifying a target object in a 3D scene based on a natural language description [19, 44]. Pioneering datasets such as ScanRefer [6] and ReferIt3D [1], built on ScanNet [10], have driven progress in this field. Recent advancements like MVT [18] address view inconsistency by developing a view-robust multi-modal representation. Other works [23, 27, 47, 53] explore multi-modal situated reasoning but lack a dedicated focus on handling the high semantic complexity of natural language in 3D grounding, particularly in disentangling intricate sentence structures.

Despite these advancements, the complexity of natural language descriptions remains a significant challenge in grounding tasks. Referring expressions often require reasoning over multiple anchor objects to precisely identify the target, making it crucial to disentangle and interpret intricate linguistic structures and spatial dependencies. Our method addresses this challenge by decoupling complex queries into simpler statements, improving the extraction of key relational information. Additionally, by leveraging view tokens, ViewSRD learns more accurate associations between textual descriptions and multi-view information.

Language Comprehension. Understanding referential language in 3DVG requires models to not only parse spatial descriptions but also interpret object relationships within a scene. Scene graphs, where objects serve as nodes and relationships form directed edges, have been widely used for tasks such as image retrieval and caption evaluation [19]. Traditional approaches employ scene graphs to enhance query comprehension [42], with efforts to convert sentences into structured representations [12, 34, 43] or generate grounded scene graphs for images [22, 24, 50]. However, these methods primarily focus on static, well-defined relationships and struggle with the dynamic, context-dependent nature of natural language. In datasets such as Nr3D [1], the complexity of interwoven spatial relationships and ambiguous references makes direct scene graph construction challenging. To address this, we propose leveraging Large Language Models (LLMs) [28, 40, 46] to enhance semantic understanding and spatial reasoning, reducing reliance on rigid structures while improving language comprehension.

3D Multi-View Learning. 3D vision research has largely

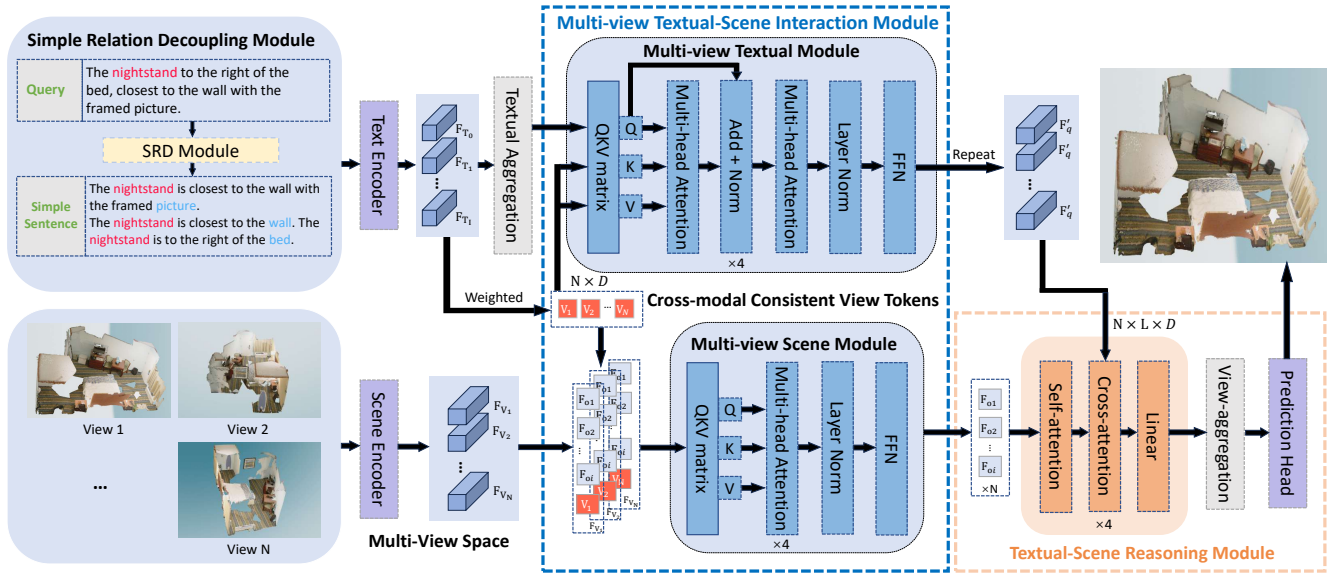


Figure 2. Overview of *ViewSRD*. We begin by employing the *Simple Relation Decoupling (SRD)* module to decompose complex multi-anchor queries into multiple simpler single-anchor queries. Next, text and scene features are extracted separately using the text encoder and scene encoder. To explicitly incorporate scene information into the model, we fuse *Cross-modal Consistent View Tokens (CCVTs)* with these extracted features. The *Multi-view Textual-Scene Interaction (Multi-TSI)* module then facilitates a comprehensive interaction between textual and scene information, the 3DVG prediction results are finally generated by the *Textual-Scene Reasoning Module*.

focused on generating 2D projections from multiple viewpoints. While LLM-based grounding methods integrate multi-view images, they struggle with accurately identifying the primary viewpoint and demonstrating reliability, as discussed in [17, 51]. MVT [18] maps 3D scenes into multiple perspectives to enhance cross-view feature aggregation but lacks a mechanism to weigh each view’s contribution, limiting performance in complex scenes. Similarly, ViewRefer [15] utilizes multi-view prototypes for cross-view interactions but lacks explicit training guidance on view importance. Mikasa [5] incorporates relative spatial coordinate information and a scene-aware module to improve object grounding but does not fully resolve view weighting challenges. In contrast, we propose Cross-modal Consistent View Tokens, which guide the model to dynamically adjust representation spaces and assess whether spatial relationships in decoupled sentences exhibit view dependency. This mechanism enables more reliable multi-view reasoning, improving performance in complex scenes.

