

# TransRefer3D: Entity-and-Relation Aware Transformer for Fine-Grained 3D Visual Grounding

Dailan He\*

hedailan@buaa.edu.cn  
School of Computer Science and Engineering, Beihang University  
Beijing, China

Yusheng Zhao\*

zhaoyusheng@buaa.edu.cn  
School of Computer Science and Engineering, Beihang University  
Beijing, China

Junyu Luo\*

luojunyu@buaa.edu.cn  
School of Computer Science and Engineering, Beihang University  
Beijing, China

Tianrui Hui

huitianrui@iie.ac.cn  
Institute of Information Engineering,  
Chinese Academy of Sciences  
Beijing, China

Shaofei Huang

huangshaofei@iie.ac.cn  
Institute of Information Engineering,  
Chinese Academy of Sciences  
Beijing, China

Aixi Zhang

aixi.zhax@alibaba-inc.com  
Alibaba Group  
Beijing, China

Si Liu<sup>†</sup>

liusi@buaa.edu.cn  
Institute of Artificial Intelligence  
Beijing, China

## ABSTRACT

Recently proposed fine-grained 3D visual grounding is an essential and challenging task, whose goal is to identify the 3D object referred by a natural language sentence from other distractive objects of the same category. Existing works usually adopt dynamic graph networks to indirectly model the intra/inter-modal interactions, making the model difficult to distinguish the referred object from distractors due to the monolithic representations of visual and linguistic contents. In this work, we exploit Transformer for its natural suitability on permutation-invariant 3D point clouds data and propose a *TransRefer3D* network to extract entity-and-relation aware multimodal context among objects for more discriminative feature learning. Concretely, we devise an Entity-aware Attention (EA) module and a Relation-aware Attention (RA) module to conduct fine-grained cross-modal feature matching. Facilitated by co-attention operation, our EA module matches visual entity features with linguistic entity features while RA module matches pair-wise visual relation features with linguistic relation features, respectively. We further integrate EA and RA modules into an Entity-and-Relation aware Contextual Block (ERCB) and stack several ERCBs to form our TransRefer3D for hierarchical multimodal context modeling. Extensive experiments on both Nr3D and Sr3D datasets demonstrate that our proposed model significantly

outperforms existing approaches by up to **10.6%** and claims the new state-of-the-art performance. To the best of our knowledge, this is the first work investigating Transformer architecture for fine-grained 3D visual grounding task.

## CCS CONCEPTS

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

## KEYWORDS

3D visual grounding, transformer, entity attention, relation attention

## ACM Reference Format:

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 (MM '21), October 20–24, 2021, Virtual Event, China*. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3474085.3475397>

## 1 INTRODUCTION

Visual-linguistic modeling, which aims to connect vision and language, is an essential task in multimedia understanding and embodied AI. Specific tasks like visual question answering [24, 32, 33], image captioning [2, 9, 14], image-text matching [20, 21, 29], visual relation detection [8] and scene graph generation [36] have been widely studied in the past few years. In this paper, we focus on another important task named visual grounding, which aims to locate referred objects (referents) given the input referring expressions. However, most of the existing works on visual grounding have been conducted on 2D images, which might fail to capture the full information in our 3D reality.

Recently, approaches and benchmarks for 3D visual-linguistic tasks emerge and attract attention [12, 25–27]. As an essential 3D

\*Equal contribution.

<sup>†</sup>Corresponding author.

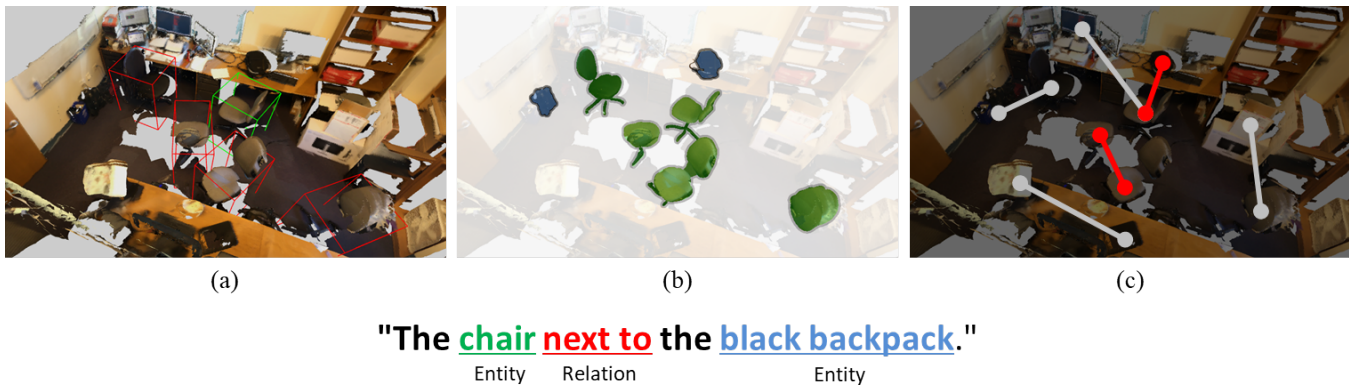
Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components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

MM '21, October 20–21, 2021, Chengdu, China

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-8651-7/21/10...\$15.00

<https://doi.org/10.1145/3474085.3475397>



**Figure 1: An example of 3D visual grounding (a). According to the input sentence (referring expression) the model is required to locate the referred object (the chair in the given case) from the 3D scene. To perform that, the model needs to identify the *entity* (b) and *relation* (c) mentioned in the referring expression and match them with the corresponding visual elements to obtain a multimodal context in an explicit or implicit way.**

visual-linguistic task, 3D visual grounding is to locate objects in 3D point cloud scenes. Two important benchmarks, ScanRefer [5] and ReferIt3D [1], have been proposed based on ScanNet dataset [10]. While the goal of ScanRefer is to predict 3D bounding boxes for the referents given natural language description, ReferIt3D are more focused on the following two aspects: (1) The objects in the scenes are annotated with **fine-grained** labels and referring expressions. For instance, given the description "choose the armchair in the center of the room", the model is required to locate only armchairs but not other types of chairs like folded chairs or office-chairs; (2) Each scene contains **multiple** objects of the same fine-grained category as the referent. As a result, the model cannot simply rely on object classification. Therefore, ReferIt3D requires more fine-grained reasoning of 3D objects and referring expressions, which is more challenging than ScanRefer. In this work, we mainly focus on this fine-grained visual grounding benchmark.

