

Grounding Commands for Autonomous Vehicles via Layer Fusion with Region-specific Dynamic Layer Attention

Hou Pong Chan^{*1}, Mingxi Guo^{*1}, and Cheng-Zhong Xu^{†1}

Abstract—Grounding a command to the visual environment is an essential ingredient for interactions between autonomous vehicles and humans. In this work, we study the problem of language grounding for autonomous vehicles, which aims to localize a region in a visual scene according to a natural language command from a passenger. Prior work only employs the top layer representations of a vision-and-language pre-trained model to predict the region referred to by the command. However, such a method omits the useful features encoded in other layers, and thus results in inadequate understanding of the input scene and command. To tackle this limitation, we present the first layer fusion approach for this task. Since different visual regions may require distinct types of features to disambiguate them from each other, we further propose the region-specific dynamic (RSD) layer attention to adaptively fuse the multimodal information across layers for each region. Extensive experiments on the Talk2Car benchmark demonstrate that our approach helps predict more accurate regions and outperforms state-of-the-art methods.

I. INTRODUCTION

In recent years, autonomous driving systems achieve significant progress due to the advances in sensing technologies and deep neural networks. However, self-driving requires a very high level of trust from users and it will only be adopted if it creates a better human experience [1]. Prior literature [2] reveals that enabling users to issue commands to autonomous vehicles via natural language can improve user experience and acceptance. In order to perform the action specified by a command (e.g., follow a particular car or stop at a specific parking spot), a vehicle needs to understand the semantic correspondence between the natural language command and the visual environment. This problem is formalized as the task of *language grounding for autonomous vehicles* [3]: given an image and a natural language command to a vehicle, the goal is to locate the region in the image that is referred to by the command. We give a sample of image, command, and the referred region in Figure 1.

Existing work typically uses an object detector to extract region proposals from the input image and casts the task as selecting the best-matched region based on the command. State-of-the-art (SOTA) methods [5] apply a vision-and-language (V&L) pre-trained model to summarize the region proposals and command words into contextualized representations through multiple Transformer [6] encoder layers, and then use the top layer representations to predict



Command: find a parking spot near the *first concrete barrier*.

Fig. 1: A sample image and the associated natural language command. The blue box indicates the ground-truth region referred to by the command. The red box shows the prediction by the UNITER model [4] when using the top layer representations to compute matching scores of regions. The green box is predicted by the UNITER model after enhanced by our RSD layer attention approach.

the matching scores of regions. Several studies [7], [8] show that the representations learned by different encoder layers embed different types of surface and semantic features. However, SOTA language grounding methods only utilize the top layer representations to compute the matching scores. In consequence, these methods ignore the useful features inside the representations in other layers, which may lead to insufficient comprehension of the input scene and command. For instance, Figure 1 shows a V&L pre-trained model, UNITER [4], outputs an incorrect region when using the top layer representations to predict the matching scores of regions.

To effectively utilize the features embedded in different encoder layers, we propose *the first encoder layer fusion approach for the language grounding task*. For each region proposal, our approach fuses its representations across all encoder layers in a V&L pre-trained model. After that, we feed the fused representation to an output layer to predict its matching score. Our work investigates the layer attention technique [7], [9] for encoder layer fusion since it has good interpretability and we can directly assess the contribution made by each encoder layer.

Current layer attention methods assign a static set of normalized attention weights to encoder layers regardless of

^{*}Both authors contributed equally to this research.

[†]Corresponding author.

¹Department of Computer and Information Science, University of Macau, Macau SAR, China. Email: {hpchan, mc05415, czxu}@um.edu.mo

the input. Then they use the attention weights to aggregate the encoder representations over layers by weighted sum. However, in the language grounding problem, different visual regions may require distinct types of linguistic and visual features to determine whether they match the command. As an example, for the trees in Figure 1, a model only needs the object category information to reject them as the correction region. Whereas for the concrete barriers in the figure, a model requires both the object category and position features to determine whether they match the command.

To address the above drawback of existing layer attention methods, we further propose the **region-specific dynamic (RSD) layer attention** mechanism, which dynamically computes a new set of layer attention weights for each individual image region. The attention weights are learned by a neural module based on the representations of the regions. Thus, our method can adaptively determine the importance of different encoder layers and acquire appropriate features for each visual region. Figure 1 illustrates that our RSD layer attention helps predict a more accurate region¹.

Comprehensive empirical studies are conducted on the Talk2Car [3] benchmark. We apply our layer fusion approach to enhance two recent V&L pre-trained models, UNITER [4] and LXMERT [10]. Experiment results demonstrate that our approach consistently improves the accuracy of both the UNITER and LXMERT models and outperforms the SOTA methods in this task. Moreover, our proposed RSD layer attention also achieves better performance than existing layer fusion methods. We then examine how our method distributes the layer attention weights. Furthermore, we give a qualitative study to illustrate why our method leads to more accurate results. Finally, we evaluate our approach in a closely related task, referring expression comprehension [11], to assess the generality of our approach.

