

Collaborative Transformers for Grounded Situation Recognition

Junhyeong Cho¹ Youngseok Yoon¹ Suha Kwak^{1,2}
 Department of CSE, POSTECH¹ Graduate School of AI, POSTECH²
 {junhyeong99, yys8646, suha.kwak}@postech.ac.kr

Abstract

Grounded situation recognition is the task of predicting the main activity, entities playing certain roles within the activity, and bounding-box groundings of the entities in the given image. To effectively deal with this challenging task, we introduce a novel approach where the two processes for activity classification and entity estimation are interactive and complementary. To implement this idea, we propose **Collaborative Glance-Gaze TransFormer (CoFormer)** that consists of two modules: *Glance* transformer for activity classification and *Gaze* transformer for entity estimation. *Glance* transformer predicts the main activity with the help of *Gaze* transformer that analyzes entities and their relations, while *Gaze* transformer estimates the grounded entities by focusing only on the entities relevant to the activity predicted by *Glance* transformer. Our CoFormer achieves the state of the art in all evaluation metrics on the SWiG dataset. Training code and model weights are available at <https://github.com/jhcho99/CoFormer>.

1. Introduction

Humans make decisions via dual systems of thinking as stated in the cognitive theory by Kahneman [14]. Those two systems are known to work in tandem and complement each other [8, 28]. Consider a comprehensive scene understanding task as a specific example of such decision making. As illustrated in Figure 1, humans cast a quick glance to figure out what is happening, and slowly gaze at details to analyze which objects are involved and how they are related. These two processes are mutually supportive, *e.g.*, understanding involved objects and their relations leads to more accurate recognition of the event depicted in the scene.

Inspired by this, we propose a collaborative framework which leverages the two processes for Grounded Situation Recognition (GSR) [29]. GSR is a comprehensive scene understanding task that is recently introduced as an extension of Situation Recognition (SR) [41]. The objective of SR is to produce a structured image summary that describes the main activity and entities playing certain roles within the

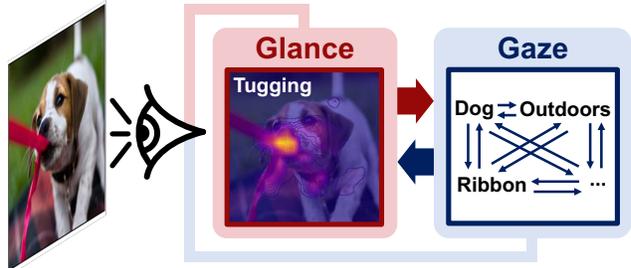


Figure 1. Two processes in comprehensive scene understanding. *Glance* figures out what is happening, and *Gaze* analyzes entities engaged in the main activity and their relations. In our CoFormer, these two processes are interactive and complementary.

activity, where the roles are predefined for each activity by a lexical database called FrameNet [7]. In GSR, those involved entities are grounded with bounding boxes; Figure 2 presents example results of GSR. Following conventions, we call an activity *verb* and an entity *noun* in this paper.

The common pipeline of SR and GSR in the literature [3, 4, 19, 27, 29, 32, 40, 41] resembles the two processes: predicting a verb (*Glance*), then estimating a noun for each role associated with the predicted verb (*Gaze*). Regarding this pipeline, correctness of the predicted verb is extremely important since noun estimation entirely depends on the predicted verb. If the result of verb prediction is incorrect, then estimated nouns cannot be correct either because the predicted verb determines the set of roles, *i.e.*, the basis of noun estimation. Moreover, verb prediction is challenging since a verb is highly abstract and situations for the same verb could significantly vary as shown in Figure 2. In spite of its importance and difficulty, verb prediction has been made in naïve ways, *e.g.*, using a single classifier on top of a convolutional neural network (CNN), which is analogous to *Glance* only. Existing methods allow *Glance* to assist *Gaze* by informing the predicted verb but not vice versa; this could limit the performance of verb prediction, and consequently, that of the entire pipeline.

We resolve the above issue by a collaborative framework that enables *Glance* and *Gaze* to interact and complement

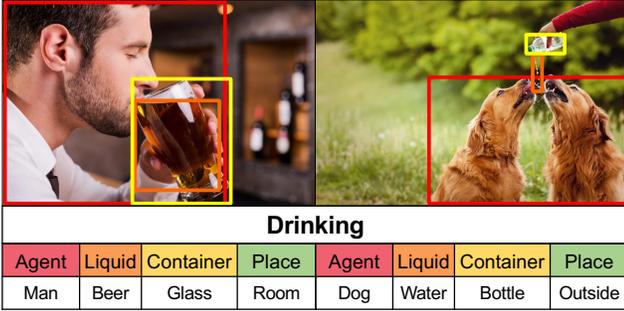


Figure 2. Two examples of Grounded Situation Recognition [29]. These show various situations for the same verb.

each other. To fully utilize this framework, we propose **Collaborative Glance-Gaze TransFormer** (CoFormer) that consists of Glance transformer and Gaze transformer as illustrated in Figure 3. Glance transformer predicts a verb by aggregating image features through self-attentions, and Gaze transformer estimates nouns and their groundings by allowing each role to focus on its relevant image region through self-attentions and cross-attentions. As shown in Figure 3, there are two steps for *Gaze* in our CoFormer. Gaze-Step1 transformer estimates nouns for all role candidates and assists Glance transformer for more accurate verb prediction. Meanwhile, Gaze-Step2 transformer estimates a noun and its grounding for each role associated with the predicted verb by exploiting the aggregated image features obtained by Glance transformer.

