

# Dense Object Grounding in 3D Scenes

Wencan Huang\*

huangwencan@stu.pku.edu.cn  
Wangxuan Institute of Computer  
Technology, Peking University  
Beijing, China

Daizong Liu\*

dzliu@stu.pku.edu.cn  
Wangxuan Institute of Computer  
Technology, Peking University  
Beijing, China

Wei Hu†

forhuwei@pku.edu.cn  
Wangxuan Institute of Computer  
Technology, Peking University  
Beijing, China

## ABSTRACT

Localizing objects in 3D scenes according to the semantics of a given natural language is a fundamental yet important task in the field of multimedia understanding, which benefits various real-world applications such as robotics and autonomous driving. However, the majority of existing 3D object grounding methods are restricted to a single-sentence input describing an individual object, which cannot comprehend and reason more contextualized descriptions of multiple objects in more practical 3D cases. To this end, we introduce a new challenging task, called 3D Dense Object Grounding (3D DOG), to jointly localize multiple objects described in a more complicated paragraph rather than a single sentence. Instead of naively localizing each sentence-guided object independently, we found that dense objects described in the same paragraph are often semantically related and spatially located in a focused region of the 3D scene. To explore such semantic and spatial relationships of densely referred objects for more accurate localization, we propose a novel Stacked Transformer based framework for 3D DOG, named 3DOGSFormer. Specifically, we first devise a contextual query-driven local transformer decoder to generate initial grounding proposals for each target object. The design of these contextual queries enables the model to capture linguistic semantic relationships of objects in the paragraph in a lightweight manner. Then, we employ a proposal-guided global transformer decoder that exploits the local object features to learn their correlation for further refining initial grounding proposals. In particular, we develop two types of proposal-guided attention layers to encode both explicit and implicit pairwise spatial relations to enhance 3D relation understanding. Extensive experiments on three challenging benchmarks (Nr3D, Sr3D, and ScanRefer) show that our proposed 3DOGSFormer outperforms state-of-the-art 3D single-object grounding methods and their dense-object variants by significant margins.

## CCS CONCEPTS

• **Computing methodologies** → **Artificial intelligence; Computer vision; Computer vision tasks.**

\*Equal Contribution.

†This work was supported by National Natural Science Foundation of China (61972009). Corresponding author: Wei Hu (forhuwei@pku.edu.cn).

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

MM '23, October 29–November 3, 2023, Ottawa, ON, Canada.

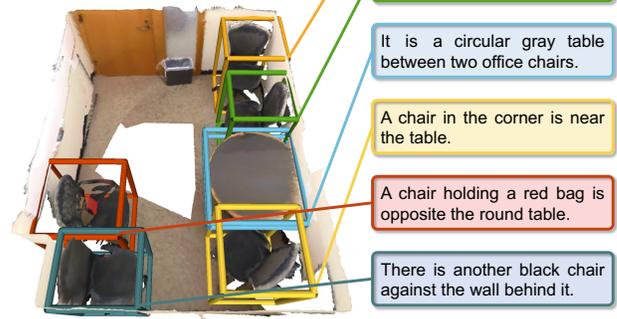
© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0108-5/23/10...\$15.00

<https://doi.org/10.1145/3581783.3611902>

**Paragraph Description:** It is a black chair to the right of the trash can. There is a black chair beside the table next to it. It is a circular gray table between two office chairs. A chair in the corner is near the table. A chair holding a red bag is opposite the round table. There is another black chair against the wall behind it.

**3D Scene:**



**Figure 1: Dense Object Grounding in 3D Scenes.** Given a paragraph description, 3D dense object grounding (3D DOG) aims to jointly localize described multiple objects in a 3D scene.

## KEYWORDS

3D Dense Object Grounding; Query-based Proposal Generation; Global Transformer

## ACM Reference Format:

Wencan Huang, Daizong Liu, and Wei Hu. 2023. Dense Object Grounding in 3D Scenes. In *Proceedings of the 31st ACM International Conference on Multimedia (MM '23)*, October 29–November 3, 2023, Ottawa, ON, Canada. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3581783.3611902>

## 1 INTRODUCTION

Grounding natural language in visual contents is a fundamental yet essential task in the multimedia understanding field. Traditional object grounding in 2D images aims to localize the object described by the given referring expression in an image, which has attracted much attention and made great progress [9, 14, 27, 28, 31, 37, 41, 42, 53]. Recently, researchers begin to explore real-world 3D object grounding [2, 7], where the target object referred by a sentence should be identified in a more complicated 3D scene. Compared to 2D object grounding, 3D object grounding is more challenging since it requires disambiguating more variant and complex spatial relationships and localizing the object referred by the text in 3D scenes with several same-class distractors.

The paradigm of previous studies on 3D object grounding is to simply localize each individual object referred by a single free-form sentence in a 3D scene. Existing 3D grounding approaches can be mainly categorized into two groups: (1) *Top-down* methods: these works directly follow a detection-then-matching two-stage framework [3, 11, 18, 19, 21, 44, 52, 56, 57], which first employ pre-trained 3D object detectors to generate multiple object proposals, then select the best proposal according to cross-modal similarity scores with the given sentence. However, these methods severely rely on the quality of the detected proposals and are often very time-consuming. (2) *Bottom-up* methods: instead of using complex proposals, these methods simply incorporate the sentence features and point-level visual representations in an early-fusion manner to directly regress the bounding box at a single stage [29, 36]. Although the above two types of methods have achieved significant progress in recent years, they are limited by the single-sentence input and are not suitable for comprehending and reasoning more contextualized descriptions of multiple objects in complicated 3D scenes. *However, such a paragraph of multiple sentences that describe several objects in a specific region of a 3D scene is natural and practical in real-world applications such as robotics and autonomous driving.* For example, as shown in Figure 1, humans may be interested in multiple objects located in a focused region of the scene. Although multiple objects with the same category like “chair” may appear in the scene, humans may only be interested in one of them. Therefore, they can utilize a paragraph consisting of multiple sentences to describe not only the most concerned object but also its context-related objects, to better avoid the ambiguity.

To localize dense objects referred by a paragraph in a 3D scene, a straightforward idea is to apply well-studied 3D single-object grounding models [3, 11, 21, 36, 44] to each individual sentence in the paragraph and integrate their results for dense grounding. However, directly applying them to the multiple-sentence setting may suffer from two critical issues: (1) *Firstly*, this solution is naive and only considers a single sentence, failing to learn the contextual semantic relationships among multiple sentences. These contextual semantic relations are vital clues for accurate paragraph comprehension and precise localization of dense objects. For instance, as shown in Figure 1, grounding sentence “*There is another black chair against the wall behind it*” requires understanding the coreference relationship between the anaphor “*it*” and the target object “*chair*” referred by another sentence “*A chair holding a red bag is opposite the round table*”. (2) *Secondly*, this solution fails to leverage the spatial correlations between locations of dense objects described in the same paragraph for cross-modal spatial reasoning. The existence of such spatial correlations is due to the fact that humans are likely to refer linguistically to multiple objects located in a focused region of a 3D scene rather than describe randomly placed ones. Ignoring such spatial relations of dense target objects may lead to inferior performance in precisely finding the location of each one.

Motivated by the above observations, in this paper, we mainly focus on addressing the challenging 3D Dense Object Grounding (3D DOG) task. Given a 3D point cloud scene and a paragraph of sentence descriptions, the goal of 3D DOG is to jointly localize dense objects described by these sentences. Rather than localizing each object independently, we explore the *semantic and spatial relationships* of densely referred objects for accurate localization.