3. ViewSRD

In the context of 3D point cloud scenes, the term *multi-view* refers to observing a shared scene representation (e.g., XYZ+RGB format) from different simulated viewpoints by rotating the scene around its central axis or camera viewpoints. Each view provides a partial observation of the same 3D environment, resulting in varying object appearances,

occlusions, and spatial configurations across views. This multi-view setup introduces significant challenges for 3D visual grounding: (1) language-grounded spatial relations must remain consistent across view-dependent variations, and (2) object referents may be partially or completely invisible in certain views.

To tackle these challenges, we propose ViewSRD, a structured multi-view 3D visual grounding framework. The overall framework of our method is illustrated in Fig. 2. ViewSRD comprises two key components. The first component is the *Simple Relation Decoupling (SRD)* module, which decomposes multi-anchor queries into a series of single-anchor queries by leveraging the powerful language processing capabilities of LLMs and predefined prompt templates. This decomposition enables more precise inference of relative relationships between objects, improving the model’s ability to capture spatial interactions. The second is the *Multi-view Textual-Scene Interaction (Multi-TSI)* module, which mitigates viewpoint dependency by integrating a shared, cross-modal consistent view token into both the language and visual models. These tokens facilitate feature interaction across perspectives, allowing the visual and textual models to align cross-modal viewpoint information more effectively.

3.1. Simple Relation Decoupling Module

The *Simple Relation Decoupling (SRD)* module is designed to structurally decompose a multi-anchor query into mul-

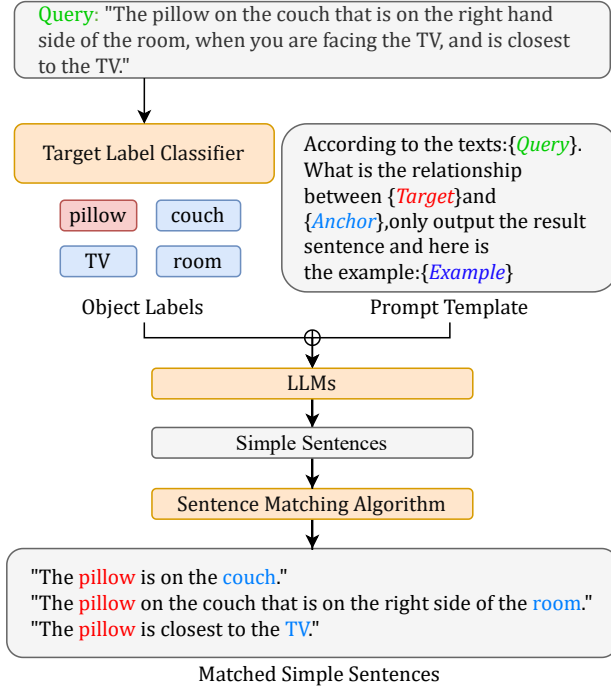


Figure 3. Overview of the SRD Module pipeline.

multiple simpler single-anchor queries, enhancing the text encoder’s ability to comprehend and process relational information. As illustrated in Fig. 3, the SRD module first predicts the target and anchor labels within a sentence, assigning them as the subject and object in the simplified sentence, respectively. This restructuring forms the foundation for generating a structured prompt, which is then fed into an LLM to produce a set of simplified queries. To maintain semantic integrity, we employ a Sentence Matching algorithm which is described in detail in the supplementary material. that filters and retains the most relevant simplified queries, ensuring that the refined queries faithfully preserve the original meaning while improving clarity and interpretability. By disentangling object relationships between the target and multiple anchors, the SRD module enables more precise relational reasoning, enhancing 3DVG performance.

Target and anchors digging. We pre-train a classifier *Clas* to identify the word in a sentence that corresponds to the $\{Target\}$ object. Given an input sentence, *Clas* first determines which word belongs to the target. Subsequently, we assess whether other words in the sentence appear in the predefined anchor set, $A_{lab} = \{A_{lab1}, A_{lab2}, \dots\}$, provided by the dataset. If a word matches an entry in this set, it is classified as an $\{Anchor\}$ object. For more details about the *Clas*, please refer to the *supplementary materials*.

Decoupled multi-anchor queries. In practical referring queries, multiple anchors frequently co-occur within the same sentence, and the spatial relationship of the target is inherently tied to the anchor labels. In such cases, spatial

descriptions involving multiple objects and their attributes often become entangled. For instance, in the $\{Query\}$ illustrated in Fig. 3, the object “couch”, which is near the target “pillow”, may dominate the spatial description, thereby weakening the relationship between the target and other anchors. To address the coupling issue in such queries, we design a set of prompt templates based on prior target and anchor digging, the process is shown in Fig. 3. Leveraging the reasoning capabilities of LLMs, we decompose complex multi-anchor queries into simpler single-anchor queries. This decoupling process clarifies the positional relationships between the target and its anchors in 3DVG, enhancing the model’s spatial understanding.

We define a total of k structured $\{Example\}$ derived from our pre-designed templates, such as “The target is on the anchor”, with additional examples provided in the supplementary materials. For each anchor, the model generates k candidate queries, where k denotes the number of generated examples. To ensure the selected sentence best aligns with the original query, we apply a sentence-matching algorithm that evaluates both label consistency and semantic consistency. The final ranking is determined by a weighted average of these two scores. For further details, please refer to the *supplementary materials*.