We summarize the contributions of this paper as follows: (1) the first encoder layer fusion approach for the language grounding task; (2) a novel RSD layer attention mechanism that dynamically aggregates the information in encoder layers for each input visual region; (3) an extensive empirical analysis of our encoder layer fusion approach; and (4) new state-of-the-art results in the Talk2Car benchmark.

II. RELATED WORK

A. Language Grounding for Autonomous Vehicles

Language grounding for autonomous vehicles [3] is an important task for human-vehicle interactions. To tackle this task, Vandenhende et al. [12] propose the C4AV-Base model that utilizes the bi-directional GRU [13] and ResNet models [14] to encode the command and the regions respectively and then compute their feature correlation. Later, the ASSMR method [15] applies an attention mechanism to strengthen the features extracted by bi-directional GRU and ResNet. To enhance the reasoning ability, the MSRR method [16] introduces a spatial memory module and a multi-step reasoning module to iteratively score the regions.

Other methods are built on the Transformer [6]. The CMTR method [17] applies the Transformer encoder-decoder model to model the command and regions separately. Dai et al. [5] use the VL-BERT pre-trained model [18] to jointly learn cross-modal representations for the input. They propose an iterative stacking algorithm, Stack-VL-BERT, to train a deeper VL-BERT model. This method only uses the top layer representations of a pre-trained model to compute matching scores, while our approach fuses all encoder layers in a pre-trained model.

B. Referring Expression Comprehension

The referring expression comprehension task aims to predict a region in an image according to a direct description of an object [11], e.g., “*girl on the left*”. In contrast, the command expressions in the language grounding for autonomous vehicles task are more complex and involve both action and object descriptions. Earlier methods [19]–[23] are built on CNN [24] and LSTM [25] models to encode the input regions and referring expression. Recent methods [4], [10], [18] use V&L pre-trained models based on the Transformer architecture and they achieve state-of-the-art performance. These methods only use the representations in the last layer of a pre-trained model to predict the matching scores, while our method fuses the representations across all encoder layers.

C. Language Grounding in Human-robot Interaction

Several methods [26]–[28] allow robots to ask clarification questions and they rely on a semantic parser [27], [28] or a probabilistic model [26] to decompose a command and then ground the resulting constituents to visual regions. Shridhar and Hsu [29] propose a robot system that picks up an object according to a command. Their method generates a caption for each region and then clusters the generated captions with the input command. Kim et al. [30] introduce a vehicle controller that accepts advice from humans. Their model uses LSTM, CNN, and visual attention [31] to ground an advice to an image. In contrast, our approach is built on V&L pre-trained models to ground a command to a visual scene.

D. Encoder Layers Fusion

Encoder layer fusion techniques can be categorized into layer aggregation [32]–[34], layer-wise coordination [35], and layer attention [7], [9]. Layer-wise coordination modifies the structure of Transformer and cannot be easily applied to pre-trained models. Some layer aggregation methods [33] dynamically fuse the representations for each input word, but they do not provide normalized weights to interpret the importance of each layer and they introduce millions of new parameters. Layer attention utilizes normalized attention weights to fuse encoder layers and allows us to directly interpret the contribution of each layer. Existing layer attention methods assign static attention weights to layers independent of the input, whereas our RSD layer attention dynamically predicts attention weights for every input region.

¹We also provide a demo of our approach in the attached video.

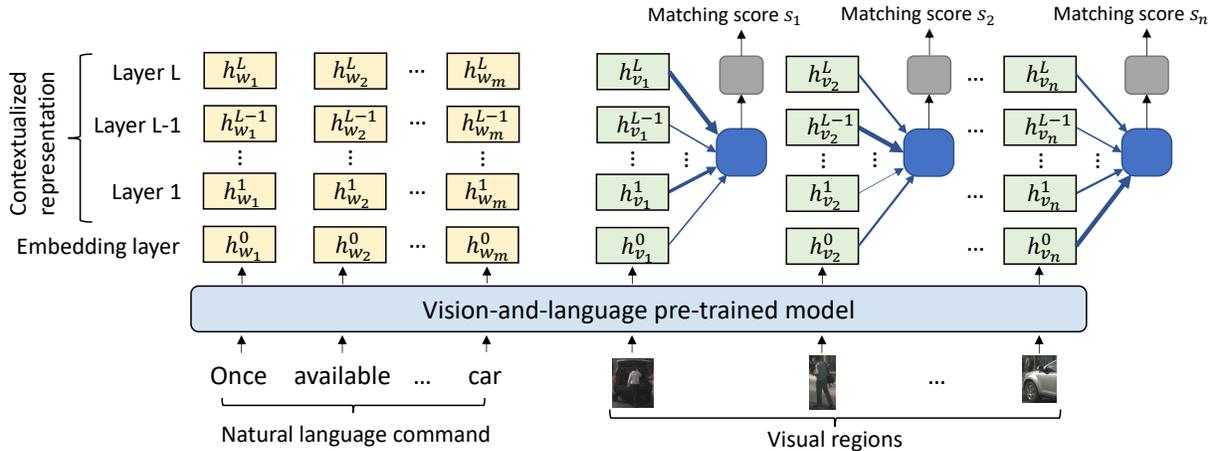


Fig. 2: Overview of our approach. A V&L pre-trained model first encodes the command and image regions into multiple layers of representation vectors. Then for each visual region, our RSD layer attention mechanism (indicated by blue rounded square) fuses the representations across encoder layers. Finally, we feed every fused representation to a linear output layer (grey rounded square) to predict the matching scores of regions.