Based on the above benchmarks, researchers have proposed a few 3D visual grounding methods. ReferIt3DNet [1] adopts a dynamic graph network (DGCNN [25]) to model the multimodal context. InstanceRefer [34] utilizes three parallel modules to match visual and linguistic features from different aspects. However, all of these methods fall short of explicit modeling of entities and their relationships in the cross-modality context. Moreover, some of the previous methods adopt simple fusion of visual feature and linguistic feature. Consequently, the performance of these methods is relatively weak on ReferIt3D benchmark.

We have two observations with a deeper analysis of previous works and the 3D visual grounding task. The first is **the similarity between point cloud encoders and Transformers** [28, 38]. In the 3D visual grounding task, the scene is represented by a set of point cloud features. These features are *permutation-invariant*, and preserving this property has been a basic guidance of designing point cloud encoders. On the other side, the multi-head self-attention in Transformers preserves this property by its nature. Since Transformers have achieved phenomenal success in both natural language processing [11] and computer vision [13], we are highly motivated to introduce Transformers in our task.

Another observation is the importance of **combining entity and relation information** in visual-linguistic modeling. *Entity* information refers to the feature within an instance, like the category, color, or size. *Relation* information refers to the feature among objects, including spatial relation, comparative relation, etc. Consider the scene "The chair next to the black backpack" (shown in Figure 1), "the chair" and "the black backpack" are linguistic entities, and the corresponding objects are visual entities. Similarly, "next to" is a linguistic relation, and the corresponding spatial relation is implied in the scene. As we can see, the entity or relation information alone is not enough for modeling the multimodal context. Therefore, combining two types of information is essential in our 3D visual grounding task.

Based on these two observations, we propose a Transformer architecture named *TransRefer3D* to better comprehend the entity-and-relation multimodal context for fine-grained 3D visual grounding. Concretely, we utilize Self-Attention modules (SA) to capture intra-modal contextual information. For inter-modal context, we propose an *Entity-aware Attention* (EA) module and a *Relation-aware Attention* (RA) module to match entity and relation features in both modalities. The EA module performs matching between object features and linguistic entity features via co-attention. The RA module first extracts visual relation features of object-object pairs via an asymmetric operator. It then performs cross-modal attention between visual relation features and linguistic relation features. With SA, EA and RA, we integrate them into a unified module called Entity-and-Relation aware Contextual Block (ERCB) for multimodal context modeling. Finally, several ERCBs are stacked together to form our *TransRefer3D* model that captures the hierarchical feature in the cross-modality context.

To summarize, we contribute to the research on 3D visual grounding from three aspects:

- To the best of our knowledge, we are the first to introduce the Transformer architecture to achieve a better cross-modal feature representation for the 3D visual grounding task.

- We propose entity-and-relation aware attention for multi-modal context comprehension and use it for fine-grained 3D visual grounding task.
- The proposed TransRefer3D model achieves state-of-the-art performance on ReferIt3D. Compared with existing cutting-edge approaches, the test accuracy improves up to 10.4% on Nr3D (with Sr3D+ augmentation).

## 2 RELATED WORKS

### 2.1 2D Visual Grounding

Visual grounding, also known as phrase grounding, aims to localize descriptive words in images. It has been widely studied in the past several years [17, 30, 31, 37]. Karpathy *et al.* [18] proposed to map image fragments and sentence fragments into a common space to match images and sentences. Chen *et al.* [7] proposed a Query-guided Regression Network with Context policy (QRC net) and introduce reinforcement learning techniques to achieve better localization. Huang *et al.* [16], proposed a Cross-Modal Progressive Comprehension module (CMPC) and a Text-Guided Feature Exchange module (TGFE) to segment referred entities in the image. Though researchers have made significant progress on 2D visual grounding, the research on 3D visual grounding is still relatively preliminary. In this work, we mainly focus on fine-grained 3D visual grounding task.

### 2.2 3D Visual Grounding

3D visual grounding is an emerging research topic. Chen *et al.* [5] released a new benchmark of 3D object localization with natural language descriptions. They proposed ScanRefer to combine visual detection and language encoder for joint inference. Achlioptas *et al.* [1] released another benchmark called ReferIt3D consisting of Nr3D, Sr3D and Sr3D+ datasets, which assumes the segmented object instances are already given. However, its task is more fine-grained where each referred object has distractors of the same category. Text-guided Graph Neural Network (TGNN [15]) are proposed for referring instance segmentation on point cloud. Yuan *et al.* [34] proposed InstanceRefer which utilizes instance attribute, relation and localization perceptions for 3d visual grounding. Different from previous works, we exploit the Transformer architecture and propose Entity-aware Attention (EA) module and Relation-aware Attention (RA) module to build the basic blocks of Transformer to achieve a finer matching between visual and linguistic features.

### 2.3 Attention Mechanism and Transformer

Transformers have been widely used in natural language processing [11, 23, 28]. Recently a lot of works also apply them on high-level vision tasks [4, 13], low-level vision tasks [6] and graph tasks [35] because of its strong ability to recognize long-term dependency.