3.2. Textual Aggregation

Given a complex multi-anchor query sentence, the SRD module decomposes it into $(I + 1)$ sentences, where I represents the number of anchors in the original sentence. Each anchor contributes to a shorter, simplified sentence, while the original complex query remains as a longer reference sentence. To extract meaningful linguistic representations, we employ BERT [11] as the text encoder to extract $(I + 1)$ sentence features, generating a language feature set $\{F_{T_0}, F_{T_1}, \dots, F_{T_I}\}$, where F_{T_0} corresponds to the original complex query, and the remaining elements represent the features of the decoupled simpler queries. To enable the model to effectively learn from diverse sentence representations, we introduce a textual feature aggregation strategy. We randomly sample one feature from the language feature set as the main feature F_{main} , while treating the remaining features as auxiliary features F_{aux} . The final aggregated feature is computed as:

$$F_{agg} = \alpha F_{main} + (1 - \alpha) \cdot \frac{1}{I} \sum_{i=1}^I F_{aux_i}, \quad (1)$$

where α is uniformly sampled from $\{0, 0.1, 0.3, 0.5\}$ during training and fixed at 0.5 during validation. This adaptive fusion strategy ensures smooth feature integration, enhancing the model’s robustness in language-conditioned 3DVG.

3.3. Multi-view Textual-Scene Interaction Module

Cross-modal Consistent View Tokens. Previous methods have largely overlooked the inconsistency in textual descriptions arising from perspective shifts in multi-view VG, making it challenging for models to accurately interpret these variations [3, 41]. To address this limitation, we introduce a series of learnable and shared *Cross-modal Consistent View Tokens (CCVTs)*, which are integrated into both the textual and scene modules. By incorporating these tokens, both models are explicitly guided with relevant perspective information, enabling them to more effectively capture and understand the transformations and interactions induced by viewpoint changes.

Formally, let $\mathcal{V} = \{V_n | n = 1, 2, \dots, N; V_n \in \mathbb{R}^D\}$ represents the set of CCVTs, where N denotes the number of viewpoints and D represents the dimensionality of CCVTs. The CCVTs are jointly optimized with our proposed textual and scene modules. Once trained, their values remain fixed during inference, serving as a stable reference that enhances the model’s ability to comprehend multi-view scenarios and resolve perspective-induced inconsistencies.

Multi-view Textual Module. To effectively integrate sentence features from text encoders with viewpoint features extracted from CCVTs, we introduce the *Multi-view Textual Module*, which employs a cross-attention mechanism [39] to seamlessly encode viewpoint features \mathcal{V} into the textual feature space through multi-head attention operation.

Since each sentence inherently carries distinct viewpoint information, it is crucial to embed perspective-aware features into textual representations effectively. To achieve this, we first compute the normalized dot product between each view token and the 0th token of each sentence’s language feature $\{F_{T_0}, F_{T_1}, \dots, F_{T_I}\}$, as the 0th token F^0 typically aggregates the most salient semantic information. We take the average of these dot products across different sentences and compute a corresponding probability distribution using the softmax function. This probability is then used to reweight the view token, adaptively increasing its contribution when the description aligns with the viewpoint and reducing it when the description does not match. The refined viewpoint token is formulated as:

$$\mathcal{V} = \text{Softmax} \left(\frac{1}{I} \sum_{i=0}^I \frac{F_{T_i}^0 \mathcal{V}^T}{\|F_{T_i}^0\| \cdot \|\mathcal{V}\|} \right) \mathcal{V}. \quad (2)$$

Subsequently, the aggregated features F_{agg} , as introduced in Section 3.2, serve as the query, while the viewpoint features \mathcal{V} act as the key and value in the attention computation. The textual feature enriched with viewpoint embeddings, denoted as F'_q , is formulated as:

$$F'_q = \text{Softmax} \left(\frac{(W_q F_{agg})(W_k \mathcal{V})^T}{\sqrt{D}} \right) W_v \mathcal{V}, \quad (3)$$

where W_q , W_k , and W_v are learnable linear projection matrices. Following this, F'_q undergoes an additional self-attention operation to further refine the textual features, ensuring that the encoded representations effectively capture perspective-dependent information.

Multi-view Scene Module. To effectively capture object features across diverse scenes, we introduce a *Multi-View Scene Module* that extracts and refines scene representations from multiple viewpoints. To achieve this, we employ PointNet++ [35] as the scene encoder, computing scene features F_{V_n} for each viewpoint, where $n \in N$ denotes the scene index across N viewpoints. Each scene feature F_{V_n} consists of object-level representations, expressed as $\{F_{o1}, F_{o2}, \dots, F_{oi}\}$, where i corresponds to the number of objects present in the scene.

To explicitly inform the model of the current scene, we concatenate our CCVTs V_n with the extracted scene features, forming the input representation:

$$\mathbf{X}_n = \{F_{V_n}, V_n\} = \{F_{o1}, F_{o2}, \dots, F_{oi}, V_n\}. \quad (4)$$

These combined feature representations are then processed through several Transformer layers [39], denoted as $\text{Trans}(\cdot)$, which enhances the relational encoding between objects and viewpoints. This mechanism ensures that both global scene context and fine-grained object details are effectively captured:

$$\mathbf{Z}_n^{(l+1)} = \text{Trans}^{(l)}(\mathbf{Z}_n^{(l)}), \quad (5)$$

where the initial input to the Transformer is $\mathbf{Z}_n^{(0)} = \mathbf{X}_n$, and $\mathbf{Z}_n^{(L)}$ represents the refined features after L Transformer layers.

At the final Transformer layer, the output consists of both *[object]* tokens and *[view]* tokens. Since the transformed features F'_{V_n} encapsulate both object-specific and viewpoint information, we retain only the *[object]* tokens for the subsequent grounding task:

$$F'_{V_n} = \{F'_{o1}, F'_{o2}, \dots, F'_{oi}\}. \quad (6)$$

This design ensures that object representations are enriched with multi-view contextual information while maintaining their distinct semantic properties for accurate 3DVG.