To achieve this, we propose a novel two-phase framework for 3D DOG based on Stacked Transformers (named 3DOGSFormer) due to the powerful relation modeling capability of the well-known transformer architecture [47]. Specifically, the proposed 3DOGSFormer consists of one standard transformer encoder for 3D scene encoding and two types of transformer decoders for the two-phase grounding pipeline. In the first phase, we devise a shared contextual query-driven local transformer decoder to generate initial grounding proposals for every single sentence in the paragraph description. We propose a scene-aware context aggregation and propagation (SCAP) module to generate contextual queries for each sentence, which enables to capture the semantic relations between multiple sentences in a lightweight fashion. In the second phase, we further develop a novel proposal-guided global transformer decoder to gather information from local proposal features and learn to refine the initial grounding proposals via cross-modal spatial reasoning. To be specific, we design a stack of interlaced proposal-guided self-attention (PGSA) and cross-attention (PGCA) layers to encode both explicit and implicit pairwise spatial relations for all object-object and object-point pairs, so as to enhance 3D spatial relation understanding of densely referred objects.

Our main contributions can be summarized as follows:

- We propose a new 3D DOG task to explore practical yet challenging 3D dense object grounding based on complicated paragraph descriptions.
- We develop a novel 3DOGSFormer to tackle this 3D DOG task in a two-phase grounding pipeline, where we design a contextual query-driven local transformer decoder to capture the semantic relationships within the paragraph efficiently, as well as a proposal-guided global transformer decoder to enhance the 3D spatial relation understanding.
- Extensive experiments on three challenging benchmarks (Nr3D [2], Sr3D [2], and ScanRefer [7]) show that the proposed 3DOGSFormer outperforms state-of-the-art 3D single-object grounding methods and their dense-object variants by significant margins.

## 2 RELATED WORK

### 2.1 3D Single Object Grounding

The 3D object grounding task [2, 7] aims to localize objects in 3D point clouds given a sentence. Existing approaches can be mainly categorized into two groups, namely top-down and bottom-up frameworks. The top-down methods [3, 11, 18, 19, 21, 44, 52, 56, 57] follow a detection-then-matching two-stage pipeline, which first obtain the features of the query sentence and candidate point cloud objects independently by a pre-trained language model [26] and a pre-trained 3D detector [33, 39] or segmentor [10, 23, 48], then employ various cross-modal fusion or matching mechanisms to select the best-matched object according to the sentence. Graph-based approaches [2, 18, 20, 56] and Transformer-based attention mechanisms [3, 11, 19, 21, 44, 52, 57] are widely adopted for the multi-module feature fusion in the matching stage. The obvious drawback of the top-down methods is that they severely rely on the quality of the detected proposals and are very time-consuming. By contrast, bottom-up models [29, 36] incorporate the sentence

features and voxel- or point-level visual representations in an early-fusion manner to directly regress the bounding box at a single stage, which are more flexible to identify various text-concerned objects. Typically, 3D-SPS [36] employs textual features to guide visual keypoint selection and progressively localizes objects. Recently, BEAUTY-DETR [22] develops a Transformer-based bottom-up top-down architecture, which combines the advantages of the above two pipelines. Although existing methods have made great attempts to solve the 3D object grounding problem, they are restricted to the single-sentence input describing an individual object. In this paper, we propose a new but challenging 3D DOG task to explore 3D dense object grounding based on complicated paragraph descriptions.

## 2.2 3D Dense Object Understanding

Understanding dense objects in 3D scenes has raised great interest among researchers in recent years. Scan2Cap [13] introduces the task of 3D dense captioning, which aims to jointly localize and describe dense objects in a 3D scene. To tackle this task, it builds a message-passing graph network to mine the relations among objects in a 3D scene. MORE [24] further takes multi-order relations into account to learn richer 3D object relation features. SpaCap3D [49] uses a spatiality-guided transformer to learn the contribution of surrounding objects to the target object for 3D dense captioning. 3DJCG [5] and UniT3D [8] develop unified transformer-based frameworks for both 3D dense captioning and 3D object grounding. X-Trans2Cap [55] introduces additional 2D prior information to improve 3D dense captioning with knowledge transfer. [58] shifts attention to contextual information for the perception of non-object information. Recently, Vote2Cap-DETR [12] proposes a one-stage model for 3D dense captioning that detects and generates captions in parallel. The 3D dense captioning task in these papers is to describe dense objects in the 3D scene with a paragraph of multiple sentences. In contrast, our 3D dense object grounding task can be viewed as the inverse problem of 3D dense captioning.

The methods most similar to our work are EDA [50], PhraseRefer [54], and ScanEnts3D [1], which propose to ground not only the target object but also all auxiliary objects mentioned in the referential utterance to improve 3D object grounding. However, these methods are still limited in the single-sentence setting and cannot deal with the challenging 3D dense object grounding task, which requires to comprehend the semantic and spatial relationships of dense objects described by multiple sentences in the same paragraph.

## 2.3 Transformers in 2D and 3D Scenes

Transformer [47] has achieved marvelous success in most 2D computer vision tasks [17, 32], such as visual grounding [25, 51], image captioning [15], and object detection [6, 60]. In particular, DETR [6] introduces a new query-based [46, 60] paradigm for object detection, which employs a set of object queries as candidates and feeds them to the Transformer decoder for parallel detection. Beyond 2D field, the DETR architecture has been extended for various 3D vision tasks such as 3D object detection [33, 38, 59], 3D instance segmentation [30, 45], 3D object grounding [22], and 3D dense captioning [12]. In our work, we extend the DETR architecture to build a new 3DOGSFormer framework for 3D dense object grounding, which employs stacked transformer decoders to localize dense

objects in parallel. Additionally, we leverage contextual queries for efficient semantic relation modeling and develop proposal-guided attention layers to enhance the 3D spatial relation understanding.

## 3 METHOD

Given a 3D scene and a paragraph description of multiple sentences, the goal of 3D DOG is to jointly localize dense objects described by these sentences. Generally, we represent the input 3D scene as a point cloud  $PC = [P_{in}; F_{in}] \in \mathbb{R}^{N \times (3+F)}$ , where  $P_{in} \in \mathbb{R}^{N \times 3}$  is the absolute locations for each point, and  $F_{in} \in \mathbb{R}^{N \times F}$  is additional input feature for each point, such as *color*, *normal*, *height*, or *multiview feature* introduced by [7]. We denote the input paragraph description containing  $K$  sentences as  $S$ , and the  $k$ -th sentence description  $S^k$  is presented as  $S^k = \{s_{k,i}\}_{i=1}^{T_k}$ , where  $s_{k,i}$  represents the  $i$ -th word and  $T_k$  denotes the total number of words. The expected output for the  $k$ -th sentence is the bounding box  $B^k$ , representing an estimated location of the target object corresponding to the semantics of this sentence in the 3D scene.

To tackle the challenging 3D DOG task, we propose a novel Stacked Transformer based framework called 3DOGSFormer to exploit both semantic and spatial relationships among objects in the scene. As shown in Figure 2, we first adopt a 3DETR [38] encoder as our scene encoder, and a BERT model [26] as the paragraph encoder. Then we develop a shared contextual query-driven local transformer decoder to generate initial grounding proposals for each input sentence. At last, we further devise a proposal-guided global transformer decoder that exploits the obtained local object features to learn their correlation to refine their bounding box (bbox) proposals for final grounding.