A Transformer is composed of several Transformer layers stacked together. Each Transformer layer consists of a multi-head self-attention (MSA) and a feed forward network. MSA is based on the self-attention mechanism. Given three feature matrices: queries

$Q \in \mathbb{R}^{n \times d}$ , keys and values  $K, V \in \mathbb{R}^{m \times d}$ , a single-head self-attention is formulated as:

$$\text{Attn}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V \quad (1)$$

introducing  $H$  parallel attention layers as different transformation heads to enlarge the model capacity, we can further obtain a multi-head attention:

$$\begin{aligned} \text{MSA}(Q, K, V) &= W_{out} [f_1, f_2, \dots, f_H] \\ f_i &= \text{Attn}\left(W_q^{(i)}Q, W_k^{(i)}K, W_v^{(i)}V\right) \end{aligned} \quad (2)$$

where  $W_q^{(i)}, W_k^{(i)}, W_v^{(i)} \in \mathbb{R}^{d_h \times d}$  are linear transformation that projects  $Q, K, V$  for the  $i$ -th head and  $W_{out} \in \mathbb{R}^{h \times H d_h}$  for the output. After the multi-head attention follows a feed forward network, which is a multi-layer perceptron (MLP). Moreover, the Transformer also utilize Layer Normalization and residual connections.

### 2.4 Co-Attention for Multimodal Learning

Multimodal learning requires a full understanding of the contents in each modality and, more importantly, learning the interactions between modalities. Co-attention mechanism [24] is proposed to modeling the image attention and question attention in VQA. Yu *et al.* [32] proposed a dense co-attention model to make full interactions between modalities and address the previous issue that each modality learns separate attention distributions. In 3D visual grounding tasks, the context of entity-and-relation has not been fully researched yet. In this paper, we expand the co-attention mechanism to guide context modeling in 3D visual grounding task. This is the first work investigating the co-attention mechanism on the 3D visual-linguistic tasks to the best of our knowledge.

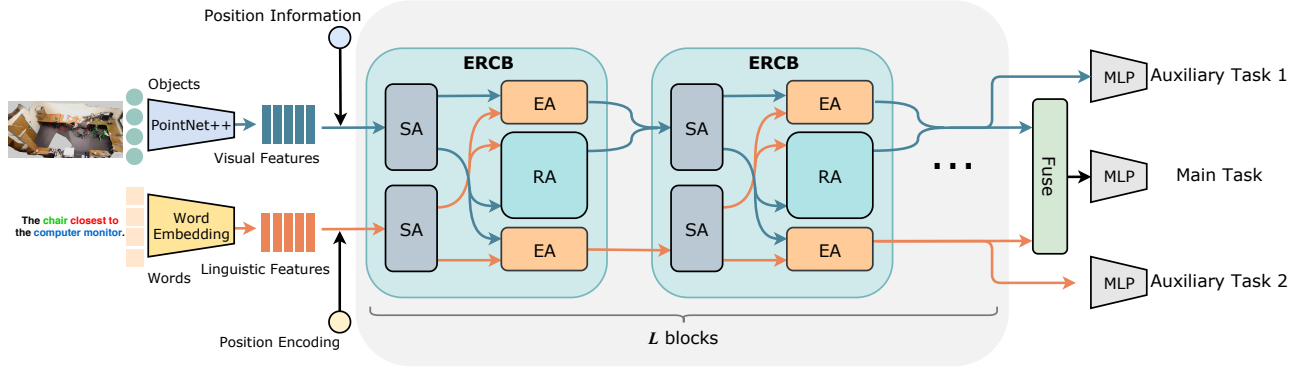
## 3 METHODOLOGY

### 3.1 Overview

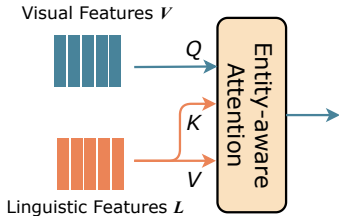
According to the above sections, scenes contain two types of information: entities' attributes and relations among entities. To better model the multimodal context for fine-grained 3D visual grounding, we propose two cross-modal attention mechanism: Entity-aware Attention (EA) and Relation-aware Attention (RA). Furthermore, we develop an Entity-and-Relation aware Contextual Block (ERCB) that integrate those attention mechanisms as a multimodal encoding unit. By stacking multiple ERCBs, we obtain the proposed TransRefer3D model.

Figure 2 shows the framework of our model. Following previous work [1], we first extract object features with a shared PointNet++ [27] and vectorize the input words with an embedding layer. After that, we fuse position information to both visual and linguistic features and feed them into the proposed TransRefer3D.

We use the same training objective as ReferIt3D [1] consisting of a main task and two auxiliary tasks. For  $N$  input objects, the main training task is a  $N$ -classification supervised by a cross-entropy loss  $\mathcal{L}_{main}$ . To obtain the features well embedded to represent  $K$  categories, we also adopt  $K$ -classification on object features and aggregated linguistic features, supervised by cross entropy loss  $\mathcal{L}_{obj}$  and  $\mathcal{L}_{lang}$ . The total loss function is  $\mathcal{L} = \mathcal{L}_{main} + \lambda_{obj}\mathcal{L}_{obj} + \lambda_{lang}\mathcal{L}_{lang}$  where we set  $\lambda_{obj} = \lambda_{lang} = 0.5$  in all of our experiments.



**Figure 2: Architecture of the proposed TransRefer3D.** The model is composed of several Entity-and-Relation aware Contextual Blocks (ERCBs) stacked together. Each block contains self-attention (SA), Entity-aware Attention (EA) and Relation-aware Attention (RA). The model predicts the referred object (Main Task) as well as an utterance classification of the referent (Auxiliary Task 1) and an object classification (Auxiliary Task 2) for better feature extraction. The blue arrows indicates the flow of visual features, while the red arrows showcases the flow of linguistic features.



**Figure 3: The Entity-aware Attention module (EA).** The figure illustrates the language-guided visual entity-aware attention, where visual features serve as queries and linguistic features are keys and values of multi-head attention. The vision-guided linguistic entity-aware attention is similar, with linguistic features as the queries and visual features as keys and values.

### 3.2 EA: Entity-aware Attention Module

Our goal is to find out the referent according to the referring expression. Usually, the referent and other objects in the scene get mentioned in the referring expression as some keywords. We name all objects appearing in the scene and their corresponding keywords as the *entities* in each modality. First of all, we establish an Entity-aware Attention (EA) module for finely understanding the multimodal context.

The importance of each entity for locating the referent is different. In a ScanNet indoor scene, often there are more than 20 objects, but most of them make a limited contribution to the visual grounding task. To better understand the referring expression, a strategy is adapt to match the entity objects with corresponding words and suppress relatively useless features. By this, the model can finely identify and establish correlations on those important entities.