The collaborative relationship between Glance and Gaze transformers lead to more accurate verb and grounded noun predictions for GSR. In CoFormer, *Gaze-Step1* supports *Glance* by analyzing involved nouns and their relations, which enables noun-aware verb prediction. *Glance* assists *Gaze-Step2* by informing the predicted verb, which reduces the role candidates considered in grounded noun prediction. **Contributions.** (i) We propose a collaborative framework where the two processes for verb prediction and noun estimation are interactive and complementary, which is novel in GSR. (ii) Our method achieves state-of-the-art accuracy in every evaluation metric on the SWiG dataset. (iii) We demonstrate the effectiveness of CoFormer by conducting extensive experiments and provide in-depth analyses.

2. Related Work

Visual reasoning such as image captioning [2, 10, 13, 35, 42], scene graph generation [15, 26, 38, 39], and human-object-interaction detection [16, 21, 36, 43] has been widely studied for comprehensive understanding of images. Given an image, image captioning aims at describing activities and entities using natural language, and scene graph generation or human-object-interaction detection aims at capturing a set of triplets ⟨subject, predicate, object⟩ or

⟨human, object, interaction⟩. However, it is not straightforward to evaluate the quality of natural language captions, and the triplets have limited expressive power. To overcome such limitations, Yatskar *et al.* [41] introduce SR along with the *imSitu* dataset. SR has more expressive power based on linguistic sources from FrameNet [7], and its quality evaluation is straightforward. GSR builds upon SR by additionally estimating bounding-box groundings.

Situation Recognition. Yatskar *et al.* [41] propose a conditional random field [17] model, and also present a tensor composition method with semantic augmentation [40]. Mallya and Lazebnik [27] employ a recurrent neural network to capture role relations in the predefined sequential order. Li *et al.* [19] propose a gated graph neural network (GGNN) [20] to capture the relations in more flexible ways. To learn context-aware role relations depending on an input image, Suhail and Sigal [32] apply a mixture kernel method to GGNN. Cooray *et al.* [4] employ inter-dependent queries to capture role relations, and present a verb model which considers nouns from the two predefined roles; they construct a query based on two nouns for verb prediction. Compared with this, CoFormer considers nouns from all role candidates for accurate verb prediction.

Grounded Situation Recognition. Pratt *et al.* [29] propose GSR along with the *SWiG* dataset, and present two models: Independent Situation Localizer (ISL) and Joint Situation Localizer (JSL). They first predict a verb using a single classifier on top of a CNN backbone, then estimate nouns and their groundings. In both models, LSTM [12] produces output features to predict nouns in the predefined sequential order, while RetinaNet [23] estimates their groundings. ISL separately predicts nouns and their groundings, and JSL jointly predicts them. Cho *et al.* [3] propose a transformer encoder-decoder architecture, where the encoder effectively captures high-level semantic features for verb prediction and the decoder flexibly learns the role relations. Compared with these models, CoFormer leverages involved nouns and their relations for accurate verb prediction via transformers.

Transformer Architecture. Transformers [34] have driven remarkable success in vision tasks [1, 2, 6, 10, 16, 18, 24, 26]. Dosovitskiy *et al.* [6] propose a transformer encoder architecture for image classification by aggregating image features using a learnable token in the encoder. Carion *et al.* [1] present a transformer encoder-decoder architecture for object detection by predicting a set of bounding boxes using a fixed number of learnable queries in the decoder. Such learnable queries have been widely used to extract features in other transformer architectures [16, 18, 24]. Compared with those transformers, CoFormer employs two learnable tokens which aggregate different kinds of features through self-attentions. In addition, CoFormer constructs a different number of learnable queries by explicitly leveraging the prediction result obtained by two encoders and a classifier.

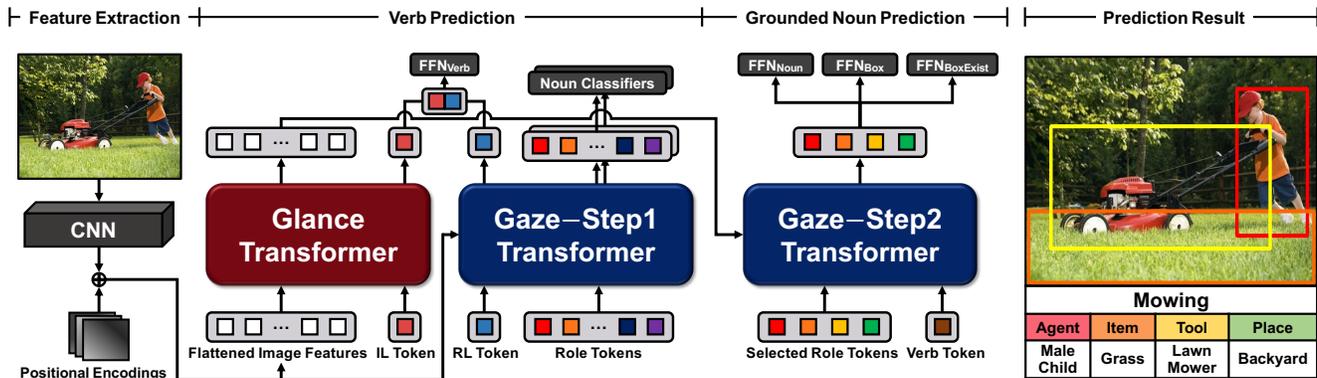


Figure 3. Overall architecture of **Collaborative Glance-Gaze TransFormer** (CoFormer). Glance transformer predicts a verb with the help of Gaze-Step1 transformer that analyzes nouns and their relations by leveraging role features, while Gaze-Step2 transformer estimates the grounded nouns for the roles associated with the predicted verb. Prediction results are obtained by feed forward networks (FFNs). The results from the two noun classifiers placed on top of Gaze-Step1 transformer are ignored at inference time.