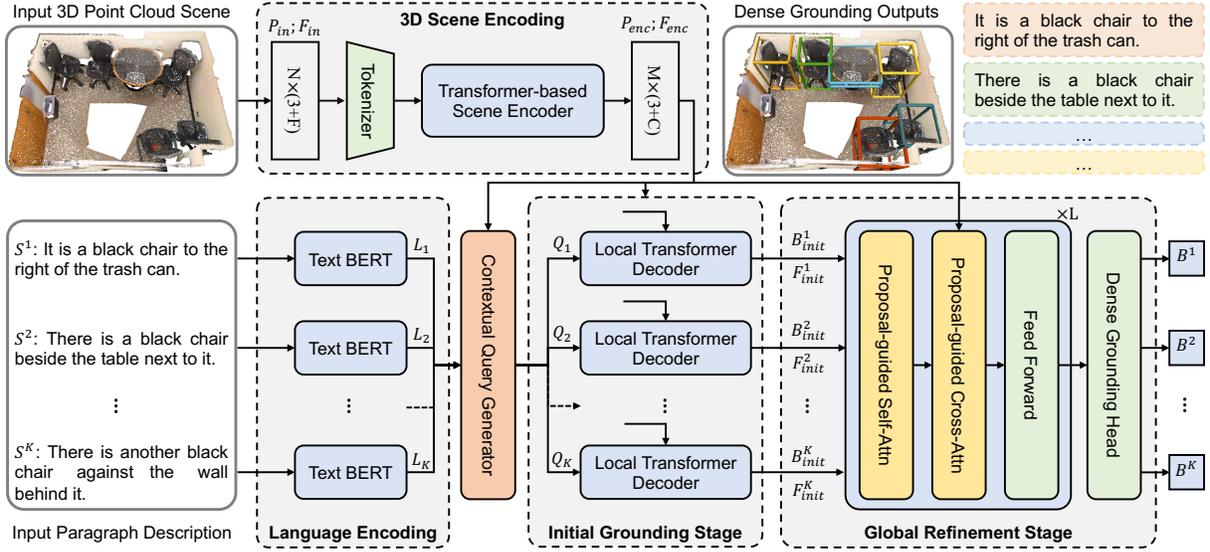
### 3.1 3D Scene and Paragraph Encoder

**3D Scene Encoder.** We exploit the powerful 3DETR [38] encoder to extract 3D visual features from the input 3D point cloud scene. During the 3D feature encoding, the input  $PC$  is first tokenized with a set-abstraction layer [40]. Then, point tokens are fed into a masked transformer encoder with a set-abstraction layer followed by another two encoder layers. We denote the encoded scene tokens as  $[P_{enc}; F_{enc}] \in \mathbb{R}^{M \times (3+C)}$ .

**Paragraph Encoder.** Following [3, 11, 44], we adopt the BERT model [26] as the paragraph encoder. Specifically, the BERT model is shared across all input sentences in the paragraph description and extracts feature vectors for each sentence separately. Given the  $k$ -th sentence  $S^k$  with  $T_k$  words, we embed them into  $D$ -dimensional feature vectors  $L_k = \left( l_s^k, l_1^k, \dots, l_{T_k}^k \right)$ , where  $l_s^k$  is the sentence-level feature,  $l_i^k$  is the feature of the  $i$ -th word.

### 3.2 Contextual Query Generator

We devise a contextual query generator to generate query embeddings for each input sentence. These queries should enable the local transformer decoder to generate the desired initial grounding proposals. To achieve this, the queries for each sentence must learn the following crucial information: (1) the complete semantics of the corresponding sentence; (2) contextual semantics among multiple sentences in the paragraph, which are vital clues for accurately understanding each individual sentence; (3) the knowledge



**Figure 2: The pipeline of our 3DOGSFormer.** After the 3D scene and paragraph encoding, we initialize the grounding proposals with a local transformer decoder, then refine them with a global transformer decoder. With the help of contextual queries and proposal-guided attention layers, the semantic and spatial relations of densely referred objects can be efficiently captured.

of the 3D spatial location and the scene-related visual information, which are proved [12, 38] to be effective in multi-modal semantic alignment and reasoning. To comply with the requirement (1), we directly utilize the encoded language feature vector  $L_k$  as the initial query embeddings for  $S^k$ . Then we design a novel scene-aware context aggregation and propagation (SCAP) mechanism to update the initial queries so as to meet the requirement (2) and (3). As shown in Figure 3, the SCAP module includes context aggregation, cross-modal interaction, and context propagation modules.

**Context Aggregation.** The high computational complexity is the biggest challenge when encoding paragraph-level contextual information into the queries for each sentence. To address this issue, we devise a context aggregation module to adaptively aggregate the contextual information in the long paragraph into a compact set of features with a much fewer number of tokens, which enables efficient contextual learning. Specifically, we first concatenate the encoded language feature vectors  $[L_1, L_2, \dots, L_K]$  to produce a paragraph feature vector  $L \in \mathbb{R}^{T \times D}$ , where  $T$  is the number of tokens in the paragraph. We then employ  $N_s$  learnable queries  $Q_{set} \in \mathbb{R}^{N_s \times C}$  to perceive the contextual information in  $L$  and aggregate it into a compact set  $F_{set} \in \mathbb{R}^{N_s \times C}$  via a plain multi-head cross-attention layer [47], as  $L' = LW_l + e_p$ ,  $F_{set} = \text{CrossAttn}(Q_{set}, L', L')$ , where  $W_l \in \mathbb{R}^{D \times C}$  is the projection matrix and  $e_p$  is the learnable positional embedding vector.

**Cross-Modal Interaction.** To inject spatial knowledge and 3D scene-related information into the contextual query features, we then use a Co-Attention module [35] to fuse the compact set of semantic features and the encoded 3D scene features. Specifically, given the input feature vectors  $F_{set}$  and  $F_{enc}$ , we conduct multi-stage co-attention in both vision-to-language and language-to-vision directions and generate a compact set of fused features

$F'_{set} \in \mathbb{R}^{N_s \times C}$  based on the language-side attention results from all stages, as:

$$F_{set}^{i+1}, F_{enc}^{i+1} = \text{CoAttn}(F_{set}^i, F_{enc}^i), i \in \{0, 1, 2\}, \quad (1)$$

$$F'_{set} = [F_{set}^0; F_{set}^1; F_{set}^2; F_{set}^3] W_f, \quad (2)$$

where  $F_{set}^i$  and  $F_{enc}^i$  are hidden features from the  $i$ -th co-attention layer,  $F_{set}^0 = F_{set}$ ,  $F_{enc}^0 = F_{enc}$ ,  $[\cdot; \cdot]$  denotes concatenation and  $W_f \in \mathbb{R}^{4C \times C}$  is the projection matrix. The cross-modality interaction is achieved by the CoAttn module, which computes the hidden vectors of one modality by attending to the other modality, given by:

$$A_i = (F_{set}^i W_{set}^i)(F_{enc}^i W_{enc}^i)^T / \sqrt{d}, \quad (3)$$

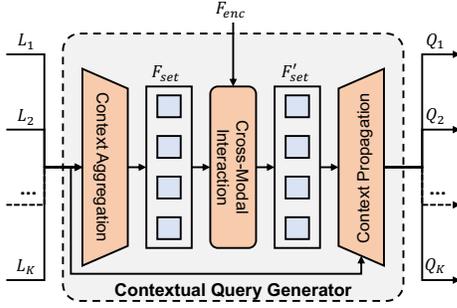
$$F_{set}^{i+1} = \text{softmax}(A_i) F_{enc}^i \hat{W}_{enc}^i, F_{enc}^{i+1} = \text{softmax}(A_i^T) F_{set}^i \hat{W}_{set}^i, \quad (4)$$

where  $d$  is the dimension of the embedding space and  $W_{set}^i, W_{enc}^i \in \mathbb{R}^{C \times d}$ ,  $\hat{W}_{enc}^i, \hat{W}_{set}^i \in \mathbb{R}^{d \times C}$  are all projection matrices.