Inspired by self-attention mechanism and co-attention widely adopted in 2D vision-language approaches [16, 32], we propose

EA module to match those features to model an entity aware multimodal context. In a 3D scene, the  $n$  objects are represented by  $d$ -dimension features. And the referring expression is handled to  $m$  word-level features. Given two feature sets in different modality  $X \in \mathbb{R}^{n \times d}$  and  $Y \in \mathbb{R}^{m \times d}$ , the EA module matches  $Y$  to  $X$ :

$$EA(X, Y) = \text{Attn}(X, Y) = \text{softmax}\left(\frac{XY^T}{\sqrt{d}}\right) Y \quad (3)$$

which is a multimodal generalization of attention. Like self-attention, we further extend EA to a multi-head form by replacing  $\text{Attn}(\cdot)$  with a multi-head version.

The proposed EA module does symmetrical matchings to visual and word-level linguistic features for better modeling entities in both modals.

### 3.3 RA: Relation-aware Attention Module

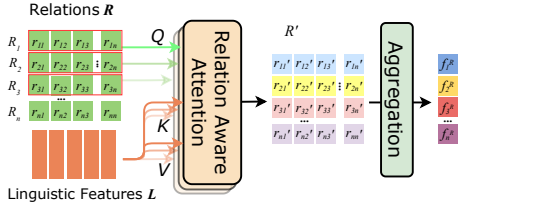
The EA module can match entity information in word features and 3D object features. However, only focusing on the entities is insufficient for complicated visual-linguistic reasoning. For example, consider the scene with multiple objects of the same category as the referent. All of them match the corresponding word feature, so their corresponding Entity-aware Attention scores can be similar. Consequently, it can be challenging for the model to distinguish the referent from others.

To address this problem, it is important to let model recognize those relational words (e.g. "next to") that appear in the referring expression and enhance entity features according to their relations. We establish a representation of abstract entity-to-entity relation for visual modality and develop a cross-modal Relation-aware Attention (RA) for finer context modeling.

Inspired by DGCNN, we define the relation representation between two entity features  $f_i, f_j$  as an asymmetric form:

$$r_{ij} = r(f_i, f_j) = H_{\Theta}(f_i - f_j) \quad (4)$$

where  $H_{\Theta}$  denotes a nonlinear layer.



**Figure 4: The Relation-aware Attention module (RA).** In RA, the visual features are first extracted and then matched with the language feature. The relation features serve as queries, and the linguistic features are keys and values. After the cross-modal attention, the language-guided relation features are then aggregated.

Like entity features, each relation features  $r_{ij} \in \mathbb{R}^d$  can also get matched with corresponding linguistic features. For each entity feature  $f_i$ , we calculate its relation features and introduce co-attention on them to obtain an enhanced cross-modal context of  $f_i$ :

$$\begin{aligned}
 R_i &= [r_{i1}, r_{i2}, \dots, r_{in}] \\
 R_i' &= \text{Attn}(R_i, L, L) \\
 f_i^R &= \text{Agg}(R_i') \\
 F_i^R &= [f_1^R, f_2^R, \dots, f_n^R]
 \end{aligned} \tag{5}$$

where  $L$  is the linguistic feature vector and  $\text{Agg}(\cdot)$  denotes a channel-wise aggregation operator gathering the  $n$  fused  $d$ -dimension features in  $R_i'$  to  $f_i^R \in \mathbb{R}^d$ . Figure 4 shows the diagram of RA. Hence for each  $f_i$  we can obtain a corresponding  $f_i^R$  as its relation-aware multimodal context representation.

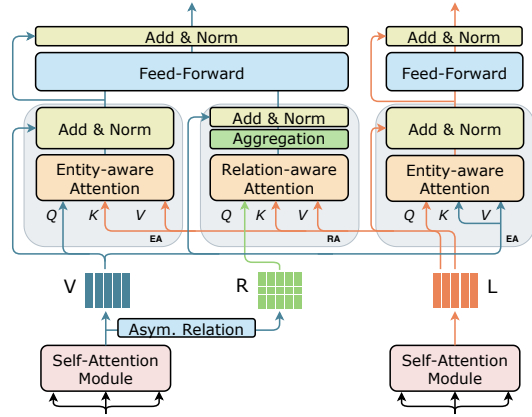
We adopt a binary relation representation calculated on each pair of entities (eq. 4). Assuming that those more complicated relations among multiple entities can be described as a set of binary relations, we stack RA modules in our model to achieve a hierarchical and progressive context comprehension with respect to relation information.

### 3.4 ERCB: Entity-and-Relation Aware Contextual Block

EA and RA can model two essential information in scenes: entities and relations. By composing EA and RA, we propose an Entity-and-relation Aware Contextual Block (ERCB), as shown in Figure 5. ERCB can form a cross-modal context comprehension with respect to both entity and relation information for visual grounding.

In each block, we first feed the visual and linguistic features to self-attention layers. Self-attention layers can effectively extract features in each modality and enhance the performance of cross-modal transformations.

Then, we match the visual features with word-level features with respect to both entity information and relation information. We calculate two types of attention representation  $F^E$  and  $F^R$  in



**Figure 5: The architecture of the proposed ERCB.** The input visual and linguistic features are first fed into two self-attention modules. Then visual entity features and linguistic entity features are matched through Entity-aware Attention (EA). The visual relation features are first extracted and then matched with linguistic features through Relation-aware Attention. Finally, the enhanced features are fused and processed by feed-forward networks.

parallel and fuse them to obtain the joint features  $f_i^J$ :

$$\begin{aligned}
 g_i^E &= \text{LN}(f_i^E + f_i) \\
 g_i^R &= \text{LN}(f_i^R + f_i) \\
 f_i^J &= \text{FFN}([g_i^E, g_i^R])
 \end{aligned} \tag{6}$$

where  $f_i^J \in \mathbb{R}^d$ .  $f_i^E$  and  $f_i^R$  are the  $i$ -th feature of  $F^E$  and  $F^R$ , respectively. FFN is a feed-forward network and LN denotes layer normalization [3]. Further performing an add-and-norm transformation, we can obtain the enhanced multimodal feature  $f_i^J$ :

$$f_i^J = \text{LN}(f_i^J + g_i^E) \tag{7}$$

We also enhance linguistic features with EA guided by visual features. For each attention layer, a feed-forward network is followed to introduce non-linearity. By stacking those ERCBs, we finally obtain a hierarchical Transformer for fine-grained multimodal context comprehension.