III. PROBLEM DEFINITION

Given a natural language command \mathbf{c} and an image \mathbf{I} , the goal is to predict the region in the image that the command is referring to. Following previous literature [12], [18], we formulate the task as selecting the best-matched region v^* from a set of region proposals $\{v_i\}_{i=1}^n$ in the image \mathbf{I} according to the command \mathbf{c} .

IV. OUR LAYER FUSION APPROACH

In our approach, we first pass the input command and image to a V&L pre-trained model to obtain multiple layers of cross-modal representations. We then propose the region-specific dynamic (RSD) layer attention mechanism to dynamically fuse representations across all layers for each individual visual region. The fused representations are then fed to a linear output layer to predict the matching scores of the regions. We display the overall architecture of our approach in Figure 2.

A. Vision-and-language (V&L) Pre-trained Model.

As a preprocessing step, we obtain the region proposals $\{v_i\}_{i=1}^n$ of the input image from the Centernet model [36]. We use the WordPieces tokenizer of BERT [37] to tokenize the input command into a sequence of tokens w_1, \dots, w_m . A V&L pre-trained model takes the region proposals and command tokens as input. Inside the pre-trained model, an embedding layer first converts the input tokens and regions into a sequence of word embeddings, $\mathbf{h}_{w_1}^0, \dots, \mathbf{h}_{w_m}^0$, and region embeddings, $\mathbf{h}_{v_1}^0, \dots, \mathbf{h}_{v_n}^0$, respectively. We refer to the embedding layer as the 0-th encoder layer of the pre-trained model.

Then, the model feeds the embeddings to a Transformer encoder with L layers. At each layer from 1 to L , the encoder uses the multi-head attention mechanism [6] to model intra-modal and/or cross-modal interactions among the input to produce a d -dimensional contextualized representation vector for each region and token. We use $\mathbf{h}_{w_i}^l$ to denote the

contextualized representation of region v_i produced by the l -th encoder layer. Several studies [7], [8] reveal that lower encoder layers extract more concrete features (e.g., position) while higher encoder layers extract more abstract features (e.g., coreference relation). The model is pre-trained on a massive dataset to learn visual and linguistic knowledge.

There are two major classes of V&L pre-trained models: (1) single-stream architecture, which models both intra-modal and cross-modal interactions in every encoder layer; and (2) dual-stream architecture, which only allows intra-modal interactions in early encoder layers and then encodes both intra-modal and cross-modal interactions in latter layers. We apply the **UNITER** [4] model from single-stream architecture and **LXMERT** [10] from dual-stream architecture because they achieve excellent performance in many V&L tasks [38]. On the other hand, SOTA methods [5] in this task use another single-stream model, VL-BERT [18]. Their model obtains lower scores than UNITER and LXMERT (see Table I for the results).

B. Region-specific Dynamic (RSD) Layer Attention.

Prior methods feed the top layer representation of a region $\mathbf{h}_{v_i}^L$ to a linear layer to predict its matching score. To effectively exploit the concrete and abstract features embedded in different layers, we propose a novel RSD layer attention mechanism to fuse the encoder representations across layers before the prediction of matching scores. The overall idea is to dynamically assign attention weights to all encoder layers based on a region's representation vectors in these layers. The weights are then used to aggregate the representation vectors across layers. Our intuition is that every region needs distinct types of features to decide whether it is referred to by the command.

More concretely, for every region proposal v_i , our RSD layer attention method passes its representation $\mathbf{h}_{v_i}^l$ at each encoder layer l into a linear layer to compute a relevance score α_i^l . Then we use the softmax function to normalize the

relevance scores over all encoder layers and obtain the layer attention weights for region v_i , as shown in the following equations:

$$\alpha_i^l = \mathbf{W}_\alpha \mathbf{h}_{v_i}^l + b_\alpha, \quad a_i^l = \frac{\exp(\alpha_i^l)}{\sum_{l'=0}^L \exp(\alpha_i^{l'})}, \quad (1)$$

where a_i^l denotes the layer attention weight at layer l by region v_i and indicates the importance of the l -th layer to region v_i . In contrast, existing layer attention methods [7], [9], [39] only allocate fixed attention weights to encoder layers regardless of the model input. We then fuse the representations at different layers by weighted sum: $\tilde{\mathbf{h}}_{v_i} = \sum_{l=0}^L a_i^l \mathbf{h}_{v_i}^l$. Our RSD layer attention is parameter-efficient since it only introduces $d+1$ new parameters, where $d = 768$ for the UNITER and LXMERT models.