3. Method

Task Definition. GSR assumes discrete sets of verbs \mathcal{V} , nouns \mathcal{N} , and roles \mathcal{R} . Each verb $v \in \mathcal{V}$ is paired with a frame derived from FrameNet [7], where the frame defines the set of roles $\mathcal{R}_v \subset \mathcal{R}$ associated with the verb. For example, a verb *Mowing* is paired with a frame which defines the set of roles $\mathcal{R}_{Mowing} = \{Agent, Item, Tool, Place\}$ as shown in Figure 3. Each role $r \in \mathcal{R}_v$ is fulfilled by a noun $n \in \mathcal{N}$ grounded by a bounding box $\mathbf{b} \in \mathbb{R}^4$, called *grounded noun*. Formally speaking, the set of fulfilled roles is $\mathcal{F}_v = \{(r_i, n_i, \mathbf{b}_i) \mid r_i \in \mathcal{R}_v, n_i \in \mathcal{N} \cup \{\emptyset_n\}, \mathbf{b}_i \in \mathbb{R}^4 \cup \{\emptyset_b\} \text{ for } i = 1, \dots, |\mathcal{R}_v|\}$; \emptyset_n and \emptyset_b denote *unknown* and *not grounded*, respectively. The output of GSR is a grounded situation denoted by $S = (v, \mathcal{F}_v)$.

3.1. Overall Architecture

CoFormer predicts a verb, then estimates grounded nouns as illustrated in Figure 3. As shown in Figure 5, our transformers consist of common building blocks, encoder and decoder, whose architectures are illustrated in Figure 6. For simplicity, we abbreviate Step1 as *S1*, and Step2 as *S2* in the remaining of this paper.

Overview. Given an image, CoFormer extracts flattened image features via a CNN backbone and flatten operation, which are fed as input to Glance transformer and Gaze-S1 transformer. From these transformers, the output features corresponding to Image-Looking (IL) and Role-Looking (RL) tokens are used for verb prediction. Considering the predicted verb, Gaze-S2 transformer estimates grounded nouns for the roles associated with the predicted verb by exploiting image features obtained by Glance transformer. Figure 4 shows the collaborative relationship between the modules; transformers for verb prediction and noun estimation are interactive and complementary in CoFormer.



Figure 4. Interactive and complementary processes in CoFormer. (a) RL token feature, (b) predicted verb, (c) loss gradients.

Glance Transformer. This transformer consists of a single encoder which takes the flattened image features and learnable IL token as input. IL token captures the essential features for verb prediction, while Glance transformer aggregates the image features through self-attentions.

Gaze-S1 Transformer. This transformer is composed of a decoder and an encoder. The decoder takes the flattened image features and learnable role tokens as input, where the role tokens correspond to all role candidates. This module extracts role features from the image features via the role tokens. Then, the encoder takes the role features and learnable RL token as input. RL token captures involved nouns and their relations for verb prediction, while the encoder aggregates the role features through self-attentions.

Gaze-S2 Transformer. This transformer consists of a single decoder, which takes learnable tokens and the aggregated image features obtained from Glance transformer as input. The input tokens correspond to the predicted verb and its associated roles. Note that a verb token is added to role tokens as shown in Figure 5; conditioning on the predicted verb significantly reduces the search space of the roles, e.g., the search space of *Mowing Tool* is much smaller than that of *Tool*. Gaze-S2 transformer extracts role features from the aggregated image features, and the extracted role features are used for grounded noun prediction.

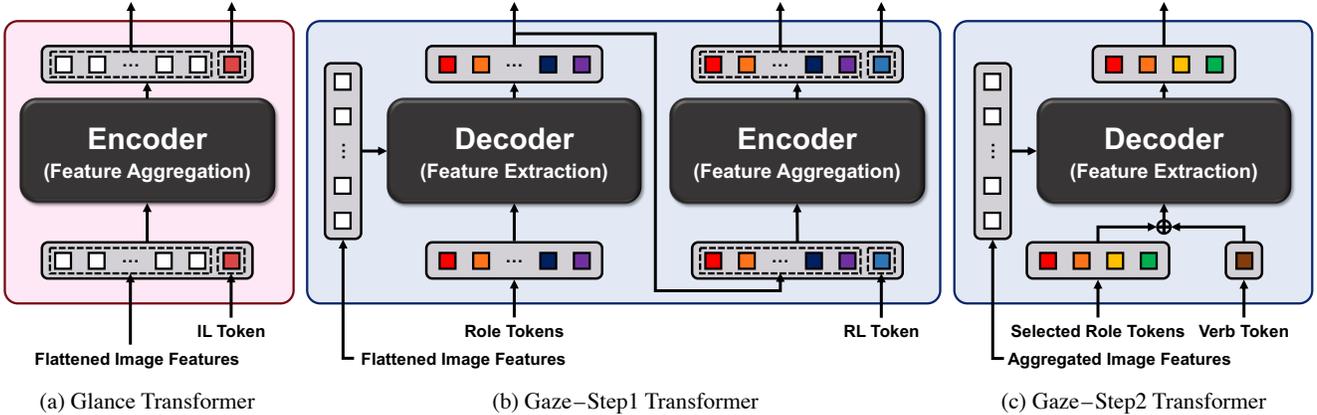


Figure 5. Transformer architectures in CoFormer are composed of common building blocks, encoder and decoder.