**Context Propagation.** At last, we selectively propagate the aggregated multi-modal global information to each sentence for the contextual query generation. In detail, given the language feature vector  $L_k$  of the  $k$ -th sentence, we leverage a multi-head cross-attention layer [47] to retrieve its contextual queries from the compact set of fused features, as  $L'_k = L_k W_q$ ,  $Q_k = \text{CrossAttn}(L'_k, F'_{set}, F'_{set}) + L'_k$ , where  $W_q \in \mathbb{R}^{D \times C}$  is the projection matrix and  $Q_k \in \mathbb{R}^{(T_k+1) \times C}$  is the output contextual queries for  $S^k$ .

### 3.3 Local Transformer Decoder

We devise a local transformer decoder to generate initial grounding proposals for each sentence independently. The local transformer decoder is based on a standard Transformer decoder [47] and is shared across all sentences. For the  $k$ -th sentence, it takes as input the encoded 3D scene features and a sequence of contextual queries



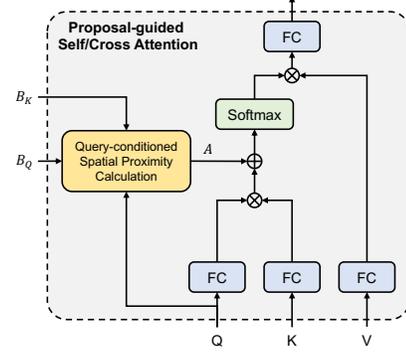
**Figure 3: The contextual query generator, which includes context aggregation, cross-modal interaction, and context propagation modules.**

$Q_k = (q_s^k, q_1^k, \dots, q_{T_k}^k)$  to produce a list of features that are then used to predict 3D-bounding boxes via a shared FFN-based grounding head. In our framework, the query embedding  $q_i^k$  represents a potential object mentioned by the  $i$ -th word in the sentence, and  $q_s^k$  is the sentence-level query embedding representing the target object described by  $s^k$ . At last, the predicted 3D-bounding box of  $q_s^k$  and the corresponding feature are obtained as the initial grounding proposal  $B_{init}^k$  and the contextual feature  $F_{init}^k$ . Following 3DETR [38] and 3D-SPS [36], 3D-bounding box estimation is formulated as box center and box size estimation. We use positional embeddings in the decoder which operates on both the 3D scene features and the query embeddings.

### 3.4 Global Transformer Decoder

To capture the 3D spatial relationships among multiple target objects and refine the initial grounding proposals, we further develop a global transformer decoder, which consists of  $L$  proposal-guided transformer layers. Each transformer layer comprises a proposal-guided self-attention layer, a proposal-guided cross-attention layer, and a feed-forward neural network (FFN). Assuming  $B_l^k \in \mathbb{R}^6$  and  $F_l^k \in \mathbb{R}^C$  are the input proposal and the contextual feature for the target object described by  $s^k$  before the  $(l+1)$ -th transformer layer, and  $B_l^k$  consists of the object center  $c_l^k \in \mathbb{R}^3$  and the object size  $s_l^k \in \mathbb{R}^3$ , we first use a linear projection layer to obtain the absolute 3D location feature as  $p^k = W_p B_l^k \in \mathbb{R}^C$  which is added to the contextual feature  $F_l^k$  to enhance the spatial information. Next, we utilize the proposal-guided self-attention layer to exploit the contextual 3D spatial relationships among dense referred objects and generate the enhanced proposal feature  $\hat{F}_l^k$  for the  $k$ -th object. After that, our novel proposal-guided cross-attention layer takes  $[B_l^k; \hat{F}_l^k]$  as queries and 3D scene features  $[P_{enc}; F_{enc}]$  as keys and values to learn the spatial relations between the initial proposals and the correct locations of each target object in the 3D scene so as to generate the refined proposal feature.

Specifically, the aforementioned proposal-guided self- or cross-attention modules are extended from the standard attention mechanism [47]. As shown in Figure 4, we use a **Query-Conditioned Spatial Proximity Calculation (QCSPC)** module to guide the



**Figure 4: The proposal-guided attention module, where the attention matrix is augmented by the query-conditioned spatial proximity matrix.**

traditional attention calculation by adding a spatial proximity matrix to the original attention matrix, given the additional input proposals  $B_Q = [c_Q; s_Q] \in \mathbb{R}^{N_Q \times 6}$  and  $B_K = [c_K; s_K] \in \mathbb{R}^{N_K \times 6}$ . Note that we set  $B_K = [P_{enc}; \varepsilon] \in \mathbb{R}^{M \times 6}$  in the proposal-guided cross-attention, where  $\varepsilon$  is a hyperparameter. To better capture the object-object or object-point pairwise 3D spatial relations, our QCSPC module computes the explicit spatial proximity matrix  $A^E$  and the implicit spatial matrix  $A^I$ , which are added together to form the spatial proximity matrix as  $A = A^E + A^I$ .

**Explicit Spatial Proximity Calculation.** We exploit a pairwise spatial feature  $f_{ij}^E \in \mathbb{R}^5$ ,  $i \in [1, N_Q]$ ,  $j \in [1, N_K]$  to model the spatial relations explicitly. For each pair of proposals  $(B_Q^i, B_K^j)$ , we compute their Euclidean distance  $d_{ij} = \|c_Q^i - c_K^j\|_2$  as well as horizontal and vertical angles  $\theta_h, \theta_v$  of the line connecting their centers  $c_Q^i$  and  $c_K^j$ . The explicit pairwise spatial feature  $f_{ij}^E$  is then defined as  $f_{ij}^E = [d_{ij}, \sin(\theta_h), \cos(\theta_h), \sin(\theta_v), \cos(\theta_v)]$ . We generate a query-conditioned weight  $g_i^E$  to select relevant spatial relations for each query proposal as  $g_i^E = Q^i W^E$ , where  $Q^i \in \mathbb{R}^C$  is the  $i$ -th input query and  $W^E \in \mathbb{R}^{C \times 5}$  is a learnable parameter. We then define the explicit spatial relevance as  $A_{ij}^E = g_i^E \cdot f_{ij}^E$ .

**Implicit Spatial Proximity Calculation.** We propose to extract implicit pairwise spatial relation features to complement the explicit spatial relations. Specifically, for each pair of 3D object proposals  $(B_Q^i, B_K^j)$ , we first obtain raw XYZ points inside these object proposals from the 3D scene as  $(O_Q^i, O_K^j)$ , and then we encode them using PointNet++ [40] networks  $\mathcal{E}_Q$  and  $\mathcal{E}_K$ , respectively, to extract their spatial relational feature vector. Next, we concatenate the pre-proposal encoding for  $B_Q^i$  and  $B_K^j$  and employ a 2-layer multi-layer perceptron to extract the  $C$ -dimensional implicit pairwise spatial relation feature for the proposal pair:

$$f_{ij}^I = \text{MLP}\left(\left[\mathcal{E}_Q(O_Q^i); \mathcal{E}_K(O_K^j)\right]\right). \quad (5)$$

A query-conditioned weight  $g_i^I = Q^i W^I$  is then generated, where  $W^I \in \mathbb{R}^{C \times C}$  is a learnable parameter. At last, we define the implicit spatial relevance for  $(B_Q^i, B_K^j)$  as  $A_{ij}^I = g_i^I \cdot f_{ij}^I$ .

**Focused Region Constraint.** Considering that humans are likely to describe multiple objects located in a focused region of a 3D scene, we propose to calculate an additional term  $A^F \in \mathbb{R}^{N_Q \times N_K}$  to adaptively constrain cross-attention within the currently focused region of the 3D scene. We propose to use the center point of all input query proposals as the focused point, *i.e.*,  $c_F = \sum_{i=1}^{N_Q} c_Q^i / N_Q$ , and compute the radius of the focused region as  $R_F = \max_{i=1}^{N_Q} \|c_F - c_Q^i\|_2$ . The additional term  $A^F$  filters points in the 3D scene by calculating their Euclidean distances to the focused point and comparing them with a threshold  $\tau R_F$ , as:

$$A_{ij}^F = \begin{cases} 0 & \text{if } \|c_F - c_K^j\|_2 < \tau R_F \\ -\infty & \text{otherwise} \end{cases}. \quad (6)$$

Empirically, we set  $\tau$  to 2. We add the term to the spatial proximity matrix when we conduct the proposal-guided cross-attention as  $A = A^E + A^I + A^F$ .