3.4. Textual-Scene Reasoning Module

With the above formulation, we obtain the view-interactive textual features F'_q and scene features F'_V , where $F'_V = \{F'_{V_n} | n = 1, 2, \dots, N\}$, each enriched with viewpoint information. These features are then processed through the proposed *Textual-Scene Reasoning Module* to generate the final prediction. This module primarily consists of a Transformer with a cross-attention mechanism [39], where F'_V serves as the query, while F'_q functions as the key and

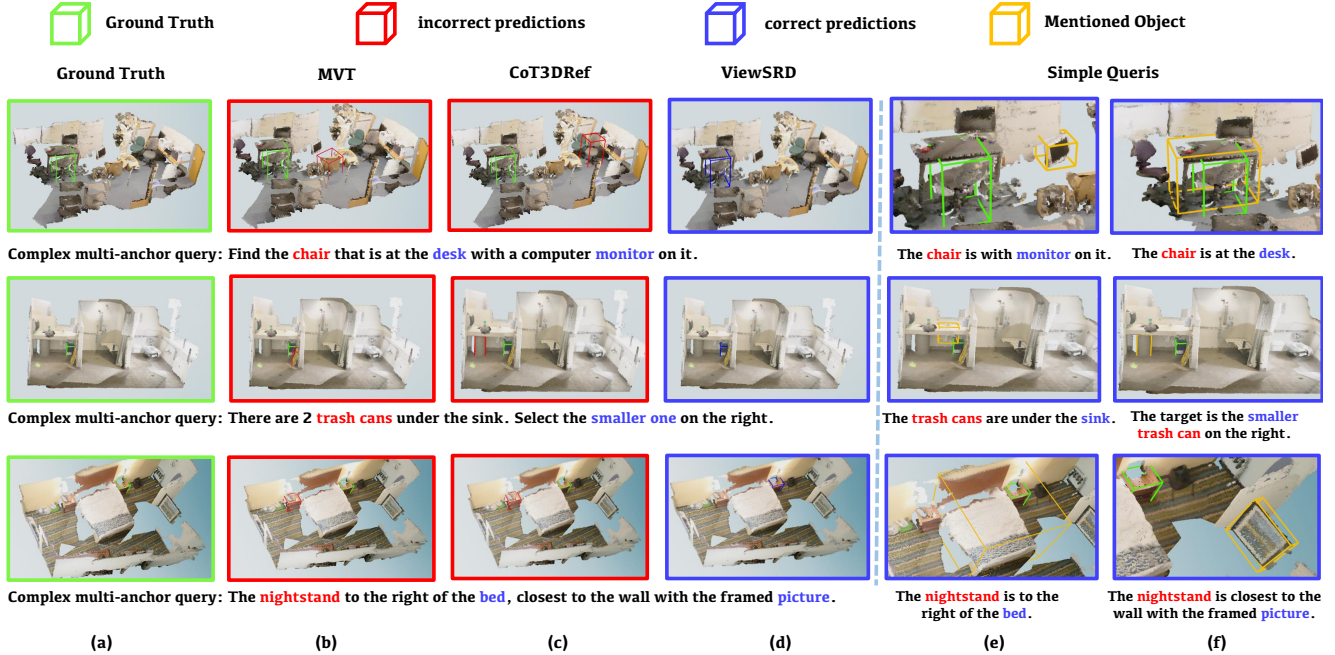


Figure 4. Visualization Results of the 3D Visual Grounding Results. For the presented 3D scenes, we utilize green, red, blue, and yellow boxes to represent the ground truth, incorrect predictions, correct predictions, and the mentioned objects, respectively. Columns (e) and (f) present the decomposed simple queries derived from the complex queries.

value, facilitating fine-grained alignment between textual and visual representations. Additionally, a View Aggregation mechanism integrates information across multiple viewpoints by computing both the average and maximum values of the output features. Finally, a Prediction Head projects the aggregated features into the result space, enabling a view-aware 3D Visual Grounding model capable of effectively reasoning across multiple perspectives.

3.5. Overall Loss Functions

Following prior research [18, 36], three distinct loss functions are applied on *ViewSRD*. These include a referential loss \mathcal{L}_{Ref} derived from grounding predictions, an object-level loss \mathcal{L}_{Object} capturing object shape and center and a sentence-level loss \mathcal{L}_{Sent} designed to identify the target and anchor phrases within the F_{agg} . Similarly to [2], we extend the referential loss \mathcal{L}_{Ref} to localize both the target and the anchors, which we term as parallel referential loss \mathcal{L}_{ref}^P , where both the target and anchors are localized simultaneously. For details of these losses, please refer to *supplementary materials*. The total loss function is defined as:

$$\mathcal{L} = \lambda_{Obj} \mathcal{L}_{Object} + \lambda_{Ref} \mathcal{L}_{Ref}^P + \lambda_{Sent} \mathcal{L}_{Sent}. \quad (7)$$

4. Experiments

4.1. Experiment Settings

Datasets. **Nr3D** [1] contains 45,503 human utterances referencing 707 indoor scenes from ScanNet [10], cover-

ing 76 object categories with multiple same-class distractors. **Sr3D** [1] includes 83,572 template-based sentences in a “target–spatial relation–anchor” format, offering a simpler setup with similar distractors. **ScanRefer** [6] provides 51,583 free-form descriptions for 11,046 objects across 800 ScanNet scenes, incorporating spatial and attribute-level references to support 3DVG.

Evaluation Metrics. For Nr3D and Sr3D, grounding accuracy is measured by the percentage of correctly matched boxes [18, 36]. For ScanRefer, we report Acc@0.25 and Acc@0.5, i.e., the percentage of predicted boxes with IoU exceeding 0.25 or 0.5, respectively [42].

Implementation Details. All experiments are implemented in PyTorch and run on a single RTX 4090 GPU. We use AdamW [30] with a learning rate of 0.0005. The number of input views is set to $N = 4$. We set $\lambda_{Ref}, \lambda_{Obj}, \lambda_{Sent} = 1.0, 0.5, 0.5$. For the SRD module, we adopt DeepSeekR1 [28], which balances performance and reproducibility, and can runs on a RTX 4090.