3.2. Feature Extraction

Given an input image, a single CNN backbone extracts image features of size $h \times w \times c$, where $h \times w$ is the resolution, and c is the number of channels. Then, a 1×1 convolution followed by a flatten operation produces flattened image features $\mathbf{X}_F \in \mathbb{R}^{hw \times d}$, where d is the number of channels. The flattened image features \mathbf{X}_F are fed as input to Glance transformer (Fig. 5(a)) and Gaze-S1 transformer (Fig. 5(b)). For the flattened image features \mathbf{X}_F , positional encodings are introduced to retain spatial information. As shown in Figure 6, positional encodings are added to the queries and keys at the self-attention layers in an encoder, and to the keys at the cross-attention layers in a decoder.

3.3. Verb Prediction

The input of the encoder in Glance transformer is obtained by the concatenation of the image features \mathbf{X}_F and learnable IL token. IL token captures the essential features for verb prediction, while the encoder aggregates the image features through self-attentions. As its output, the encoder produces aggregated image features $\mathbf{X}_A \in \mathbb{R}^{hw \times d}$ and IL token feature. For the aggregated image features \mathbf{X}_A , positional encodings are applied.

Gaze-S1 transformer supports Glance transformer for more accurate verb prediction, while predicting nouns for all role candidates. To be specific, the decoder of Gaze-S1 transformer takes the flattened image features \mathbf{X}_F and learnable role tokens corresponding to all predefined roles; each role token embedding is denoted by $\mathbf{w}_r \in \mathbb{R}^d$, where $r \in \mathcal{R}$. This decoder extracts role features through self-attentions on the role tokens and cross-attentions between the tokens and the image features \mathbf{X}_F . The input of the encoder in Gaze-S1 transformer is obtained by the concatenation of the extracted role features and learnable RL token. RL token captures involved nouns and their relations from all role candidates, while the encoder aggregates the

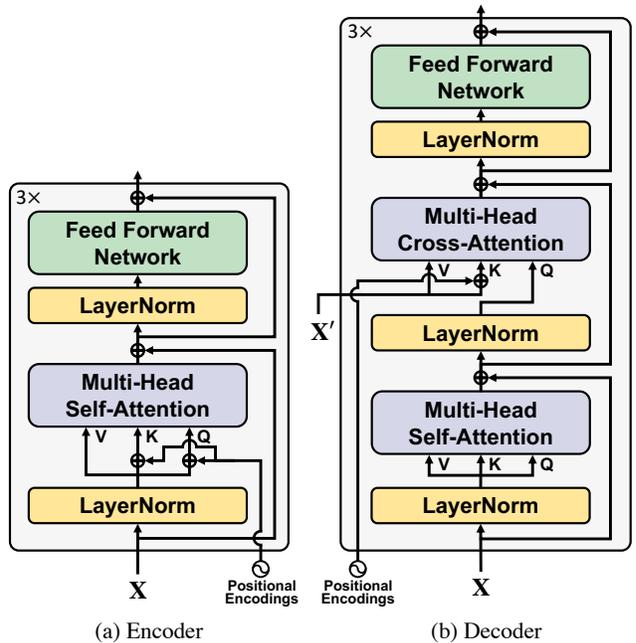


Figure 6. Detailed architectures of encoder and decoder. We use Pre-Layer Normalization [37] for these two modules. An encoder performs feature aggregation through self-attentions on \mathbf{X} , and a decoder performs feature extraction through self-attentions on \mathbf{X} and cross-attentions between \mathbf{X} and \mathbf{X}' .

role features through self-attentions. For this encoder, positional encodings are not added to the queries and keys at the self-attention layers since roles are permutation-invariant in GSR. Regarding to Gaze-S1 transformer, the extracted and aggregated role features are fed as input to noun classifiers; these classifiers are auxiliary modules and their results are ignored at inference time. Note that Gaze-S1 transformer assists Glance transformer via RL token feature which is aware of involved nouns and their relations.

IL token feature and RL token feature are concatenated, then fed as input to the feed forward network (FFN) for verb classification, which consists of learnable linear layers with activation functions. The verb classifier FFN_{Verb} followed by a softmax function produces a verb probability distribution \mathbf{p}_v , which is used to estimate the most probable verb $\hat{v} = \arg \max_v \mathbf{p}_v$. The predicted verb \hat{v} supports Gaze-S2 transformer so that the transformer concentrates only on the roles associated with the predicted verb and estimates their grounded nouns more accurately in consequence.

3.4. Grounded Noun Prediction

The aggregated image features \mathbf{X}_A from Glance transformer are fed as input to Gaze-S2 transformer (Fig. 5(c)). The decoder in this transformer takes the image features \mathbf{X}_A and *frame-role queries* as input. Specifically, for each role r in the frame of the predicted verb \hat{v} , its frame-role query $\mathbf{q}_r \in \mathbb{R}^d$ is constructed by an addition of the learnable role token embedding $\mathbf{w}_r \in \mathbb{R}^d$ and the learnable verb token embedding $\mathbf{w}_{\hat{v}} \in \mathbb{R}^d$, *i.e.*, $\mathbf{q}_r = \mathbf{w}_r + \mathbf{w}_{\hat{v}}$ for $r \in \mathcal{R}_{\hat{v}}$. The decoder extracts role features through self-attentions on the frame-role queries and cross-attentions between the queries and the image features \mathbf{X}_A to capture the involved nouns and their relations from roles relevant to the verb \hat{v} . Those extracted role features are used for grounded noun prediction. Note that this task requires to predict a noun, a bounding box, and a box existence for each role. Accordingly, we employ three feed forward networks FFN_{Noun} , FFN_{Box} , and $\text{FFN}_{\text{BoxExist}}$ that take the role features as input for noun classification, bounding box estimation, and box existence prediction, respectively. Each of these FFNs consists of learnable linear layers with activation functions.