In Section 4.3, we will discuss the ability of each module and the structure of our model. Experiments show that TransRefer3D and its component EACB are much more powerful in modeling multimodal context than other methods.

## 4 EXPERIMENTS

In this section, we present experiments to demonstrate the effectiveness of the proposed TransRefer3D compared with other methods. We also conduct quantitative/qualitative studies to analyze our method in details.

### 4.1 Implementation details and Datasets

For a fair comparison, we follow the same settings as ReferIt3D [1]. We train the model until convergence with Adam optimizer [19]

Arch.	Overall	Easy	Hard	View-dep.	View-indep.
<b>Nr3D</b>					
ReferIt3DNet [1]	35.6% ± 0.7%	43.6% ± 0.8%	27.9% ± 0.7%	32.5% ± 0.7%	37.1% ± 0.8%
TGNN [15]	37.3% ± 0.3%	44.2% ± 0.4%	30.6% ± 0.2%	35.8% ± 0.2%	38.0% ± 0.3%
InstanceRefer [34]	38.8% ± 0.4%	46.0% ± 0.5%	31.8% ± 0.4%	34.5% ± 0.6%	41.9% ± 0.4%
TransRefer3D (ours)	<b>42.1% ± 0.2%</b>	<b>48.5% ± 0.2%</b>	<b>36.0% ± 0.4%</b>	<b>36.5% ± 0.6%</b>	<b>44.9% ± 0.3%</b>
<b>Nr3D w/ Sr3D</b>					
ReferIt3DNet	37.2% ± 0.3%	44.0% ± 0.6%	30.6% ± 0.3%	33.3% ± 0.6%	39.1% ± 0.2%
TransRefer3D (ours)	<b>47.2% ± 0.3%</b>	<b>55.4% ± 0.5%</b>	<b>39.3% ± 0.5%</b>	<b>40.3% ± 0.4%</b>	<b>50.6% ± 0.2%</b>
<b>Nr3D w/ Sr3D+</b>					
ReferIt3DNet	37.6% ± 0.4%	45.4% ± 0.6%	30.0% ± 0.4%	33.1% ± 0.5%	39.8% ± 0.4%
TransRefer3D (ours)	<b>48.0% ± 0.2%</b>	<b>56.7% ± 0.4%</b>	<b>39.6% ± 0.2%</b>	<b>42.5% ± 0.2%</b>	<b>50.7% ± 0.4%</b>

**Table 1: The performance of TransRefer3D on Nr3D trained with or without Sr3D/Sr3D+, compared with previous works.**

and a learning rate of 0.0005. The default feature dimension ( $d$ ) is 128 and the default batch size is 32. We set the depth of ERCBs ( $L$ ) to 4 and use 4-head attentions if not particularly indicated.

We conduct our experiments on three datasets, Nr3D, Sr3D and Sr3D+ [1]. The details of these datasets are listed as follows:

- **Nr3D:** Nr3D (Natural Reference in 3D) consists of 41.5K human descriptions collected using a referring game (ReferIt Game). It describes objects in 707 ScanNet scenes. In each scene, there are multiple objects of the same category as the referent.
- **Sr3D:** Sr3D (Spatial Reference in 3D) contains 83.5K synthetic descriptions. It categorizes spatial relations into 5 types: horizontal proximity, vertical proximity, between, allocentric and support [1], and then generates descriptions using language templates. Similar with Nr3D, there are always more than one objects of the same category as the referred object.
- **Sr3D+:** Sr3D+ is an augmentation of Sr3D. It uses the same method as Sr3D to generate synthetic descriptions. The difference is the referent could be the only object belonging to a fine-grained category in each scene.

Sr3D and Sr3D+ are designed to enhance the performance of the model on Nr3D [1]. The model is first jointly trained on both Nr3D training set and Sr3D/Sr3D+ and then evaluated on Nr3D test set. Nevertheless, for the completeness of experiments, we also train and test our model only on Sr3D.

## 4.2 Quantitative Performance on Nr3D and Sr3D/Sr3D+

In this section, we report the quantitative performance of our model on Nr3D, Sr3D and Sr3D+.

### 4.2.1 Evaluation on Nr3D.

Firstly, we train our model on Nr3D training set and test the model on Nr3D test set. As shown in Table 1 (upper part), TransRefer3D outperforms previous models by a large margin. Besides, we also measure our performance on different types of scenarios, defined in [1]:

- **Easy vs. Hard:** easy cases are those having only one distractor (the object of the same class as the referent but not

the referent itself) and hard cases are those having multiple distractors.

- **View dependency:** view dependent cases (View-dep.) are those whose descriptions ask the viewers to place themselves facing certain objects, and other cases are view independent (View-indep.) .

The results show that TransRefer3D performs especially well in hard cases, suggesting that our model can better understand the complicated referring context.

Then, we also train our model on Nr3D jointly with Sr3D/Sr3D+, and test the model on Nr3D test set. The results are shown in Table 1 (middle part and bottom part). The results indicate that TransRefer3D may perform even better on larger datasets. Compared with training only on Nr3D, the performance of ReferIt3DNet improves by 1.6% by joint training on Nr3D and Sr3D. However, we find that this improvement enlarges to 5.1% on TransRefer3D. A similar trend is also observed on our model by joint training with Sr3D+. One possible explanation for this phenomenon is that training Transformers require relatively more data. Recent studies on Transformers designed for other tasks also report similar conclusions [11, 13]. Therefore, TransRefer3D may have the potential to achieve better performance on visual grounding if a larger dataset is available in the future.