4.2. 3D Visual Grounding Results

We compare *ViewSRD* with recent state-of-the-art approaches to evaluate its effectiveness on 3DVG. Fig. 4 illustrates complex query cases from Nr3D [1], including ground truth boxes, predictions from MVT [18], CoT3DRef [2], and *ViewSRD*, along with the original queries and the simplified sentences produced by the SRD module. In multi-anchor scenarios (e.g., involving “bed”, “table”, and “chair”), MVT and CoT3DRef often misalign predictions

Table 1. Performance (%) comparison on Nr3D [1] and Sr3D [1].

Method	Nr3D					Sr3D				
	Overall	Easy	Hard	View Dep.	View Indep.	Overall	Easy	Hard	View Dep.	View Indep.
3DVG-Transformer [54]	40.8	48.5	34.8	34.8	43.7	51.4	54.2	44.9	44.6	51.7
LanguageRefer [36]	43.9	51.0	36.6	41.7	45.0	56.0	58.9	49.3	49.2	56.3
TransRefer3D [16]	42.1	48.5	36.0	36.5	44.9	57.4	60.5	50.2	49.9	57.7
SAT [48]	49.2	56.3	42.4	46.9	50.4	57.9	61.2	50.0	49.2	58.3
MVT [18]	55.1	61.3	49.1	54.3	55.4	64.5	66.9	58.8	58.4	64.7
ViewRefer [15]	56.0	63.0	49.7	55.1	56.8	67.0	68.9	62.1	52.2	67.7
MiKASA [5]	64.4	69.7	59.4	65.4	64.0	75.2	78.6	67.3	70.4	75.4
CoT3DRef [2]	64.4	70.0	59.2	61.9	65.7	73.2	75.2	67.9	67.6	73.5
ViewSRD (ours)	69.9	75.3	64.8	68.6	70.6	76.0	78.3	70.6	69.0	76.2

due to challenges in spatial reasoning. In contrast, *ViewSRD* correctly grounds targets by decomposing complex queries and leveraging robust spatial relationships between target-anchor pairs. Moreover, under viewpoint shifts, CoT3DRef struggles to maintain alignment, whereas *ViewSRD* reliably grounds targets by capturing spatial relations invariant to viewpoint changes (e.g., “The trash cans are under the sink”).

Quantitative results on Nr3D (Table 1) show that *ViewSRD* achieves a 5.2% accuracy gain over the best prior method, CoT3DRef, under identical settings. Under viewpoint-dependent evaluation, it further outperforms CoT3DRef by 6.7%, demonstrating the effectiveness of CCVTs in aligning textual and visual spaces and modeling viewpoint-sensitive relations through query decomposition. To assess generalization, we also evaluate on Sr3D [1] (Table 1). *ViewSRD* achieves the highest accuracy of 76.2% in the View-Independent setting, with additional gains of +2.8% and +2.7% in the View-Independence and Hard scenarios, respectively. These results confirm the robustness and generalizability of our approach across diverse scenario.

4.3. Analysis of Anchors

In this section, we analyze the impact of the number of anchors in a query on 3DVG performance. The results presented in Table 2 underscore the effectiveness of our approach, particularly in multi-anchor scenarios, where our method successfully disentangles spatial relationships by explicitly modeling target-anchor interactions. In contrast, existing methods such as MVT [18] and CoT3DRef [2], which do not account for the necessity of spatial relationship decoupling, exhibit a notable performance decline in multi-anchor queries compared to single-anchor cases. Notably, our approach achieves higher accuracy in multi-anchor queries than in single-anchor ones, demonstrating that when properly processed, multi-anchor information enhances 3DVG performance rather than introducing ambiguity. These findings validate the efficacy of ViewSRD in effectively leveraging complex spatial relationships for improved grounding accuracy.

Table 2. Performance (%) comparison on Nr3D [1] with new criteria Multi-Anc and Single-Anc.

Model	Multi-Anc	Single-Anc	Overall
MVT [18]	52.6	56.6	55.1
CoT3DRef [2]	63.1	65.2	64.4
ViewSRD	71.5	69.5	69.9

Table 3. Performance (%) of SRD module improves MVT [18], BUTD-DETR [21] and EDA [42] on ScanRefer [6] dataset.

Method	Unique (19%) Multiple (81%)				Overall	
	0.25	0.5	0.25	0.5	0.25	0.5
MVT [18]	77.7	66.5	31.9	25.3	40.8	33.3
MVT+SRD	78.6	67.2	34.1	27.1	42.1 (3.2%↑)	34.3(3.0%↑)
BUTD-DETR [21]	82.8	64.9	44.7	33.9	50.4	38.6
BUTD-DETR+SRD	85.0	66.2	45.3	34.2	57.9 (14.9%↑)	45.7 (18.4%↑)
EDA [42]	80.4	65.3	35.6	25.1	43.6	32.3
EDA+SRD	81.0	67.3	36.4	28.3	44.4 (1.8%↑)	35.3 (9.3%↑)
ViewSRD	82.1	68.2	37.4	29.0	45.4	36.0

4.4. SRD Enhances Other 3DVG Methods.

Our SRD module is inherently model-agnostic, operating independently of the training process by focusing exclusively on decoupling complex multi-anchor queries into simpler single-anchor queries. This decoupling mechanism reduces ambiguity in multi-anchor descriptions, enhances target grounding, and serves as a model-independent pre-processing step, ensuring seamless compatibility with various 3DVG methods to improve performance without modifying existing architectures. As demonstrated in Table 3, integrating SRD into MVT [18], BUTD-DETR [21] and EDA [42] consistently leads to performance improvements. These results highlight SRD’s ability to refine query interpretation by effectively disentangling target-anchor relationships, thereby reducing errors introduced by complex linguistic structures. These improvements reinforce the critical role of SRD module in enhancing accuracy of 3DVG.