For each role $r \in \mathcal{R}_{\hat{v}}$, FFN_{Noun} followed by a softmax function produces a noun probability distribution \mathbf{p}_{n_r} . FFN_{Box} followed by a sigmoid function produces a bounding box $\hat{\mathbf{b}}_r \in [0, 1]^4$ which indicates the center coordinates, height and width relative to the input image size. The predicted box $\hat{\mathbf{b}}_r$ can be transformed into the top-left and bottom-right coordinates $\hat{\mathbf{b}}'_r \in \mathbb{R}^4$. $\text{FFN}_{\text{BoxExist}}$ followed by a sigmoid function produces a box existence probability $p_{b_r} \in [0, 1]$. If $p_{b_r} < 0.5$, the predicted box $\hat{\mathbf{b}}'_r$ is ignored. Note that the predicted verb \hat{v} assists Gaze-S2 transformer via the construction of frame-role queries, while the loss gradients propagated from Gaze-S2 transformer through the aggregated image features \mathbf{X}_A enable Glance transformer to implicitly consider involved nouns.

3.5. Training CoFormer

The predicted verb, nouns and bounding boxes are used for computing losses to train CoFormer. At training time, we construct frame-role queries based on the ground-truth verb for stable training of Gaze-S2 transformer. Please refer to the supplementary material for more training details.

Verb Classification Loss. The verb classification loss is the cross-entropy between the verb probability distribution \mathbf{p}_v and the ground-truth verb distribution.

Noun Classification Losses. As illustrated in Figure 3, CoFormer has three noun classifiers; two of them are placed on top of Gaze-S1 transformer and the other is incorporated with Gaze-S2 transformer. For each noun classifier, we compute the cross-entropy between the estimated noun probability distribution and the ground-truth noun distribution for each role $r \in \mathcal{R}_{\tilde{v}}$, where \tilde{v} is the ground-truth verb. The computed cross-entropy loss is averaged over roles $R_{\tilde{v}}$. Note that we only train role tokens for the roles in the frame of the ground-truth verb \tilde{v} , since noun annotations are given for the roles associated with the verb \tilde{v} in the dataset.

Box Existence Prediction Loss. To deal with roles which have no ground-truth boxes (*i.e.*, \emptyset_b), *e.g.*, by occlusion, CoFormer estimates a box existence probability p_{b_r} for each role $r \in \mathcal{R}_{\tilde{v}}$. The box existence prediction loss is the cross-entropy between the probability p_{b_r} and the ground-truth box existence, which is averaged over roles $R_{\tilde{v}}$.

Box Regression Losses. We employ $L1$ loss and GIoU loss [30] for box regression. Let \mathbf{b}_r denote the ground-truth box in the form of the center coordinates, height and width relative to the given image size. In the computation of box regression losses, we ignore the roles which have no ground-truth boxes (*i.e.*, \emptyset_b). The $L1$ box regression loss \mathcal{L}_{L1} is computed by

$$\mathcal{L}_{L1} = \frac{1}{|\tilde{\mathcal{R}}|} \sum_{r \in \tilde{\mathcal{R}}} \|\mathbf{b}_r - \hat{\mathbf{b}}_r\|_1, \quad (1)$$

where $\tilde{\mathcal{R}} = \{r \mid \mathbf{b}_r \neq \emptyset_b \text{ for } r \in R_{\tilde{v}}\}$. To compute the GIoU loss, $\text{GIoU}(\cdot)$ is first computed by

$$\begin{aligned} \text{GIoU}(\mathbf{b}'_r, \hat{\mathbf{b}}'_r) &= \frac{|\mathbf{b}'_r \cap \hat{\mathbf{b}}'_r|}{|\mathbf{b}'_r \cup \hat{\mathbf{b}}'_r|} - \frac{|C(\mathbf{b}'_r, \hat{\mathbf{b}}'_r) \setminus (\mathbf{b}'_r \cup \hat{\mathbf{b}}'_r)|}{|C(\mathbf{b}'_r, \hat{\mathbf{b}}'_r)|}, \quad (2) \end{aligned}$$

where \mathbf{b}'_r indicates the top-left and bottom-right coordinates transformed from \mathbf{b}_r , and $C(\mathbf{b}'_r, \hat{\mathbf{b}}'_r)$ denotes the smallest box which encloses \mathbf{b}'_r and $\hat{\mathbf{b}}'_r$. The GIoU box regression loss $\mathcal{L}_{\text{GIoU}}$ is then computed by

$$\mathcal{L}_{\text{GIoU}} = \frac{1}{|\tilde{\mathcal{R}}|} \sum_{r \in \tilde{\mathcal{R}}} \left(1 - \text{GIoU}(\mathbf{b}'_r, \hat{\mathbf{b}}'_r)\right). \quad (3)$$

4. Experiments

CoFormer is evaluated on the SWiG dataset [29], which is constructed by adding box annotations to the imSitu dataset [41]. The imSitu dataset contains 75K, 25K and 25K images for train, development and test set, respectively. This dataset contains 504 verbs, 11K nouns and 190 roles.