C. Output Layer and Loss Function.

Finally, we feed the fused encoder representation $\tilde{\mathbf{h}}_{v_i}$ of each region v_i to a linear output layer to predict a matching score: $s_i = \sigma(\mathbf{W}_s \tilde{\mathbf{h}}_{v_i} + b_s)$, where $\mathbf{W}_s \in \mathbb{R}^{d \times 1}$ and $b_s \in \mathbb{R}$, σ is the sigmoid function that normalizes the output to the range of $[0, 1]$. Following previous work [40], we use the **intersection-over-union (IoU)** score between a region v_i and the ground-truth region as the ground-truth matching score s_i^* . The IoU score of a predicted region is the intersection between that region and the ground-truth region divided by the union of them. The training objective is the **binary cross-entropy loss**: $s_i^* \log \sigma(s_i) + (1 - s_i^*) \log(1 - \sigma(s_i))$. During inference, we use the region proposal that has the highest predicted matching score as the model output.

D. Model Ablation.

To verify the importance of region-specific attention weights in our method, we make an ablation in our RSD layer attention to construct a baseline called **sample-specific layer attention**. Instead of predicting layer attention weights for each individual region, we compute layer attention weights for the entire input sample using the mean-pooled representation of the regions. Then, each region in the sample uses the same set of attention weights to fuse encoder layers.

Specifically, at each encoder layer l , we perform mean-pooling over the representations of all regions: $\mathbf{h}^l = 1/n \sum_{i=1}^n \mathbf{h}_{v_i}^l$. Next, we feed the mean-pooled representation \mathbf{h}^l to a linear layer to learn a relevance score. All the relevance scores are then normalized by softmax to yield the layer attention weights shared by all the regions.

V. EXPERIMENTAL SETUP

A. Datasets

We use the **Talk2Car** [3] dataset to conduct our experiments. Talk2Car is the standard benchmark for the language grounding for autonomous vehicles task. The images are urban scenes captured by cameras on a car. Each input text is a natural language command referring to a particular object in the input image, e.g., “*get a parking spot near the second car on the left side*”. A command expression contains both an action instruction and an object description. The averaged

TABLE I: IoU_{0.5} scores of different models on the Talk2Car dataset. We bold the best results and underline the second-best results.

Model	Val. set	Test set
C4AV-Base	43.5	44.1
MAC	-	50.5
MSRR	60.3	60.1
ASSMR	67.9	66.4
CMTR	68.2	69.1
Stack-VL-BERT	68.2	71.0
LXMERT	72.7	73.1
RSD-LXMERT (Ours)	<u>74.7</u>	<u>73.7</u>
UNITER	74.0	73.2
RSD-UNITER (Ours)	74.9	73.9

length of expressions is 11.0. The training, validation, and test splits contain 8,349/1,163/2,447 samples.

We further evaluate our method on the **RefCOCO+** [11] and **RefCOCog** [19] datasets, which are popular benchmarks for referring expression comprehension. The input text in these two datasets is a description of an object in the image (e.g., “*giraffe with lowered head*”) rather than a command. The average expression length of RefCOCO+ and RefCOCog datasets are 3.5 and 8.4 respectively. The images in these datasets are collected from the MSCOCO benchmark [41]. For RefCOCO+, we divide the data according to Bugliarello et al. [40] and obtain the split of 287,113/13,368/11,490 for training, validation, and test. For RefCOCog, we use the UMD split of 42,226/2,573/5,023.

B. Evaluation Metric

During evaluation, we compute the IoU score of the predicted region. If the IoU score is larger than 0.5, we consider the prediction correct. Following [3], we report the averaged number of correct predictions and refer to it as the IoU_{0.5} score.

C. Baselines and Comparison

We adopt recent methods of language grounding for autonomous vehicles as baselines, including C4AV-Base [12], ASSMR [15], MSRR [16], CMTR [17], and Stack-VL-BERT [5]. Moreover, we consider the UNITER [4] and LXMERT [10] models as baselines. For the UNITER model, we adopt the base model with 12 layers. Furthermore, we compare with previous layer fusion methods: **coarse-grained layer attention** [7], [9], which assigns static attention weights to encoder layers, **fine-grained layer attention** [39], which assigns static attention weights to the elements of the representations in all layers, **dynamic combination** [33], which uses L feedforward networks to aggregate the representations in different layers, **dynamic routing** [33], which iteratively refines the fused representation based on the agreement between each layer representation and the fused representation. We use **RSD-UNITER** and **RSD-LXMERT** to denote the UNITER and LXMERT models after being enhanced by our layer fusion approach.