After the proposal-guided self- and cross-attention, the final FFN uses two fully connected layers to further encode the output tokens. Given the output embedding  $F_L^k$  from the last transformer layer, we use a two-layer feed-forward neural network as the dense grounding head to predict a bounding box around every referred object. The prediction consists of a center offset  $\Delta c_L^k$  and a size offset  $\Delta s_L^k$ , which are added to the initial proposal, *i.e.*,  $B^k = [c_{L-1}^k + \Delta c_L^k; s_{L-1}^k + \Delta s_L^k]$ .

### 3.5 Training Objectives

**Initial Grounding Loss.** In the initial grounding stage, we directly obtain the predicted 3D bounding boxes  $\{\hat{B}_{init}\}$  for each referring sentence in the paragraph from the local transformer decoder. To supervise the training, we use a weighted sum of an L1 loss and a Generalized IOU loss [43] following the previous work [22], as  $\mathcal{L}_{init} = \lambda_{iou} \mathcal{L}_{iou}(\hat{B}_{init}, B_{init}) + \lambda_{L1} \|\hat{B}_{init} - B_{init}\|_1$ . The  $\lambda_{iou}$  and  $\lambda_{L1}$  control the relative weighting of the two losses in the initial grounding objective.

**Refinement Loss with Iterative Prediction.** Considering the slow convergence of transformer-based model [6], we feed every layer output features  $F_l^k$  of the global transformer decoder into the shared dense grounding head to generate proposals  $B_l^k = B_{l-1}^k + \Delta B_l^k$ , and feed the generated proposals into the next transformer decoder layer as its input. During training, we use a weighted sum of a L1-based center offset regression loss  $\mathcal{L}_{cent\_reg}$  and a L1-based size offset regression loss  $\mathcal{L}_{size\_reg}$  to supervise each layer’s output, as  $\mathcal{L}_{refine} = \lambda_{cent} \mathcal{L}_{cent\_reg} + \lambda_{size} \mathcal{L}_{size\_reg}$ .

**Pretraining the Local Transformer.** Although modules in 3DOGSFormer can be trained end-to-end from scratch, we found that a simple pretrain-then-finetune strategy can stabilize the training process and improve the performance. Specifically, we first pre-train the encoders and the local transformer decoder without the refinement module and the contextual query generator. We then train the contextual query generator and finetune the encoders and the local decoder without the refinement module. At last, we train the global transformer decoder and finetune all other modules in 3DOGSFormer with the total loss  $\mathcal{L} = \lambda_{refine} \mathcal{L}_{refine} + \lambda_{init} \mathcal{L}_{init}$ .

**Data Augmentation.** During training, we use data augmentation strategies to alleviate the overfitting issue. Concretely, following previous work [57], we randomly erase some words in each sentence of the input paragraph before the BERT encoder to alleviate the issue that the grounding model is mainly decided by the prominent parts of the sentences. Additionally, we add Gaussian noise to the initial grounding proposals before every layer of the global transformer decoder to learn more robust features for more accurate grounding.

## 4 EXPERIMENTS

### 4.1 Experimental Settings

**Datasets.** We evaluate the performance on the commonly used datasets in 3D single-object grounding, ScanRefer [7] and Nr3D/Sr3D [2], and adopt the same evaluation metrics for ease of comparisons. **ScanRefer** [7] is built on 3D scenes from ScanNet [16]. ScanRefer has 36,665 free-form language annotations describing 7,875 objects from 562 3D scenes for training, and evaluates on 9,508 sentences for 2,068 objects from 141 3D scenes. According to whether the target object is a unique object class in the scene, the dataset is divided into a “unique” and a “multiple” subset in evaluation. The evaluation metric of the dataset is the  $\text{Acc}@m\text{IoU}$ , which means the fraction of descriptions whose predicted box overlaps the ground truth with  $\text{IoU} > m$ , where  $m \in \{0.25, 0.5\}$ . **Nr3D/Sr3D** [2] is also proposed based on ScanNet [16], with Nr3D containing 41,503 human-written sentences similar to ScanRefer’s text annotation and Sr3D including 83,572 synthetic expressions generated by templates. Sentences in Nr3D and Sr3D are split into “easy” and “hard” subsets in evaluation based on whether the target object contains more than one same-class distractor. GT boxes for all candidate objects in the scene are provided by these datasets. The metric is the accuracy of selecting the target bounding box among the proposals.

**Implementation Details.** For the model architecture, we set the dimension  $C = 256$  and use 8 heads for all the transformer layers. The text encoding module is a three-layer transformer initialized from BERT [26], and both the transformer encoder and decoders contain 4 attention blocks. We use the point cloud tokenizer to subsample  $M = 1024$  points. We set the size of compact set  $N_s$  to 64. We adopt the pretrain-then-finetune training process. In each training stage, we utilize rotation augmentation to increase the viewpoint invariance. The hyper-parameters  $\lambda_{iou}$ ,  $\lambda_{L1}$ ,  $\lambda_{cent}$ ,  $\lambda_{size}$ ,  $\lambda_{refine}$ , and  $\lambda_{init}$  are empirically set to 1.0, 1.0, 1.0, 1.0, 1.0, and 0.05, respectively. The hyper-parameter  $\varepsilon$  in the proposal-guided cross-attention is set to 0.01. During the end-to-end training, we use the AdamW algorithm [34] to optimize the loss function with the initial learning rate of  $5 \times 10^{-4}$  and the batch size of 2.

With the existing 3D single-object grounding datasets, we develop a strategy to simulate the proposed 3D dense object grounding setting where humans are likely to describe multiple nearby objects in a focused region. Specifically, for each training step, we first randomly sample a 3D scene from the training set. Next, we select a random object in the scene as the focused object and randomly choose its  $K - 1$  nearest objects as the concerned objects. At last, we sample a referring sentence for each of the  $K$  ( $2 \leq K \leq 12$ ) selected objects, and then we arrange these sentences into a paragraph in such a way that referring sentences for objects closer to the center