Set	Method	Top-1 Predicted Verb					Top-5 Predicted Verbs					Ground-Truth Verb			
		verb	value	value-all	grnd value	grnd value-all	verb	value	value-all	grnd value	grnd value-all	value	value-all	grnd value	grnd value-all
<i>Methods for Situation Recognition</i>															
dev	CRF [41]	32.25	24.56	14.28	–	–	58.64	42.68	22.75	–	–	65.90	29.50	–	–
	CRF w/ DataAug [40]	34.20	26.56	15.61	–	–	62.21	46.72	25.66	–	–	70.80	34.82	–	–
	RNN w/ Fusion [27]	36.11	27.74	16.60	–	–	63.11	47.09	26.48	–	–	70.48	35.56	–	–
	GraphNet [19]	36.93	27.52	19.15	–	–	61.80	45.23	29.98	–	–	68.89	41.07	–	–
	CAQ w/ RE-VGG [4]	37.96	30.15	18.58	–	–	64.99	50.30	29.17	–	–	73.62	38.71	–	–
	Kernel GraphNet [32]	43.21	35.18	19.46	–	–	68.55	56.32	30.56	–	–	73.14	41.68	–	–
<i>Methods for Grounded Situation Recognition</i>															
	ISL [29]	38.83	30.47	18.23	22.47	7.64	65.74	50.29	28.59	36.90	11.66	72.77	37.49	52.92	15.00
	JSL [29]	39.60	31.18	18.85	25.03	10.16	67.71	52.06	29.73	41.25	15.07	73.53	38.32	57.50	19.29
	GSRTTR [3]	41.06	32.52	19.63	26.04	10.44	69.46	53.69	30.66	42.61	15.98	74.27	39.24	58.33	20.19
	CoFormer (Ours)	44.41	35.87	22.47	29.37	12.94	72.98	57.58	34.09	46.70	19.06	76.17	42.11	61.15	23.09
<i>Methods for Situation Recognition</i>															
test	CRF [41]	32.34	24.64	14.19	–	–	58.88	42.76	22.55	–	–	65.66	28.96	–	–
	CRF w/ DataAug [40]	34.12	26.45	15.51	–	–	62.59	46.88	25.46	–	–	70.44	34.38	–	–
	RNN w/ Fusion [27]	35.90	27.45	16.36	–	–	63.08	46.88	26.06	–	–	70.27	35.25	–	–
	GraphNet [19]	36.72	27.52	19.25	–	–	61.90	45.39	29.96	–	–	69.16	41.36	–	–
	CAQ w/ RE-VGG [4]	38.19	30.23	18.47	–	–	65.05	50.21	28.93	–	–	73.41	38.52	–	–
	Kernel GraphNet [32]	43.27	35.41	19.38	–	–	68.72	55.62	30.29	–	–	72.92	42.35	–	–
<i>Methods for Grounded Situation Recognition</i>															
	ISL [29]	39.36	30.09	18.62	22.73	7.72	65.51	50.16	28.47	36.60	11.56	72.42	37.10	52.19	14.58
	JSL [29]	39.94	31.44	18.87	24.86	9.66	67.60	51.88	29.39	40.60	14.72	73.21	37.82	56.57	18.45
	GSRTTR [3]	40.63	32.15	19.28	25.49	10.10	69.81	54.13	31.01	42.50	15.88	74.11	39.00	57.45	19.67
	CoFormer (Ours)	44.66	35.98	22.22	29.05	12.21	73.31	57.76	33.98	46.25	18.37	75.95	41.87	60.11	22.12

Table 1. Quantitative evaluations of methods in SR and GSR. SR models are evaluated on the imSitu dataset, and GSR models are evaluated on the SWiG dataset. The only difference between the two datasets is the existence of bounding box annotation.

Method	Top-1 Predicted Verb			Top-5 Predicted Verbs			Ground-Truth Verb			
	verb	value	grnd value	verb	value	grnd value	value	value-all	grnd value	grnd value-all
w/o Gaze-S1 Transformer	42.46	34.21	28.23	70.89	55.47	45.34	76.02	41.96	61.21	23.15
w/o Gaze-S2 Transformer	43.02	31.24	23.27	71.17	51.70	36.59	69.68	32.94	48.44	13.05
w/o Noun Classifiers on Gaze-S1 Transformer	41.30	33.33	27.50	69.76	55.05	44.96	75.97	41.94	61.32	23.39
w/o Gradient Flow from Gaze-S2 Transformer to Glance Transformer	42.96	33.82	25.77	70.97	54.59	41.11	73.91	38.59	55.10	17.10
w/o Verb Token in Gaze-S2 Transformer	44.36	35.57	29.16	72.84	56.79	46.19	74.53	39.83	60.07	21.83
CoFormer (Ours)	44.41	35.87	29.37	72.98	57.58	46.70	76.17	42.11	61.15	23.09

Table 2. Ablation study of CoFormer on the SWiG dev set. The contributions of different components used in our model are evaluated.

The number of roles in the frame of a verb ranges from 1 to 6. Each image is paired with the annotation of a verb, and three nouns from three different annotators for each role. In addition to this annotation, the SWiG dataset provides a box annotation for each role (except role *Place*).