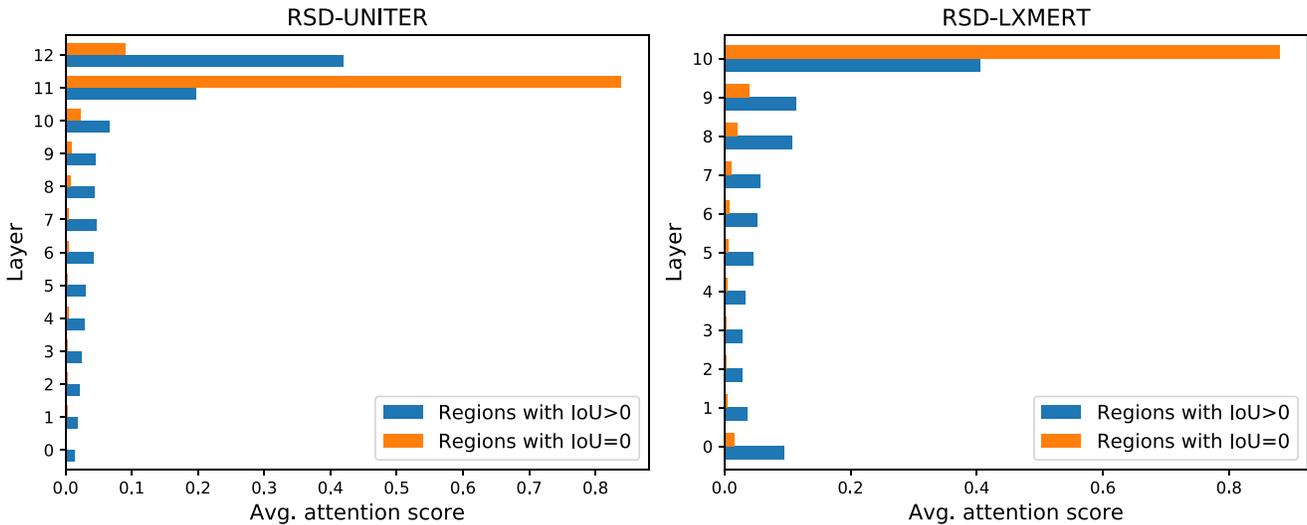


Fig. 3: Attention weights over all encoder layers learned by our RSD layer attention. The attention weights are averaged over all regions in the validation set of Talk2Car. The 0-th layer indicates the embedding layer.

VI. EXPERIMENT RESULTS

A. Comparison Results with SOTA methods

We report the main comparison results on the Talk2Car dataset in Table I. It is observed that our layer fusion approach significantly and consistently improves the $IOU_{0.5}$ scores of both LXMERT and UNITER models in validation and test sets, which show that our RSD layer attention can increase the accuracy of both single-stream and dual-stream V&L pre-trained models. The reason is that our approach is agnostic to the architecture of the underlying pre-trained model. Moreover, both of our methods (RSD-LXMERT and RSD-UNITER) outperform the state-of-the-art models in this benchmark. The above results indicate that it is important to effectively utilize the features embedded in different layers of a V&L pre-trained model to locate the referred region.

B. Comparison of Different Layer Fusion Methods

We compare the performance of different encoder layer fusion methods on UNITER since it is the most accurate pre-trained model as shown in the previous section. From Table II, we observe that our RSD layer attention mechanism achieves better performance than existing layer fusion methods. Moreover, we can see that our constructed baseline, sample-specific layer attention, obtains lower $IOU_{0.5}$ scores than our RSD layer attention, which indicate that it is crucial to learn a distinct set of layer attention weights for each visual region. Surprisingly, our RSD layer attention substantially outperforms previous dynamic layer fusion methods (rows 4&5) that introduce millions of new parameters. The results suggest that the combination of region specific layer attention weights and weighted sum fusion operation is effective in this task. Previous dynamic aggregation methods may introduce unnecessary parameters which make the model more difficult to train.

TABLE II: $IOU_{0.5}$ scores of different layer fusion methods in the UNITER model on the Talk2Car dataset. # Param. denotes the number of parameters.

#	Model	Val	Test	# Param.
1	UNITER	74.0	73.2	112.0M
2	FineGrainedAttn-UNITER	73.5	73.0	+9984
3	CoarseGrainedAttn-UNITER	74.0	73.6	+13
4	DynamicCombin-UNITER	74.1	72.3	+107.3M
5	DynamicRouting-UNITER	72.4	72.2	+7.7M
6	SampleSpecificAttn-UNITER	74.6	72.6	+769
7	RSD-UNITER	74.9	73.9	+769

C. Analysis of Layer Attention Weights

We analyze the distribution of layer attention weights learned by our RSD layer attention for different regions. We partition the input regions into the following two groups according to their IoU with the ground-truth region. $IoU > 0$: contains all the regions that have IoU greater than 0 (i.e., the regions that overlap the ground-truth region). $IoU = 0$: includes all the regions with IoU equals to 0. Figure 3 shows the distribution of layer attention weights in two groups of regions. We draw the following two observations.