**Table 1: Performance evaluation results on the ScanRefer, Nr3D, and Sr3D datasets.**

Method	Type	ScanRefer						Nr3D			Sr3D			
		Unique		Multiple		Overall		Overall	Easy	Hard	Overall	Easy	Hard	
		@0.25	@0.5	@0.25	@0.5	@0.25	@0.5							
ScanRefer [7]	Single	67.64	46.19	32.06	21.26	38.97	26.10	34.2	41.0	23.5	-	-	-	
ReferIt3D [2]		53.80	37.50	21.00	12.80	26.40	16.90	35.6	43.6	27.9	40.8	44.7	31.5	
TGNN [20]		68.61	56.80	29.84	23.18	37.37	29.70	37.3	44.2	30.6	45.0	48.5	36.9	
InstanceRefer [56]		77.45	66.83	31.27	24.77	40.23	32.93	38.8	46.0	31.8	48.0	51.1	40.5	
SAT [52]		73.21	50.83	37.64	25.16	44.54	30.14	49.2	56.3	42.4	57.9	61.2	50.0	
FFL-3DOG [18]		78.80	67.94	35.19	25.70	41.33	34.01	41.7	48.2	35.0	-	-	-	
3DVG-Transformer [57]		77.16	58.47	38.38	28.70	45.90	34.47	40.8	48.5	34.8	51.4	54.2	44.9	
TransRefer3D [19]		-	-	-	-	-	-	42.1	48.5	36.0	57.4	60.5	50.2	
LanguageRefer [44]		-	-	-	-	-	-	43.9	51.0	36.6	56.0	58.9	49.3	
LAR [3]		-	-	-	-	-	-	48.9	58.4	42.3	59.4	63.0	51.2	
3DJCG [5]		78.75	61.30	40.13	30.08	47.62	36.14	-	-	-	-	-	-	
3D-SPS [36]		81.63	64.77	39.48	29.61	47.65	36.43	51.5	58.1	45.1	62.6	56.2	65.4	
MVT [21]		77.67	66.45	31.92	25.26	40.80	33.26	55.1	61.3	49.1	64.5	66.9	58.8	
ViL3DRel [11]		81.58	68.62	40.30	30.71	47.94	37.73	64.4	70.2	57.4	72.8	74.9	67.9	
BUTD-DETR † [22]		81.93	67.11	42.61	31.84	50.24	38.68	-	-	-	-	-	-	
ViL3DRel + BS		Dense	82.51	69.50	42.94	31.95	50.26	38.90	66.0	71.6	59.3	74.1	76.4	68.7
BUTD-DETR + BS			84.63	68.90	45.41	33.68	53.02	40.51	-	-	-	-	-	-
ViL3DRel + DepNet	87.34		71.04	47.56	35.93	54.92	42.43	70.6	75.8	64.4	78.0	80.5	72.2	
BUTD-DETR + DepNet	87.75		70.66	50.70	37.81	57.89	44.18	-	-	-	-	-	-	
3DOGSFormer (Ours)	<b>90.11</b>		<b>73.08</b>	<b>58.62</b>	<b>43.57</b>	<b>64.73</b>	<b>49.29</b>	<b>73.5</b>	<b>77.9</b>	<b>68.2</b>	<b>80.8</b>	<b>83.4</b>	<b>74.7</b>	

object have higher probability to appear earlier in the paragraph. Paragraphs with fewer than 12 sentences are padded with zeros to 12 sentences. Such a random sampling strategy reduces the model to overfit spatial relation priors of dense objects. During the evaluation, we divide sentences in the dataset into paragraphs in the same way so that the target objects of the sentences in the same paragraph form a KNN cluster in the 3D scene. We mainly sample 12 sentences per paragraph for evaluation.

**Comparison Baselines.** We compare the proposed 3DOGSFormer model with most of the existing 3D single-object grounding methods since we can adapt these methods to the dense grounding setting straightforwardly by applying them to each sentence in the paragraph. To better understand our 3DOGSFormer’s performance, we further extend two of the state-of-the-art 3D single-object grounding models, **ViL3DRel** [11] and **BUTD-DETR** [22], to the 3D DOG setting as more competitive baselines. We devise two strategies to extend these models by leveraging **Beam Search (BS)** and **DepNet** [4] (a dense grounding method in the 2D field), respectively. In specific, the BS model first localizes each sentence in the 3D scene independently with the base model, then applies beam search on the top 12 grounding results of each object as post-processing, such that the final dense object grounding results are most concentrated in the 3D scene, measured by center variances. In the DepNet-based strategy, we first extract features before the grounding head of the base model for each referring sentence, and then apply DepNet’s DEAP module [4] to produce enhanced features via global reasoning, which are fed into the grounding head to generate grounding results. We adopt a pretrain-then-finetune strategy to optimize the extended model in this case. At last, we obtain 4 combined strong baselines: **ViL3DRel+BS**, **BUTD-DETR+BS**, **ViL3DRel+DepNet**, and **BUTD-DETR+DepNet**.

## 4.2 Comparison Results

Table 1 shows the comparison results for all methods on ScanRefer, Nr3D, and Sr3D. † denotes that we reevaluate BUTD-DETR in a fair setting since its reported performance ignores some challenging samples. We only evaluate BUTD-DETR and its variants on ScanRefer because it needs to regress the target bounding box while models in the Nr3D/Sr3D setting identify the target among GT boxes. We modify our 3DOGSFormer accordingly to evaluate on Nr3D/Sr3D by replacing the 3D scene encoder with a GT box encoder and modifying the decoders to select the target GT boxes. The comparison results reveal some interesting points. (1) All 3D DOG methods outperform 3D single-object grounding methods with a clear margin, and even a simple BS strategy can improve the SOTA single-object grounding models. It is due to the fact that we can access more information about the densely referred objects in the 3D DOG setting to design our systems for more accurate grounding, which verifies the superiority of the proposed 3D DOG setting. (2) The  $\{\cdot\}$ +DepNet methods achieve clearly superior results than their base models and the  $\{\cdot\}$ +BS methods. This is because the additional global reasoning module enables the base model to capture the paragraph-level information while the SOTA single grounding methods can only conduct contextual modeling of a single target object. The performance gap validates the necessity of jointly modeling multiple target objects in a paragraph in 3D DOG. (3) Our 3DOGSFormer outperforms all baselines, especially the competitive  $\{\cdot\}$ +DepNet methods, with a significant margin, which suggests that the proposed 3DOGSFormer framework can efficiently capture the semantic relations among multiple sentences by the contextual query generator and our proposal-guided global transformer decoder can predict the target boxes precisely via modeling the 3D spatial relations among densely referred objects.

**Table 2: Ablation study of the 3DOGSFormer components.**

	CQG	LTD	GTD	Unique @0.5	Multiple @0.5	Overall @0.5
R1		✓		66.51	31.67	38.43
R2	✓	✓		68.92	36.59	42.86
R3	✓		✓	70.88	39.53	45.61
R4		✓	✓	71.85	41.64	47.50
R5	✓	✓	✓	<b>73.08</b>	<b>43.57</b>	<b>49.29</b>

**Table 3: Ablation study of the Global Transformer Decoder.**

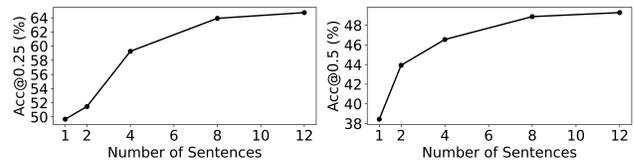
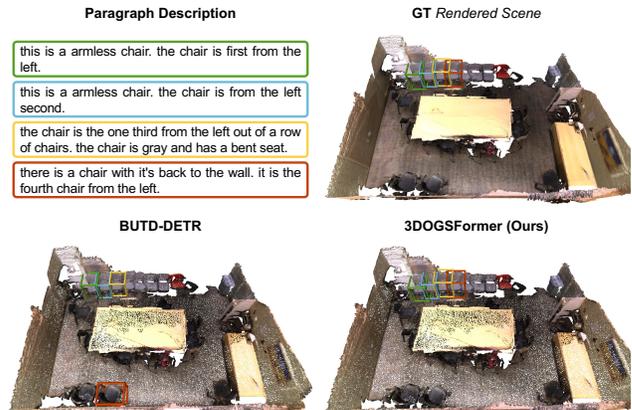
	$A^E$	$A^I$	$A^F$	Unique @0.5	Multiple @0.5	Overall @0.5
R1				70.97	39.80	45.85
R2	✓			71.81	41.13	47.08
R3		✓		71.29	40.20	46.23
R4			✓	71.67	40.52	46.56
R5	✓	✓		72.01	42.21	47.99
R6		✓	✓	72.15	42.34	48.12
R7	✓		✓	72.65	42.94	48.70
R8	✓	✓	✓	<b>73.08</b>	<b>43.57</b>	<b>49.29</b>

### 4.3 Ablation Study

To investigate the effectiveness of each component in 3DOGSFormer, we conduct ablation studies on ScanRefer. Concretely, our 3DOGSFormer includes the Contextual Query Generator (CQG), Local Transformer Decoder (LTD), and Global Transformer Decoder (GTD). We selectively discard them to generate ablation models and report the results in Table 2. From these results, we can find that the full model outperforms all ablation models, validating each component is helpful for 3D dense object grounding. The baseline R1 using only LTD achieves comparable results as BUTD-DETR since they have similar DETR-based architectures. Rows R2-R4 show that removing GTD causes the worst performance degradation, demonstrating the 3D spatial relation understanding captured by GTD is the most essential in 3D DOG. Comparing R2 to R1 and R5 to R4, we find that adding CQG can consistently improve the performance, showing the effectiveness of modeling semantic relations of multiple sentences.