4.1. Evaluation Metric

Metric Details. The prediction accuracy of verb is measured by *verb*, that of noun is evaluated by *value* and *value-all*, and that of grounded noun is assessed by *grounded-value* and *grounded-value-all*. Regarding to the noun metrics, *value* measures whether a noun is correct for each role, and *value-all* measures whether all nouns are correct for entire roles in a frame simultaneously. The noun prediction is considered correct if the predicted noun matches any of the three noun annotations given by three annotators. For the grounded noun metrics, *grounded-value* measures whether a noun and its grounding are correct for each role, and *grounded-value-all* measures whether all nouns and their groundings are correct for entire roles in a

frame simultaneously. The grounding prediction is considered correct if the predicted box existence is correct and the predicted bounding box has Intersection-over-Union (IoU) value at least 0.5 with the box annotation. Note that the above metrics are calculated per verb and then averaged over all verbs, since the number of roles in a frame depends on a verb and each verb might be associated with a different number of samples in the dataset.

Evaluation Settings. Three evaluation settings are proposed for comprehensive evaluation: *Top-1 Predicted Verb*, *Top-5 Predicted Verbs*, and *Ground-Truth Verb*. In *Top-1 Predicted Verb* setting, the predicted nouns and their groundings are considered incorrect if the top-1 verb prediction is incorrect. In *Top-5 Predicted Verbs* setting, the predicted nouns and their groundings are considered incorrect if the ground-truth verb is not contained in the top-5 predicted verbs. In *Ground-Truth Verb* setting, the predicted nouns and their groundings are obtained by conditioning on the ground-truth verb.

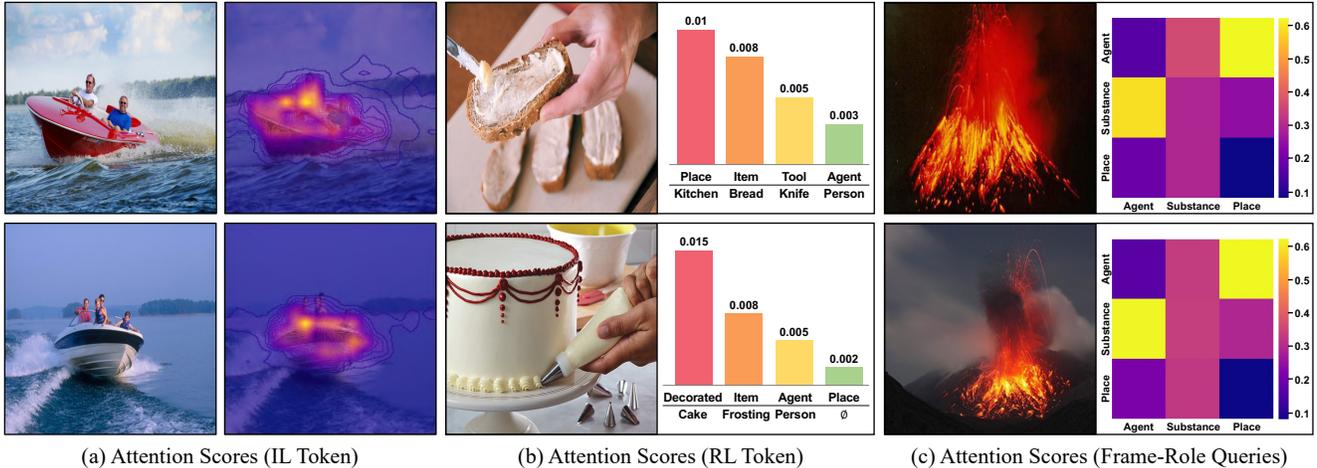


Figure 7. Attention scores from IL token to image features, from RL token to role features, and on frame-role queries. We visualize the attention scores computed from the last self-attention layer of the encoder in Glance transformer, the encoder in Gaze-S1 transformer, and the decoder in Gaze-S2 transformer, respectively. Higher attention scores are highlighted in red color on images.

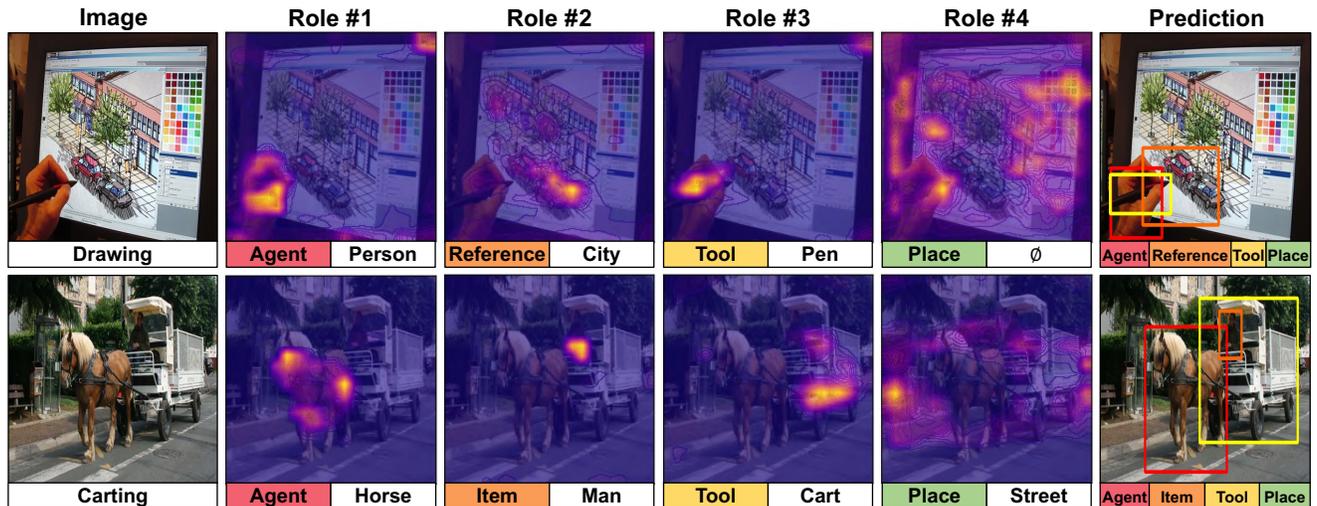


Figure 8. Attention scores from frame-role queries to image features. We visualize the attention scores computed from the last cross-attention layer of the decoder in Gaze-S2 transformer. Higher attention scores are highlighted in red color on images.