- Both models allocate a larger proportion of attention weights to higher encoder layers. It is because the representations in higher encoder layers are exposed to more cross-modal interactions. Thus, they usually capture more cross-modal semantic features which are essential to command and scene understanding.
- Compared with the regions with $IoU = 0$, the regions with $IoU > 0$ assign more attention weights to lower encoder layers. It is because most of the regions with $IoU = 0$ are obviously irrelevant to the command. Hence, the abstract features embedded in higher encoder layers can provide sufficient information for the model to predict a matching score of zero. On the other hand, the regions with $IoU > 0$ overlap the target region and they are more difficult to be disambiguated. Hence,

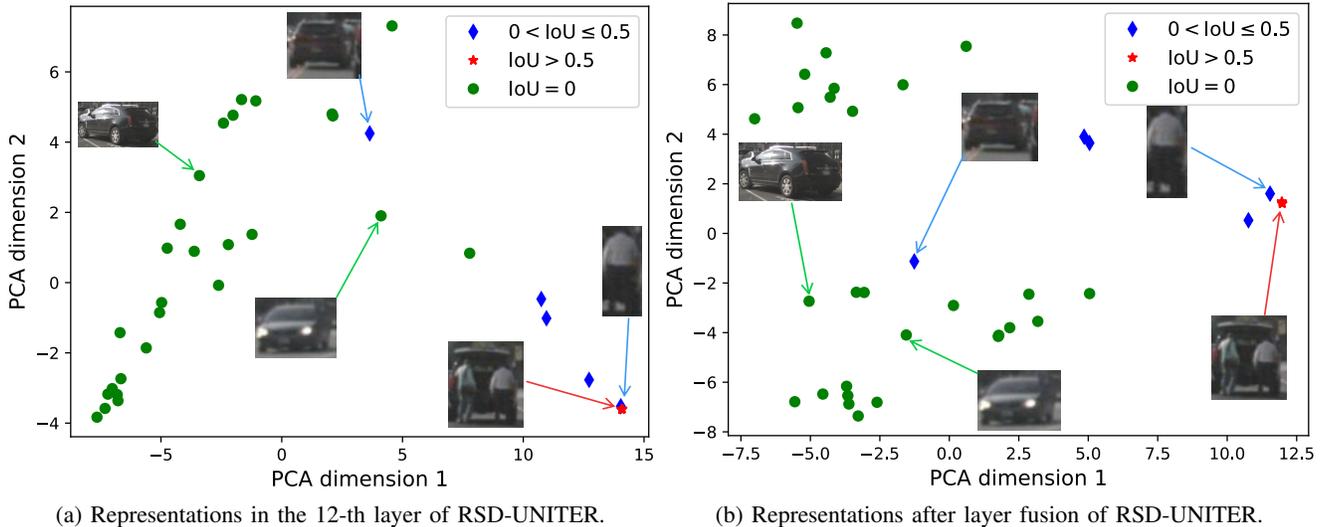


Fig. 4: Regions’ representations of a validation sample in the RSD-UNITER model after projected by PCA. The input command is “when it is safe, slowly pass the car that is standing still up ahead”. Red star denotes the ground-truth region.

the model requires more surface-level features (e.g., position) of the regions in lower encoder layers to decide their matching scores.

D. Qualitative Analysis of Region Representations

To illustrate how our layer fusion method results in more accurate predictions, we qualitatively analyze the representations of visual regions. We first collect the regions’ representations learned by our RSD-UNITER model. We then utilize the Principal Component Analysis (PCA) [42] technique to project the representation vectors into a 2-dimensional vector space for the sake of visualization.

Figure 4 visualizes the regions’ representations learned by our RSD-UNITER model. The ground-truth region’s representation is denoted by a red star. From Figure 4a, we observe that in the top layer, the ground-truth region’s representation is far away from the representations of regions that do not overlap the ground-truth ($\text{IoU} = 0$). Previous methods feed the top layer representation of a region to a linear layer to predict its matching score. However, the ground-truth region is extremely close to a region that has intersection with the ground-truth ($0 < \text{IoU} \leq 0.5$). It is difficult for a linear layer to separate the target region from its closest neighbor. Then in Figure 4b, we observe that after applying our RSD layer attention to fuse the representations across layers, the ground-truth region’s representation is pushed slightly further from its closest neighbor. Thus, a linear layer can separate out the ground-truth region more easily, which demonstrates the advantage of our approach.

E. Results in Referring Expression Comprehension

We further evaluate the performance of our approach on the referring expression comprehension task. Table III presents the results in the RefCOCO+ and RefCOCOg datasets. We observe that our approach increases the performance of UNITER and LXMERT models in most of the

TABLE III: $\text{IoU}_{0.5}$ results in the RefCOCO+ and RefCOCOg datasets.

Model	RefCOCO+		RefCOCOg	
	Val	Test	Val	Test
UNITER	72.0	70.8	73.1	73.2
RSD-UNITER	72.6	71.5	74.0	73.0
LXMERT	71.3	70.0	71.9	72.1
RSD-LXMERT	71.4	70.1	72.3	72.4

cases but the improvements are less significant than that in the Talk2Car dataset. We analyze the reason as follow. The command expressions in the autonomous driving setting are complicated and contain both action and object descriptions. The model requires features from multiple layers of cross-modal representations to learn the alignment between text and image. Thus, our encoder layer fusion approach significantly improves the accuracy of predicted regions. On the other hand, the referring expressions in RefCOCO+ and RefCOCOg datasets are straight-forward descriptions of a target object. Hence, the top layer representations often capture enough information to find the target object and encoder layer fusion provides a smaller benefit to the model.