Table 4. Ablation studies on Nr3D [1]. All components contribute to final performance(%).

Component	Overall	Easy	Hard	View Dep.	View Indep.
w/o CCVTs.	62.2	68.5	56.1	60.1	63.2
w/o Textual M.	68.0	73.5	62.6	67.6	68.1
w/o Scene M.	64.6	70.5	58.9	63.8	64.9
w/o SRD M.	68.6	73.0	64.8	66.5	70.0
w/o Weight.	69.0	74.2	64.0	66.5	70.2
LLM-Aug.	69.1	74.5	63.7	68.0	69.5
ViewSRD	69.9	75.3	64.8	68.6	70.6

Table 5. Ablation of view numbers on Nr3D [1].

View Number		Overall	Easy	Hard	View Dep.	View Indep.
Train	Test					
4	1	66.0	71.7	60.5	64.0	67.0
4	2	68.9	75.1	63.0	66.9	69.9
4	4	69.9	75.3	64.8	68.6	70.6
1	1	64.4	70.9	58.1	60.8	66.2
2	2	67.7	73.0	62.5	66.1	68.4
8	8	68.4	74.1	63.0	67.4	68.9

4.5. Ablation Study

Analysis of ViewSRD Components. To assess the contribution of each component within ViewSRD, we conducted detailed ablation studies on the Nr3D dataset [1]. Starting from the full model, we systematically removed key modules one at a time to evaluate their individual impact. The results, presented in Table 4, demonstrate that each component plays a crucial role in enhancing model performance across different scenarios. Notably, the removal of the CCVTs leads to the most significant performance degradation. This is primarily because, without the view token, the model lacks explicit viewpoint information, impairing its ability to distinguish between different perspectives. Similarly, removing either the textual module or the scene module results in a noticeable decline, underscoring the necessity of cross-modal interaction. When view-alignment weighting is disabled (*w/o Weight*), performance drops by 0.9%, showing that dynamic alignment of view features is critical for performance under view-dependent conditions. Removing the SRD module leads to performance degradation, confirming the benefit of multi-anchor query decoupling. We also compare it with an LLM-based augmentation method from Multi3DRefer [52] and find that SRD achieves greater gains, highlighting the advantage of structured query decomposition over generic augmentation.

Analysis of Multi-View Modeling. We evaluate the effect of varying view counts on 3DVG performance using the Nr3D dataset. As shown in Table 5, testing with more views consistently improves accuracy when the model is trained with four views, highlighting the benefit of aggregating complementary spatial cues from multiple perspectives. When training and testing with the same number of

Table 6. Accuracy comparison when replacing different LLMs in SRD module on Nr3D [1].

LLM decoupler	Accuracy
OpenChat [40]	69.6%
DeepSeek-R1 [28]	69.9%
Qwen-Plus [46]	70.5%
Qwen-Turbo [46]	70.7%

views, performance improves from 64.4% (1 view) to 67.7% (2 views), but plateaus at 68.4% with 8 views, suggesting diminishing returns. Notably, four views offer a strong trade-off, capturing diverse spatial information with minimal redundancy and maintaining computational efficiency. This also suggests that uniformly attending to many views may dilute focus on key perspectives. Future work will explore adaptive view selection.

Analysis of SRD’s LLM Decoupler. In this paper, we employ the open-source DeepSeek-R1 [28] as the LLM in the SRD module and further investigate the impact of different LLMs on the final performance of 3DVG. As shown in Table 6, different LLM decouplers exhibit varying levels of effectiveness in the sentence decoupling task. Models with stronger decoupling capabilities yield better results. For instance, OpenChat [40] and DeepSeek-R1 [28] achieve accuracies of 69.6% and 69.9%, respectively, while models designed with enhanced sentence decoupling capabilities, such as Qwen-Plus [46] and Qwen-Turbo [46], achieve 70.5% and 70.7%, with Qwen-Turbo demonstrating the highest performance. These results indicate that as an LLM’s ability to disentangle complex sentence structures improves, it becomes more effective at isolating and extracting relevant information, ultimately leading to significant gains in 3DVG accuracy.

5. Conclusion

In this paper, we introduce ViewSRD, a framework that disentangles target-anchor relationships via the Simple Relation Decoupling (SRD) module and enhances multi-view understanding through the Multi-view Textual-Scene Interaction (Multi-TSI) module. By decomposing complex multi-anchor queries into simpler single-anchor sentences, SRD clarifies positional relationships, while Multi-TSI integrates textual and scene features across viewpoints using cross-modal consistent view tokens (CCVTs) to capture spatial interactions. Extensive experiments demonstrate ViewSRD’s state-of-the-art performance in 3DVG.

A limitation of ViewSRD is its assumption that complex queries can be fully decomposed without overlapping relationships. While the decomposition into overlapping relations does not degrade performance, it diminishes the intended benefits of simplification. Future work will explore adaptive query to better preserve contextual dependencies.

Acknowledgements. This work is supported by Guangdong Provincial Natural Science Foundation for Outstanding Youth Team Project (No. 2024B1515040010), NSFC Key Project (No. U23A20391), China National Key R&D Program (Grant No. 2023YFE0202700, 2024YFB4709200), Key-Area Research and Development Program of Guangzhou City (No. 2023B01J0022), Guangdong Natural Science Funds for Distinguished Young Scholars (Grant 2023B1515020097), the National Research Foundation, Singapore under its AI Singapore Programme (AISG Award No.: AISG3-GV-2023-011), the Singapore Ministry of Education AcRF Tier 1 Grant (Grant No.: MSS25C004), and the Lee Kong Chian Fellowships.