### 4.2.2 Evaluation on Sr3D.

We also train and test the proposed model only on Sr3D, and the results are shown in Table 2. The results show that our model outperforms all the previous works with more than 9.4% accuracy improvement. Since the data in Sr3D is generated from language templates that describe spatial relations, it is relatively easy for TransRefer3D that can explicitly model entity and relation in a cross-modality context.

## 4.3 Ablation Study

In this subsection, we further discuss the effectiveness of each module in TransRefer3D and the architecture design.

To verify the effectiveness of each module in the proposed model, we train and test TransRefer3D without each of the modules: self-attention (SA), vision-guided linguistic entity aware attention (EA,  $V \rightarrow L$ ), language-guided visual entity aware attention (EA,  $L \rightarrow V$ )

Arch.	Overall	Easy	Hard	View-dep.	View-indep.
ReferIt3DNet [1]	40.8% $\pm$ 0.2%	44.7% $\pm$ 0.1%	31.5% $\pm$ 0.4%	39.2% $\pm$ 1.0%	40.8% $\pm$ 0.1%
TGNN [15]	45.0% $\pm$ 0.2%	48.5% $\pm$ 0.2%	36.9% $\pm$ 0.5%	45.0% $\pm$ 1.1%	45.0% $\pm$ 0.2%
InstanceRefer [34]	48.0% $\pm$ 0.3%	51.1% $\pm$ 0.2%	40.5% $\pm$ 0.3%	45.4% $\pm$ 0.9%	48.1% $\pm$ 0.3%
TransRefer3D (ours)	57.4% $\pm$ 0.2%	60.5% $\pm$ 0.2%	50.2% $\pm$ 0.2%	49.9% $\pm$ 0.6%	57.7% $\pm$ 0.2%

Table 2: The performance of TransRefer3D trained and evaluated on Sr3D only, compared with previous works.

Arch.	Overall	Easy	Hard	View-dep.	View-indep.
w/o SA	38.6% $\pm$ 0.2%	45.5% $\pm$ 0.6%	32.0% $\pm$ 0.4%	33.3% $\pm$ 0.7%	41.3% $\pm$ 0.3%
w/o EA (V $\rightarrow$ L)	41.7% $\pm$ 0.2%	48.7% $\pm$ 0.5%	34.9% $\pm$ 0.1%	37.7% $\pm$ 0.5%	43.7% $\pm$ 0.2%
w/o EA (L $\rightarrow$ V)	41.1% $\pm$ 0.4%	48.9% $\pm$ 0.6%	33.6% $\pm$ 0.4%	37.9% $\pm$ 0.2%	42.7% $\pm$ 0.5%
w/o RA	40.8% $\pm$ 0.1%	47.9% $\pm$ 0.2%	34.0% $\pm$ 0.1%	36.7% $\pm$ 0.5%	42.9% $\pm$ 0.1%
EA + RA (stacked)	41.4% $\pm$ 0.4%	47.7% $\pm$ 0.6%	35.3% $\pm$ 0.3%	37.0% $\pm$ 0.6%	43.5% $\pm$ 0.4%
EA + RA (TransRefer3D)	42.1% $\pm$ 0.2%	48.5% $\pm$ 0.2%	36.0% $\pm$ 0.4%	36.5% $\pm$ 0.6%	44.9% $\pm$ 0.3%

Table 3: Ablation study of our model. We train and test the proposed models without self-attention (SA) and the proposed modules (vision-guided EA, language guided EA and RA). We also try different module architectures EA and RA modules stacked together or in parallel (TransRefer3D).

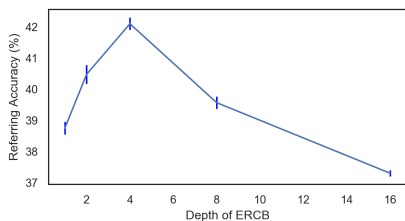


Figure 6: Referring Accuracy with different depth of ERCB.

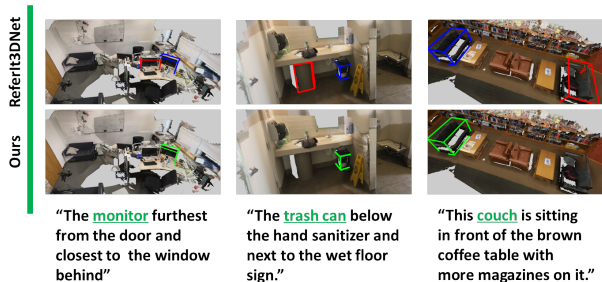


Figure 7: Visualization of the visual grounding results compared with ReferIt3DNet. The blue bounding boxes are the targets, the red ones are misclassifications, and the green ones are the correct classifications.

and relation aware attention (RA). The results are shown in Table 3. The model without SA performs worst in all the experiments. This indicates that the Transformer architecture is suitable for our 3D visual grounding task and that the contextual awareness within a single modal is the basis of cross-modal understanding. Furthermore, removing all of the proposed modules hinders the performance, demonstrating the effectiveness of the proposed EA and RA modules.

Moreover, we train and test our model with different architectures. By default, the proposed ERCB consists of an EA module and a RA module in parallel. We have also tried stacking EA and RA modules, in which the data is first fed into EA module and then RA module. A possible explanation for this is that the entity information and relation information are relatively independent, and a parallel structure might be more reasonable.

As for the depth of our model, we tested TransRefer3D with different depth of ERCB layers on Nr3D dataset. As shown in Figure 6, 4 layers of ERCB achieve the best performance. Shallower models with fewer layers might fail to represent the complicated relations as analyzed in Section 3.3. Simultaneously, deeper Transformers have been observed hard to train because of unbalanced gradients and amplification effect [22].