VII. CONCLUSION

In this work, we present the first encoder layer fusion approach for the language grounding for autonomous vehicles task. In our approach, we apply a V&L pre-trained model to learn contextualized representations for the input command and regions. Then we propose a novel RSD layer attention mechanism to dynamically aggregate the representations across all encoder layers for each individual region proposal. Experiment results on the real-world Talk2Car benchmark show that our approach consistently improves the performance of UNITER and LXMERT pre-trained models and achieves the new SOTA results in this task.

ACKNOWLEDGMENT

This paper is supported by National Key Research and Development Program of China (No. 2019YFB2102100), the Science and Technology Development Fund of Macau SAR (File no. 0015/2019/AKP), and Guangdong-Hong Kong-Macao Joint Laboratory of Human-Machine Intelligence-Synergy Systems (No. 2019B121205007).

REFERENCES

- [1] L. Fridman, “Self-driving cars state of the art (2019),” *MIT Deep Learning and Artificial Intelligence Lectures*, p. 11, 2019.
- [2] J. Kim, S. Moon, A. Rohrbach, T. Darrell, and J. F. Canny, “Advisable learning for self-driving vehicles by internalizing observation-to-action rules,” in *CVPR*, 2020, pp. 9658–9667.
- [3] T. Deruyttere, S. Vandenhende, D. Grujicic, L. V. Gool, and M. Moens, “Talk2car: Taking control of your self-driving car,” in *EMNLP-IJCNLP*, 2019, pp. 2088–2098.
- [4] Y. Chen, L. Li, L. Yu, A. E. Kholy, F. Ahmed, Z. Gan, Y. Cheng, and J. Liu, “UNITER: universal image-text representation learning,” in *ECCV*, ser. Lecture Notes in Computer Science, vol. 12375, 2020, pp. 104–120.
- [5] H. Dai, S. Luo, Y. Ding, and L. Shao, “Commands for autonomous vehicles by progressively stacking visual-linguistic representations,” in *ECCV*, ser. Lecture Notes in Computer Science, vol. 12536, 2020, pp. 27–32.
- [6] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” in *NeurIPS*, 2017, pp. 5998–6008.
- [7] M. E. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, and L. Zettlemoyer, “Deep contextualized word representations,” in *NAACL-HLT*, 2018, pp. 2227–2237.
- [8] I. Tenney, D. Das, and E. Pavlick, “BERT rediscovers the classical NLP pipeline,” in *ACL*, 2019, pp. 4593–4601.
- [9] A. Bapna, M. X. Chen, O. Firat, Y. Cao, and Y. Wu, “Training deeper neural machine translation models with transparent attention,” in *EMNLP*, 2018, pp. 3028–3033.
- [10] H. Tan and M. Bansal, “LXMERT: learning cross-modality encoder representations from transformers,” in *EMNLP-IJCNLP*, 2019, pp. 5099–5110.
- [11] S. Kazemzadeh, V. Ordonez, M. Matten, and T. L. Berg, “Referitgame: Referring to objects in photographs of natural scenes,” in *EMNLP*, 2014, pp. 787–798.
- [12] S. Vandenhende, T. Deruyttere, and D. Grujicic, “A baseline for the commands for autonomous vehicles challenge,” *CoRR*, vol. abs/2004.13822, 2020.
- [13] K. Cho, B. van Merriënboer, Ç. Gülçehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, “Learning phrase representations using RNN encoder-decoder for statistical machine translation,” in *EMNLP*, 2014, pp. 1724–1734.
- [14] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *CVPR*, 2016, pp. 770–778.
- [15] J. Ou and X. Zhang, “Attention enhanced single stage multimodal reasoner,” in *ECCV*, ser. Lecture Notes in Computer Science, vol. 12536, 2020, pp. 51–61.
- [16] T. Deruyttere, G. Collell, and M. Moens, “Giving commands to a self-driving car: A multimodal reasoner for visual grounding,” *CoRR*, vol. abs/2003.08717, 2020.
- [17] S. Luo, H. Dai, L. Shao, and Y. Ding, “C4AV: learning cross-modal representations from transformers,” in *ECCV*, ser. Lecture Notes in Computer Science, vol. 12536, 2020, pp. 33–38.
- [18] W. Su, X. Zhu, Y. Cao, B. Li, L. Lu, F. Wei, and J. Dai, “VL-BERT: pre-training of generic visual-linguistic representations,” in *ICLR*, 2020.
- [19] J. Mao, J. Huang, A. Toshev, O. Camburu, A. L. Yuille, and K. Murphy, “Generation and comprehension of unambiguous object descriptions,” in *CVPR*, 2016, pp. 11–20.
- [20] R. Hu, H. Xu, M. Rohrbach, J. Feng, K. Saenko, and T. Darrell, “Natural language object retrieval,” in *CVPR*, 2016, pp. 4555–4564.