Since GTD is critical in 3DOGSFormer, we further conduct detailed ablation studies on the key modules in GTD, the proposal-guided attention modules, consisting of the explicit spatial matrix  $A^E$ , implicit spatial matrix  $A^I$ , and focused region constraint matrix  $A^F$ . We discard them selectively to create ablation models. As shown in Table 3, the full model achieves the best performance, verifying all the proximity and constraint matrices are effective in 3D DOG. If only one matrix is applied, the model with  $A^E$  is the best, demonstrating the explicit spatial modeling is the most important to capture 3D spatial relations among dense objects. And if two matrices are used, the model with  $A^E$  and  $A^F$  outperforms other models, suggesting the focused region modeling is crucial in 3D spatial relation understanding and high-quality dense grounding.

Moreover, we explore the effect of the number of sentences in paragraphs as shown in Figure 5. Specifically, we evaluate the well-trained 3DOGSFormer on reconstructed ScanRefer validation

**Figure 5: Effect of the number of sentences in paragraphs.****Figure 6: An example of 3D Dense Object Grounding results.**

sets by sampling with different numbers of sentences in paragraph descriptions. We report the overall performances. As can be seen, the metrics achieve higher values when the number of sentences increases, validating that descriptions with more sentences can provide richer semantic and spatial information for our model to obtain more precise grounding.

### 4.4 Qualitative Results

We qualitatively compare 3DOGSFormer to the baseline model BUTD-DETR on ScanRefer and display a typical example in Figure 6. By capturing additional semantic and spatial relation information of dense target objects, 3DOGSFormer produces more precise results.

## 5 CONCLUSION

In this paper, we propose a novel 3D DOG task to explore the 3D dense object grounding. We devise a 3DOGSFormer model to tackle 3D DOG in a two-phase grounding pipeline. In the first phase, it initializes the grounding proposals with a local transformer decoder, where a contextual query generator is developed to efficiently capture the semantic relationships among the paragraph. In the second phase, it refines the initial proposals with a global transformer decoder, which contains newly designed proposal-guided attention layers to improve 3D spatial relation understanding via encoding pairwise spatial relations both explicitly and implicitly. Extensive experiments on three benchmarks demonstrate the significance of our 3DOGSFormer framework.

## REFERENCES

- [1] Ahmed Abdelreheem, Kyle Olszewski, Hsin-Ying Lee, Peter Wonka, and Panos Achlioptas. 2022. ScanEnts3D: Exploiting Phrase-to-3D-Object Correspondences for Improved Visio-Linguistic Models in 3D Scenes. *arXiv preprint arXiv:2212.06250* (2022).
- [2] Panos Achlioptas, Ahmed Abdelreheem, Fei Xia, Mohamed Elhoseiny, and Leonidas Guibas. 2020. ReferIt3d: Neural listeners for fine-grained 3d object identification in real-world scenes. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*. Springer, 422–440.
- [3] Eslam Bakr, Yasmeen Alsaedy, and Mohamed Elhoseiny. 2022. Look around and refer: 2d synthetic semantics knowledge distillation for 3d visual grounding. *Advances in Neural Information Processing Systems* 35 (2022), 37146–37158.
- [4] Peijun Bao, Qian Zheng, and Yadong Mu. 2021. Dense events grounding in video. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 920–928.
- [5] Daigang Cai, Lichen Zhao, Jing Zhang, Lu Sheng, and Dong Xu. 2022. 3djc: 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*. 16464–16473.
- [6] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. 2020. End-to-end object detection with transformers. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*. Springer, 213–229.
- [7] Dave Zhenyu Chen, Angel X Chang, and Matthias Nießner. 2020. Scanrefer: 3d object localization in rgb-d scans using natural language. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XX*. Springer, 202–221.
- [8] Dave Zhenyu Chen, Ronghang Hu, Xinlei Chen, Matthias Nießner, and Angel X Chang. 2022. UniT3D: A Unified Transformer for 3D Dense Captioning and Visual Grounding. *arXiv preprint arXiv:2212.00836* (2022).
- [9] Howard Chen, Alane Suhr, Dipendra Misra, Noah Snively, and Yoav Artzi. 2019. Touchdown: Natural language navigation and spatial reasoning in visual street environments. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12538–12547.
- [10] Shaoyu Chen, Jiemin Fang, Qian Zhang, Wenyu Liu, and Xinggang Wang. 2021. Hierarchical aggregation for 3d instance segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 15467–15476.
- [11] Shizhe Chen, Pierre-Louis Guhur, Makarand Tapaswi, Cordelia Schmid, and Ivan Laptev. 2022. Language Conditioned Spatial Relation Reasoning for 3D Object Grounding. In *NeurIPS 2022-36th Conference on Neural Information Processing Systems*.
- [12] Sijin Chen, Hongyuan Zhu, Xin Chen, Yinjie Lei, Tao Chen, and Gang Yu. 2023. End-to-End 3D Dense Captioning with Vote2Cap-DETR. *arXiv preprint arXiv:2301.02508* (2023).
- [13] Zhenyu Chen, Ali Gholami, Matthias Nießner, and Angel X Chang. 2021. Scan2cap: Context-aware dense captioning in rgb-d scans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3193–3203.
- [14] Zhenfang Chen, Peng Wang, Lin Ma, Kwan-Yee K Wong, and Qi Wu. 2020. Cops-ref: A new dataset and task on compositional referring expression comprehension. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10086–10095.
- [15] Marcella Cornia, Matteo Stefanini, Lorenzo Baraldi, and Rita Cucchiara. 2020. Meshed-memory transformer for image captioning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 10578–10587.
- [16] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. 2017. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 5828–5839.
- [17] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiuhua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*.
- [18] Mingtao Feng, Zhen Li, Qi Li, Liang Zhang, Xiangdong Zhang, Guangming Zhu, Hui Zhang, Yaonan Wang, and Ajmal Mian. 2021. Free-form description guided 3d visual graph network for object grounding in point cloud. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 3722–3731.
- [19] Dailan He, Yusheng Zhao, Junyu Luo, Tianrui Hui, Shaofei Huang, Aixi Zhang, and Si Liu. 2021. Transrefer3d: Entity-and-relation aware transformer for fine-grained 3d visual grounding. In *Proceedings of the 29th ACM International Conference on Multimedia*. 2344–2352.
- [20] Pin-Hao Huang, Han-Hung Lee, Hwann-Tzong Chen, and Tyng-Luh Liu. 2021. Text-guided graph neural networks for referring 3d instance segmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 1610–1618.
- [21] Shijia Huang, Yilun Chen, Jiaya Jia, and Liwei Wang. 2022. Multi-view transformer for 3d visual grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 15524–15533.
- [22] Ayush Jain, Nikolaos Gkanatsios, Ishita Mediratta, and Katerina Fragkiadaki. 2022. Bottom up top down detection transformers for language grounding in images and point clouds. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXVI*. Springer, 417–433.
- [23] Li Jiang, Hengshuang Zhao, Shaoshuai Shi, Shu Liu, Chi-Wing Fu, and Jiaya Jia. 2020. Pointgroup: Dual-set point grouping for 3d instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and Pattern recognition*. 4867–4876.
- [24] Yang Jiao, Shaoxiang Chen, Zequn Jie, Jingjing Chen, Lin Ma, and Yu-Gang Jiang. 2022. More: Multi-order relation mining for dense captioning in 3d scenes. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXV*. Springer, 528–545.
- [25] Aishwarya Kamath, Mannat Singh, Yann LeCun, Gabriel Synnaeve, Ishan Misra, and Nicolas Carion. 2021. Mdetr-modulated detection for end-to-end multi-modal understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 1780–1790.
- [26] Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of NAACL-HLT*. 4171–4186.
- [27] Daizong Liu, Xiaoye Qu, Jianfeng Dong, Pan Zhou, Yu Cheng, Wei Wei, Zichuan Xu, and Yulai Xie. 2021. Context-aware biaffine localizing network for temporal sentence grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11235–11244.
- [28] Daizong Liu, Xiaoye Qu, Xiao-Yang Liu, Jianfeng Dong, Pan Zhou, and Zichuan Xu. 2020. Jointly cross-and self-modal graph attention network for query-based moment localization. In *Proceedings of the 28th ACM International Conference on Multimedia*. 4070–4078.
- [29] Haolin Liu, Anran Lin, Xiaoguang Han, Lei Yang, Yizhou Yu, and Shuguang Cui. 2021. Refer-it-in-rgb-d: A bottom-up approach for 3d visual grounding in rgb-d images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6032–6041.
- [30] Jiaheng Liu, Tong He, Honghui Yang, Rui Su, Jiayi Tian, Junran Wu, Hongcheng Guo, Ke Xu, and Wanli Ouyang. 2022. 3D-QueryIS: A Query-based Framework for 3D Instance Segmentation. *arXiv preprint arXiv:2211.09375* (2022).
- [31] Runtao Liu, Chenxi Liu, Yutong Bai, and Alan L Yuille. 2019. Clevr-ref+: Diagnosing visual reasoning with referring expressions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 4185–4194.
- [32] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 2021. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*. 10012–10022.
- [33] Ze Liu, Zheng Zhang, Yue Cao, Han Hu, and Xin Tong. 2021. Group-free 3d object detection via transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2949–2958.
- [34] Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101* (2017).
- [35] Jiasen Lu, Jianwei Yang, Dhruv Batra, and Devi Parikh. 2016. Hierarchical question-image co-attention for visual question answering. *Advances in neural information processing systems* 29 (2016).
- [36] Junyu Luo, Jiahui Fu, Xianghao Kong, Chen Gao, Haibing Ren, Hao Shen, Huaxia Xia, and Si Liu. 2022. 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*. 16454–16463.
- [37] Junhua Mao, Jonathan Huang, Alexander Toshev, Oana Camburu, Alan L Yuille, and Kevin Murphy. 2016. Generation and comprehension of unambiguous object descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 11–20.
- [38] Ishan Misra, Rohit Girdhar, and Armand Joulin. 2021. An end-to-end transformer model for 3d object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2906–2917.
- [39] Charles R Qi, Or Litany, Kaiming He, and Leonidas J Guibas. 2019. Deep hough voting for 3d object detection in point clouds. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 9277–9286.
- [40] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. 2017. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in neural information processing systems* 30 (2017).
- [41] Yuankai Qi, Qi Wu, Peter Anderson, Xin Wang, William Yang Wang, Chunhua Shen, and Anton van den Hengel. 2020. Reverie: Remote embodied visual referring expression in real indoor environments. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9982–9991.
- [42] Yanyuan Qiao, Chaorui Deng, and Qi Wu. 2020. Referring expression comprehension: A survey of methods and datasets. *IEEE Transactions on Multimedia* 23 (2020), 4426–4440.
- [43] Hamid Rezaatofighi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese. 2019. Generalized intersection over union: A metric and a loss for bounding box regression. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 658–666.