## 4.4 Qualitative Analysis

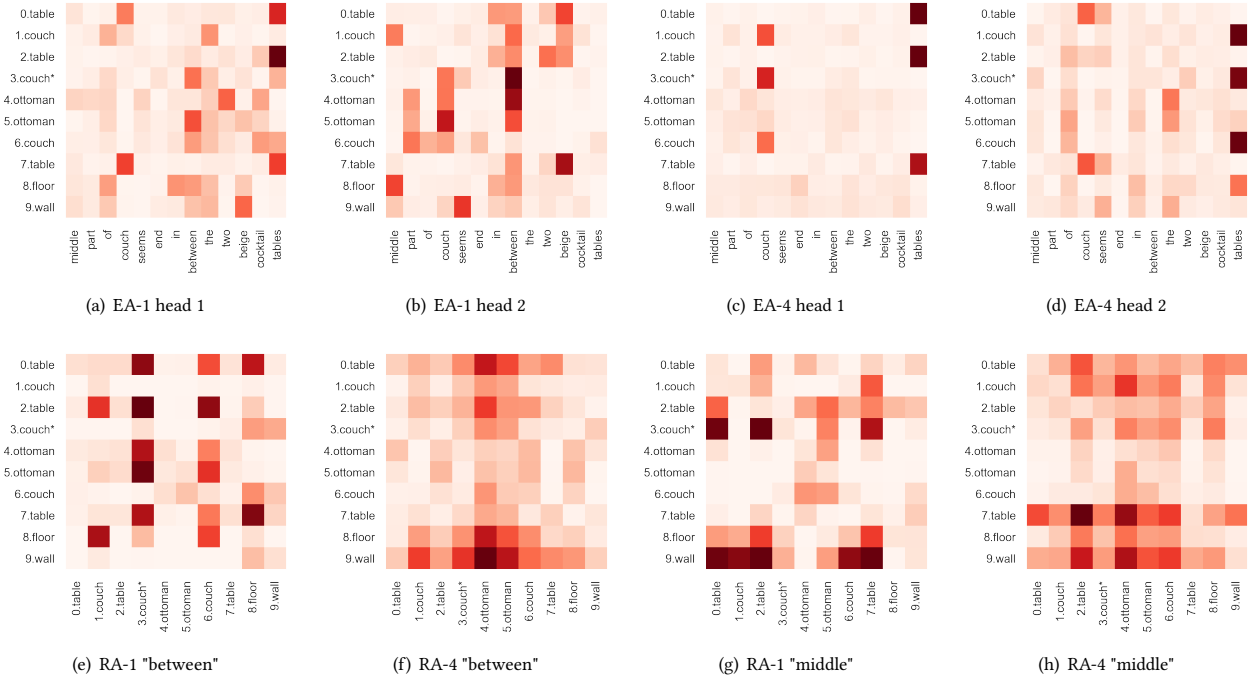
**4.4.1 Case Study.** We render the 3D scenes sampled from Nr3D test set and visualize the visual grounding results of our model. Figure 7 shows part of the results. Compared with ReferIt3DNet, our TransRefer3D performs better on both simple and hard cases because of a finely multi-modal context modeling. The first column in the figure indicates that, even in a complicated scene with a great number of distractors (5 monitors are replaced closely on a table), our model can still identify the referent from the context.

**4.4.2 Attention visualization of EA and RA.** Figure 8 shows the visualization of our proposed attention modules enhancing visual features. In this case, the referring expression is "Middle part of couch, seams end in between the 2 beige cocktail tables". For both EA and RA, we visualize the attention maps extracted from the first and the last ERCB (represented as EA/RA-1/4 in the figure).

According to the EA-4 maps showed in Figure 8(c) and Figure 8(d), it is very clear that in the last layer, EA matches the entity words with corresponding visual features (the couches and the tables in this case). We also observe that another EA head in the last



**Middle part of couch, seams end in between the 2 beige cocktail tables**



**Figure 8: Visualization of EA and RA. The attention maps are titled EA/RA- $x$  with  $x$  denoting the layer depth. The third word "couch" marked with an asterisk denotes the referent. (a-d) Attention maps of EA. For each EA module, we report attention maps from 2 of the 4 heads. (e-h) Attention maps of RA responding to the words "between" and "middle".**

layer matches the visual and linguistic features for couches and tables across. This behavior of EA perhaps helps establish a more expressive entity aware feature for the later grounding. See EA-1 maps in Figure 8(a) and Figure 8(b), EA behaviour in the first layer also maps tables and couches to different modality but not such obviously, compared with the EA-4 maps.

In this case, we investigate RA responding to two relational words, "between" and "middle". The attention maps of RA show a similar phenomenon that the maps from the last layer are more understandable than ones from the first layer. See Figure 8(f), the *0.table-4.ottoman* and *0.table-5.ottoman* relation representations strongly respond to the word feature "between". This makes sense since the table is exactly located between the two ottomans in the scene. Notice that "between" is a ternary relation describing a spatial state of an entity relative to the other two entities. In Figure 8(h), the referred part of couch (id. 3) is located in the "middle" of most of strongly responded object pairs (e.g. *7.table-2.table*, *9.wall-2.table*, *9.wall-4.ottoman*). The model learns to find objects surrounding the referent. This instances imply that the proposed RA module can also capture relations more complicated than binary relations.

## 5 CONCLUSION

We propose TransRefer3D consisting of ERCBs which integrate EA and RA modules for entity-and-relation aware multimodal context modeling. The proposed model performs effectively on fine-grained 3D visual grounding. On both Nr3D and Sr3D, the model significantly achieves superior performance compared with previous approaches. This proves that Transformer is suitable to connect language and 3D vision and demonstrates the effectiveness of our entity-and-relation aware modeling. Hence, the proposed EA and RA modules and ERCB may have a positive influence to future works on 3D multimodal machine learning. In the future, we will further dig into Transformers for 3D modal tasks. As 3D vision is rapidly developing today, we believe that the proposed TransRefer3D will become an important reference.

## 6 ACKNOWLEDGEMENT

This research is supported in part by National Natural Science Foundation of China (Grant 61876177), Beijing Natural Science Foundation (4202034), Fundamental Research Funds for the Central Universities.