- [21] R. Luo and G. Shakhnarovich, “Comprehension-guided referring expressions,” in *CVPR*, 2017, pp. 3125–3134.
- [22] J. Liu, L. Wang, and M. Yang, “Referring expression generation and comprehension via attributes,” in *ICCV*, 2017, pp. 4866–4874.
- [23] L. Yu, Z. Lin, X. Shen, J. Yang, X. Lu, M. Bansal, and T. L. Berg, “Mattnet: Modular attention network for referring expression comprehension,” in *CVPR*, 2018, pp. 1307–1315.
- [24] Y. LeCun, B. E. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. E. Hubbard, and L. D. Jackel, “Backpropagation applied to handwritten zip code recognition,” *Neural Computation*, vol. 1, no. 4, pp. 541–551, 1989. [Online]. Available: <https://doi.org/10.1162/neco.1989.1.4.541>
- [25] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997. [Online]. Available: <https://doi.org/10.1162/neco.1997.9.8.1735>
- [26] R. Deits, S. Tellex, P. Thaker, D. Simeonov, T. Kollar, and N. Roy, “Clarifying commands with information-theoretic human-robot dialog,” *Journal of Human-Robot Interaction*, vol. 2, no. 2, pp. 58–79, 2013.
- [27] J. Thomason, A. Padmakumar, J. Sinapov, N. Walker, Y. Jiang, H. Yedidsion, J. Hart, P. Stone, and R. Mooney, “Jointly improving parsing and perception for natural language commands through human-robot dialog,” *Journal of Artificial Intelligence Research*, vol. 67, pp. 327–374, 2020.
- [28] J. Thomason, A. Padmakumar, J. Sinapov, N. Walker, Y. Jiang, H. Yedidsion, J. W. Hart, P. Stone, and R. J. Mooney, “Improving grounded natural language understanding through human-robot dialog,” in *ICRA*, 2019, pp. 6934–6941.
- [29] M. Shridhar and D. Hsu, “Interactive visual grounding of referring expressions for human-robot interaction,” in *RSS*, 2018.
- [30] J. Kim, T. Misu, Y. Chen, A. Tawari, and J. F. Canny, “Grounding human-to-vehicle advice for self-driving vehicles,” in *CVPR*, 2019, pp. 10591–10599.
- [31] K. Xu, J. Ba, R. Kiros, K. Cho, A. C. Courville, R. Salakhutdinov, R. S. Zemel, and Y. Bengio, “Show, attend and tell: Neural image caption generation with visual attention,” in *ICML*, ser. JMLR Workshop and Conference Proceedings, vol. 37, 2015, pp. 2048–2057.
- [32] F. Yu, D. Wang, E. Shelhamer, and T. Darrell, “Deep layer aggregation,” in *CVPR*, 2018, pp. 2403–2412.
- [33] Z. Dou, Z. Tu, X. Wang, L. Wang, S. Shi, and T. Zhang, “Dynamic layer aggregation for neural machine translation with routing-by-agreement,” in *AAAI*, 2019, pp. 86–93.
- [34] Q. Wang, B. Li, T. Xiao, J. Zhu, C. Li, D. F. Wong, and L. S. Chao, “Learning deep transformer models for machine translation,” in *ACL*, 2019, pp. 1810–1822.
- [35] T. He, X. Tan, Y. Xia, D. He, T. Qin, Z. Chen, and T. Liu, “Layer-wise coordination between encoder and decoder for neural machine translation,” in *NeurIPS*, 2018, pp. 7955–7965.
- [36] X. Zhou, D. Wang, and P. Krähenbühl, “Objects as points,” *CoRR*, vol. abs/1904.07850, 2019.
- [37] J. Devlin, M. Chang, K. Lee, and K. Toutanova, “BERT: pre-training of deep bidirectional transformers for language understanding,” in *NAACL-HLT*, 2019, pp. 4171–4186.
- [38] J. Cao, Z. Gan, Y. Cheng, L. Yu, Y. Chen, and J. Liu, “Behind the scene: Revealing the secrets of pre-trained vision-and-language models,” in *ECCV*, ser. Lecture Notes in Computer Science, vol. 12351, 2020, pp. 565–580.
- [39] X. Liu, L. Wang, D. F. Wong, L. Ding, L. S. Chao, and Z. Tu, “Understanding and improving encoder layer fusion in sequence-to-sequence learning,” in *ICLR*, 2021.
- [40] E. Bugliarello, R. Cotterell, N. Okazaki, and D. Elliott, “Multimodal pretraining unmasked: A meta-analysis and a unified framework of vision-and-language bert,” *TACL*, vol. 9, pp. 978–994, 2021.
- [41] T. Lin, M. Maire, S. J. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft COCO: common objects in context,” in *ECCV*, ser. Lecture Notes in Computer Science, vol. 8693, 2014, pp. 740–755.
- [42] K. Pearson, “Liii. on lines and planes of closest fit to systems of points in space,” *The London, Edinburgh, and Dublin philosophical magazine and journal of science*, vol. 2, no. 11, pp. 559–572, 1901.