- [44] Junha Roh, Karthik Desingh, Ali Farhadi, and Dieter Fox. 2022. Language-refer: Spatial-language model for 3d visual grounding. In *Conference on Robot Learning*. PMLR, 1046–1056.
- [45] Jiahao Sun, Chunmei Qing, Junpeng Tan, and Xiangmin Xu. 2022. Superpoint Transformer for 3D Scene Instance Segmentation. *arXiv preprint arXiv:2211.15766* (2022).
- [46] Peize Sun, Rufeng Zhang, Yi Jiang, Tao Kong, Chenfeng Xu, Wei Zhan, Masayoshi Tomizuka, Lei Li, Zehuan Yuan, Changhu Wang, et al. 2021. Sparse r-cnn: End-to-end object detection with learnable proposals. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 14454–14463.
- [47] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- [48] Thang Vu, Kookhoi Kim, Tung M Luu, Thanh Nguyen, and Chang D Yoo. 2022. Softgroup for 3d instance segmentation on point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2708–2717.
- [49] Heng Wang, Chaoyi Zhang, Jianhui Yu, and Weidong Cai. 2022. Spatiality-guided Transformer for 3D Dense Captioning on Point Clouds. In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI 2022, Vienna, Austria, 23-29 July 2022*. 1393–1400.
- [50] Yanmin Wu, Xinhua Cheng, Renrui Zhang, Zesen Cheng, and Jian Zhang. 2022. EDA: Explicit Text-Decoupling and Dense Alignment for 3D Visual and Language Learning. *arXiv preprint arXiv:2209.14941* (2022).
- [51] Antoine Yang, Antoine Miech, Josef Sivic, Ivan Laptev, and Cordelia Schmid. 2022. Tubedetr: Spatio-temporal video grounding with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16442–16453.
- [52] Zhengyuan Yang, Songyang Zhang, Liwei Wang, and Jiebo Luo. 2021. Sat: 2d semantics assisted training for 3d visual grounding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 1856–1866.
- [53] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. 2016. Modeling context in referring expressions. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part II 14*. Springer, 69–85.
- [54] Zhihao Yuan, Xu Yan, Zhuo Li, Xuhao Li, Yao Guo, Shuguang Cui, and Zhen Li. 2022. Toward Explainable and Fine-Grained 3D Grounding through Referring Textual Phrases. *arXiv preprint arXiv:2207.01821* (2022).
- [55] Zhihao Yuan, Xu Yan, Yinghong Liao, Yao Guo, Guanbin Li, Shuguang Cui, and Zhen Li. 2022. X-trans2cap: Cross-modal knowledge transfer using transformer for 3d dense captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8563–8573.
- [56] Zhihao Yuan, Xu Yan, Yinghong Liao, Ruimao Zhang, Sheng Wang, Zhen Li, and Shuguang Cui. 2021. 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*. 1791–1800.
- [57] Lichen Zhao, Daigang Cai, Lu Sheng, and Dong Xu. 2021. 3DVG-Transformer: Relation modeling for visual grounding on point clouds. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2928–2937.
- [58] Yufeng Zhong, Long Xu, Jiebo Luo, and Lin Ma. 2022. Contextual Modeling for 3D Dense Captioning on Point Clouds. *arXiv preprint arXiv:2210.03925* (2022).
- [59] Benjin Zhu, Zhe Wang, Shaoshuai Shi, Hang Xu, Lanqing Hong, and Hongsheng Li. 2022. ConQueR: Query Contrast Voxel-DETR for 3D Object Detection. *arXiv preprint arXiv:2212.07289* (2022).
- [60] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. 2021. Deformable DETR: Deformable Transformers for End-to-End Object Detection. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3–7, 2021*.