## REFERENCES

- [1] 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 *ECCV*.
- [2] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. 2018. Bottom-up and top-down attention for image captioning and visual question answering. In *CVPR*.
- [3] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. 2016. Layer normalization. *arXiv preprint arXiv:1607.06450* (2016).
- [4] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. 2020. End-to-end object detection with transformers. In *ECCV*.
- [5] Dave Zhenyu Chen, Angel X Chang, and Matthias Nießner. 2019. Scanrefer: 3d object localization in rgb-d scans using natural language. *arXiv preprint arXiv:1912.08830* (2019).
- [6] Hanting Chen, Yunhe Wang, Tianyu Guo, Chang Xu, Yiping Deng, Zhenhua Liu, Siwei Ma, Chunjing Xu, Chao Xu, and Wen Gao. 2020. Pre-trained image processing transformer. *arXiv preprint arXiv:2012.00364* (2020).
- [7] Kan Chen, Rama Kovvuri, and Ram Nevatia. 2017. Query-guided regression network with context policy for phrase grounding. In *ICCV*.
- [8] Meng-Jiun Chiou, Roger Zimmermann, and Jiashi Feng. 2021. Visual Relationship Detection With Visual-Linguistic Knowledge From Multimodal Representations. *IEEE Access* 9 (2021), 50441–50451.
- [9] Marcella Cornia, Matteo Stefanini, Lorenzo Baraldi, and Rita Cucchiara. 2020. Meshed-memory transformer for image captioning. In *CVPR*.
- [10] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. 2017. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *CVPR*.
- [11] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [12] Zhipeng Ding, Xu Han, and Marc Niethammer. 2019. VoteNet: A deep learning label fusion method for multi-atlas segmentation. In *MICCAI*.
- [13] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xi-aohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020).
- [14] Lun Huang, Wenmin Wang, Jie Chen, and Xiao-Yong Wei. 2019. Attention on attention for image captioning. In *ICCV*.
- [15] 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 *AAAI*.
- [16] Shaofei Huang, Tianrui Hui, Si Liu, Guanbin Li, Yunchao Wei, Jizhong Han, Luoqi Liu, and Bo Li. 2020. Referring image segmentation via cross-modal progressive comprehension. In *CVPR*.
- [17] Chenchen Jing, Yuwei Wu, Mingtao Pei, Yao Hu, Yunde Jia, and Qi Wu. 2020. Visual-Semantic Graph Matching for Visual Grounding. In *ACM MM*.
- [18] Andrej Karpathy, Armand Joulin, and Li Fei-Fei. 2014. Deep fragment embeddings for bidirectional image sentence mapping. *arXiv preprint arXiv:1406.5679* (2014).
- [19] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [20] Kunpeng Li, Yulun Zhang, Kai Li, Yuanyuan Li, and Yun Fu. 2019. Visual semantic reasoning for image-text matching. In *ICCV*.
- [21] Chunxiao Liu, Zhendong Mao, Tianzhu Zhang, Hongtao Xie, Bin Wang, and Yongdong Zhang. 2020. Graph structured network for image-text matching. In *CVPR*.
- [22] Liyuan Liu, Xiaodong Liu, Jianfeng Gao, Weizhu Chen, and Jiawei Han. 2020. Understanding the difficulty of training transformers. *arXiv preprint arXiv:2004.08249* (2020).
- [23] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).
- [24] Jiasen Lu, Jianwei Yang, Dhruv Batra, and Devi Parikh. 2016. Hierarchical question-image co-attention for visual question answering. *arXiv preprint arXiv:1606.00061* (2016).
- [25] Anh Viet Phan, Minh Le Nguyen, Yen Lam Hoang Nguyen, and Lam Thu Bui. 2018. Dgcn: A convolutional neural network over large-scale labeled graphs. *Neural Networks* (2018).
- [26] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. 2017. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *CVPR*.
- [27] Charles R Qi, Li Yi, Hao Su, and Leonidas J Guibas. 2017. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *arXiv preprint arXiv:1706.02413* (2017).
- [28] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *arXiv preprint arXiv:1706.03762* (2017).
- [29] Haoran Wang, Ying Zhang, Zhong Ji, Yanwei Pang, and Lin Ma. 2020. Consensus-aware visual-semantic embedding for image-text matching. In *ECCV*.
- [30] Zhengyuan Yang, Boqing Gong, Liwei Wang, Wenbing Huang, Dong Yu, and Jiebo Luo. 2019. A fast and accurate one-stage approach to visual grounding. In *ICCV*.
- [31] Tianyu Yu, Tianrui Hui, Zhihao Yu, Yue Liao, Sansi Yu, Faxi Zhang, and Si Liu. 2020. Cross-Modal Omni Interaction Modeling for Phrase Grounding. In *ACM MM*.
- [32] Zhou Yu, Jun Yu, Yuhao Cui, Dacheng Tao, and Qi Tian. 2019. Deep modular co-attention networks for visual question answering. In *CVPR*.
- [33] Zhou Yu, Jun Yu, Chenchao Xiang, Jianping Fan, and Dacheng Tao. 2018. Beyond bilinear: Generalized multimodal factorized high-order pooling for visual question answering. *TNNLS* (2018).
- [34] Zhihao Yuan, Xu Yan, Yinghong Liao, Ruimao Zhang, Zhen Li, and Shuguang Cui. 2021. InstanceRefer: Cooperative Holistic Understanding for Visual Grounding on Point Clouds through Instance Multi-level Contextual Referring. *arXiv preprint arXiv:2103.01128* (2021).
- [35] Seongjun Yun, Minbyul Jeong, Raehyun Kim, Jaewoo Kang, and Hyunwoo J Kim. 2019. Graph transformer networks. *arXiv preprint arXiv:1911.06455* (2019).
- [36] Rowan Zellers, Mark Yatskar, Sam Thomson, and Yejin Choi. 2018. Neural motifs: Scene graph parsing with global context. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 5831–5840.
- [37] Chao Zhang, Weiming Li, Wanli Ouyang, Qiang Wang, Woo-Shik Kim, and Sunghoon Hong. 2019. Referring Expression Comprehension with Semantic Visual Relationship and Word Mapping. In *ACM MM*.
- [38] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip Torr, and Vladlen Koltun. 2020. Point transformer. *arXiv preprint arXiv:2012.09164* (2020).