References

- [1] 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 *ECCV*, pages 422–440, 2020. [2](#), [6](#), [7](#), [8](#)
- [2] Eslam Mohamed Bakr, Mohamed Ayman, Mahmoud Ahmed, Habib Slim, and Mohamed Elhoseiny. Cot3dref: Chain-of-thoughts data-efficient 3d visual grounding. *The Twelfth International Conference on Learning Representations*, 2024. [6](#), [7](#)
- [3] Christopher Beckham, Martin Weiss, Florian Golemo, Sina Honari, Derek Nowrouzezahrai, and Christopher Pal. Visual question answering from another perspective: Clevr mental rotation tests. *Pattern Recognition*, 136:109209, 2023. [5](#)
- [4] Francisco Maria Calisto, João Fernandes, Margarida Morais, Carlos Santiago, João Maria Abrantes, Nuno Nunes, and Jacinto C Nascimento. Assertiveness-based agent communication for a personalized medicine on medical imaging diagnosis. In *Proceedings of the 2023 CHI conference on human factors in computing systems*, pages 1–20, 2023. [1](#)
- [5] 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. [1](#), [3](#), [7](#)
- [6] Dave Zhenyu Chen, Angel X Chang, and Matthias Nießner. Scanrefer: 3d object localization in rgb-d scans using natural language. In *ECCV*, pages 202–221. Springer, 2020. [2](#), [6](#), [7](#)
- [7] Weihong Chen, Xuemiao Xu, Haoxin Yang, Yi Xie, Peng Xiao, Cheng Xu, Huaidong Zhang, and Pheng-Ann Heng. Scjd: Sparse correlation and joint distillation for efficient 3d human pose estimation. *arXiv preprint arXiv:2503.14097*, 2025. [2](#)
- [8] Tushar Choudhary, Vikrant Dewangan, Shivam Chandhok, Shubham Priyadarshan, Anushka Jain, Arun K Singh, Siddharth Srivastava, Krishna Murthy Jatavallabhula, and K Madhava Krishna. Talk2bev: Language-enhanced bird’s-eye view maps for autonomous driving. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 16345–16352. IEEE, 2024. [1](#)
- [9] Can Cui, Yunsheng Ma, Xu Cao, Wenqian Ye, and Ziran Wang. Drive as you speak: Enabling human-like interaction with large language models in autonomous vehicles. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 902–909, 2024. [1](#)
- [10] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *CVPR*, pages 5828–5839, 2017. [2](#), [6](#)
- [11] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. [4](#)
- [12] Azade Farshad, Yousef Yeganeh, Yu Chi, Chengzhi Shen, Björn Ommer, and Nassir Navab. Scenegenie: Scene graph guided diffusion models for image synthesis. In *ICCV*, pages 88–98, 2023. [2](#)
- [13] Shuo Feng, Haowei Sun, Xintao Yan, Haojie Zhu, Zhengxia Zou, Shengyin Shen, and Henry X Liu. Dense reinforcement learning for safety validation of autonomous vehicles. *Nature*, 615(7953):620–627, 2023. [1](#)
- [14] Liang Geng and Jianqin Yin. Viewinfer3d: 3d visual grounding based on embodied viewpoint inference. *IEEE Robotics and Automation Letters*, 2024. [1](#)
- [15] 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 *ICCV*, pages 15372–15383, 2023. [1](#), [3](#), [7](#)
- [16] Dailan He, Yusheng Zhao, Junyu Luo, Tianrui Hui, Shaofei Huang, Aixi Zhang, and Si Liu. Transrefer3d: Entity-and-relation aware transformer for fine-grained 3d visual grounding. In *ACM MM*, pages 2344–2352, 2021. [7](#)
- [17] 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](#)
- [18] Shijia Huang, Yilun Chen, Jiaya Jia, and Liwei Wang. Multi-view transformer for 3d visual grounding. In *CVPR*, pages 15524–15533, 2022. [2](#), [3](#), [6](#), [7](#)
- [19] Yufeng Huang, Jiji Tang, Zhuo Chen, Rongsheng Zhang, Xinfeng Zhang, Weijie Chen, Zeng Zhao, Zhou Zhao, Tangjie Lv, Zhipeng Hu, et al. Structure-clip: Towards scene graph knowledge to enhance multi-modal structured representations. In *AAAI*, pages 2417–2425, 2024. [1](#), [2](#)
- [20] Zixin Huang, Xuesong Tao, and Xinyuan Liu. Nan-detr: noising multi-anchor makes detr better for object detection. *Frontiers in Neurorobotics*, 18:1484088, 2024. [1](#)
- [21] 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. [7](#)
- [22] Xinjie Jiang, Chenxi Zheng, Xuemiao Xu, Bangzhen Liu, Weiyang Zheng, Huaidong Zhang, and Shengfeng He. Vr-done: One-stage video visual relation detection. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 1437–1446, 2024. [2](#)
- [23] Jingdan Kang, Haoxin Yang, Yan Cai, Huaidong Zhang, Xuemiao Xu, Yong Du, and Shengfeng He. Sita: Structurally

- imperceptible and transferable adversarial attacks for stylized image generation. *IEEE Transactions on Information Forensics and Security*, 2025. 2
- [24] Sanjoy Kundu and Sathyanarayanan N Aakur. Is-ggt: Iterative scene graph generation with generative transformers. In *CVPR*, pages 6292–6301, 2023. 2
- [25] Jialu Li and Mohit Bansal. Panogen: Text-conditioned panoramic environment generation for vision-and-language navigation. *Advances in Neural Information Processing Systems*, 36, 2024. 1
- [26] Yihong Lin, Xuemiao Xu, Huaidong Zhang, Cheng Xu, Weijie Li, Yi Xie, Jing Qin, and Shengfeng He. Delving into invisible semantics for generalized one-shot neural human rendering. *IEEE Transactions on Visualization and Computer Graphics*, 2025. 2
- [27] Xiongkun Linghu, Jiangyong Huang, Xuesong Niu, Xiaojian Shawn Ma, Baoxiong Jia, and Siyuan Huang. Multi-modal situated reasoning in 3d scenes. *Advances in Neural Information Processing Systems*, 37:140903–140936, 2025. 2
- [28] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024. 2, 6, 8
- [29] Bangzhen Liu, Chenxi Zheng, Xuemiao Xu, Cheng Xu, Huaidong Zhang, and Shengfeng He. Rotation-adaptive point cloud domain generalization via intricate orientation learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025. 2
- [30] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 6
- [31] Ruiyuan Lyu, Jingli Lin, Tai Wang, Xiaohan Mao, Yilun Chen, Runsen Xu, Haifeng Huang, Chenming Zhu, Dahua Lin, and Jiangmiao Pang. Mmscan: A multi-modal 3d scene dataset with hierarchical grounded language annotations. *Advances in Neural Information Processing Systems*, 37:50898–50924, 2025. 1
- [32] Yunze Man, Liang-Yan Gui, and Yu-Xiong Wang. Situational awareness matters in 3d vision language reasoning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13678–13688, 2024. 1
- [33] Honghan Pan, Bangzhen Liu, Xuemiao Xu, Chenxi Zheng, Yongwei Nie, and Shengfeng He. Gaussian prompter: Linking 2d prompts for 3d gaussian segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025. 2
- [34] Itthisak Phueaksri, Marc A Kastner, Yasutomo Kawanishi, Takahiro Komamizu, and Ichiro Ide. An approach to generate a caption for an image collection using scene graph generation. *IEEE Access*, 2023. 2
- [35] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in neural information processing systems*, 30, 2017. 5
- [36] Junha Roh, Karthik Desingh, Ali Farhadi, and Dieter Fox. Language-refer: Spatial-language model for 3d visual grounding. In *Conference on Robot Learning*, pages 1046–1056. PMLR, 2022. 6, 7
- [37] 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. 1
- [38] Yashar Talebirad and Amirhossein Nadiri. Multi-agent collaboration: Harnessing the power of intelligent llm agents. *arXiv preprint arXiv:2306.03314*, 2023. 1
- [39] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 5
- [40] Guan Wang, Sijie Cheng, Xianyuan Zhan, Xiangang Li, Sen Song, and Yang Liu. Openchat: Advancing open-source language models with mixed-quality data. *arXiv preprint arXiv:2309.11235*, 2023. 2, 8
- [41] 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. 5
- [42] Yanmin Wu, Xinhua Cheng, Renrui Zhang, Zesen Cheng, and Jian Zhang. Eda: Explicit text-decoupling and dense alignment for 3d visual grounding. In *CVPR*, pages 19231–19242, 2023. 1, 2, 6, 7
- [43] Cheng Xu, Wei Qu, Xuemiao Xu, and Xueting Liu. Multi-scale flow-based occluding effect and content separation for cartoon animations. *IEEE Transactions on Visualization and Computer Graphics*, 29(9):4001–4014, 2022. 2
- [44] Can Xu, Yuehui Han, Rui Xu, Le Hui, Jin Xie, and Jian Yang. Multi-attribute interactions matter for 3d visual grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17253–17262, 2024. 1, 2
- [45] Yingjie Xu, Bangzhen Liu, Hao Tang, Bailin Deng, and Shengfeng He. Learning with unreliability: Fast few-shot voxel radiance fields with relative geometric consistency. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20342–20351, 2024. 2
- [46] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024. 2, 8
- [47] Haoxin Yang, Xuemiao Xu, Cheng Xu, Huaidong Zhang, Jing Qin, Yi Wang, Pheng-Ann Heng, and Shengfeng He. G2face: High-fidelity reversible face anonymization via generative and geometric priors. *IEEE Transactions on Information Forensics and Security*, 2024. 2
- [48] Zhengyuan Yang, Songyang Zhang, Liwei Wang, and Jiebo Luo. Sat: 2d semantics assisted training for 3d visual grounding. In *ICCV*, pages 1856–1866, 2021. 7
- [49] Zhao Yang, Jiakuan Liu, Yucheng Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang Yu. Appagent: Multimodal agents as smartphone users. *arXiv preprint arXiv:2312.13771*, 2023. 1
- [50] Qifan Yu, Juncheng Li, Yu Wu, Siliang Tang, Wei Ji, and Yueting Zhuang. Visually-prompted language model for fine-grained scene graph generation in an open world. In *ICCV*, pages 21560–21571, 2023. 2

- [51] 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. [1](#), [3](#)
- [52] Yiming Zhang, ZeMing Gong, and Angel X Chang. Multi3drefer: Grounding text description to multiple 3d objects. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15225–15236, 2023. [8](#)
- [53] Yuqi Zhang, Han Luo, and Yinjie Lei. Towards clip-driven language-free 3d visual grounding via 2d-3d relational enhancement and consistency. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13063–13072, 2024. [2](#)
- [54] Lichen Zhao, Daigang Cai, Lu Sheng, and Dong Xu. 3dvg-transformer: Relation modeling for visual grounding on point clouds. In *ICCV*, pages 2928–2937, 2021. [7](#)
- [55] Chenxi Zheng, Yihong Lin, Bangzhen Liu, Xuemiao Xu, Yongwei Nie, and Shengfeng He. Recdreamer: Consistent text-to-3d generation via uniform score distillation. In *The Thirteenth International Conference on Learning Representations*. [2](#)
- [56] Chenxi Zheng, Bangzhen Liu, Xuemiao Xu, Huaidong Zhang, and Shengfeng He. Learning an interpretable stylized subspace for 3d-aware animatable artforms. *IEEE Transactions on Visualization and Computer Graphics*, 31(2):1465–1477, 2024. [2](#)
- [57] Gengze Zhou, Yicong Hong, and Qi Wu. Navgpt: Explicit reasoning in vision-and-language navigation with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 7641–7649, 2024. [1](#)
- [58] 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. [1](#)