4.2. Implementation Details

We use ResNet-50 [11] pretrained on ImageNet [5] as a CNN backbone following existing models [3, 29] in GSR. Given an image, the CNN backbone extracts image features of size $h \times w \times c$, where $h = w = 22$ and $c = 2048$. The embedding dimension of each token is $d = 512$. We employ AdamW Optimizer [25] with 10^{-4} weight decay, $\beta_1 = 0.9$, and $\beta_2 = 0.999$. We train CoFormer with 10^{-4} learning rate (10^{-5} for the CNN backbone) which decreases by a factor of 10 at epoch 30. Training CoFormer with batch size of 16 for 40 epochs takes about 30 hours on four RTX 3090 GPUs. Complete details including loss coefficients are provided in the supplementary material.

4.3. Quantitative Evaluations

CoFormer achieves the state of the art in all evaluations as shown in Table 1. Existing SR models [4, 19, 27, 32] use at least two VGG-16 [31] backbones, and GSR models [29] employ two ResNet-50 [11] backbones for verb and noun prediction, while CoFormer only employs a single ResNet-50 backbone. Compared with GSRTR [3], the improvements in the verb prediction accuracies range from 3.35%p to 4.03%p. Regarding to the noun prediction accuracies, the improvements range from 1.84%p to 3.89%p, and those in the grounded noun prediction accuracies range from 2.11%p to 4.09%p. These results demonstrate that the proposed collaborative framework is effective for GSR.

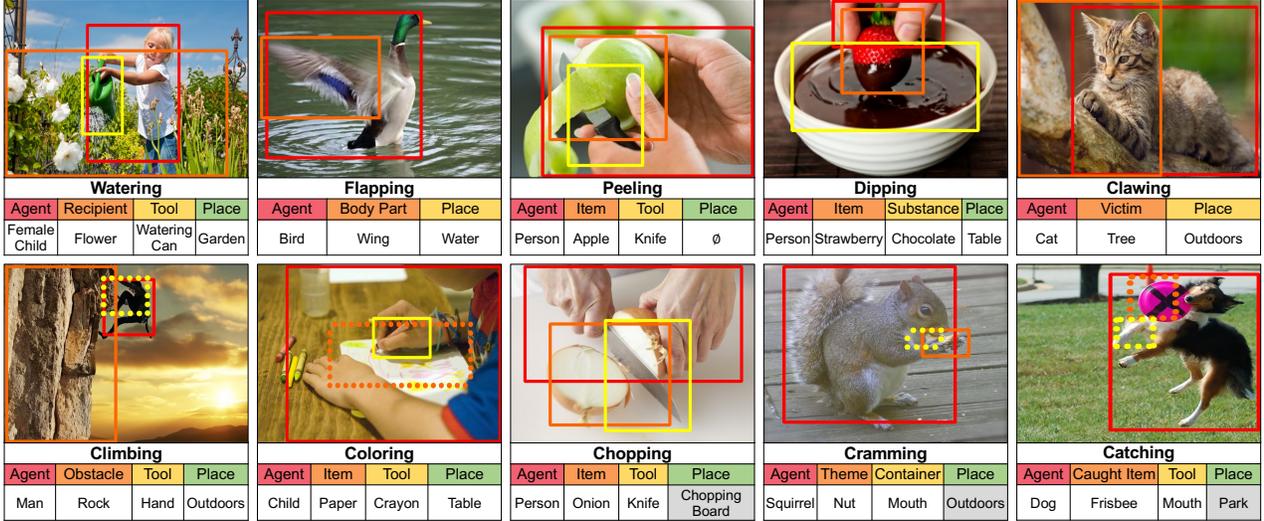


Figure 9. Prediction results. Dashed boxes denote incorrect grounding predictions. Incorrect noun predictions are highlighted in gray color.

Ablation Study. We analyze the effects of different components in CoFormer as shown in Table 2. When we train our model without using Gaze-S1 transformer or Gaze-S2 transformer, the accuracies in verb prediction or grounded noun prediction largely decrease, which demonstrates the effectiveness of the collaborative framework. Training our CoFormer without using the two noun classifiers placed on top of Gaze-S1 transformer leads to significant drops in the verb prediction accuracies. In this case, it is difficult for role features to learn involved nouns and their relations, while the encoder in Gaze-S1 transformer aggregates the role features through self-attentions. To figure out whether Gaze-S2 transformer assists Glance transformer by forcing it to implicitly consider involved nouns, we train CoFormer by restricting the flow of loss gradients through the aggregated image features from Gaze-S2 transformer to Glance transformer. As shown in the fourth row of Table 2, the verb prediction accuracies drop, which demonstrates that Gaze-S2 transformer supports Glance transformer via loss gradients through the aggregated image features. In CoFormer, each frame-role query is constructed by an addition of a role token embedding and a verb token embedding. We study how effective it is by training CoFormer without using a verb token embedding for the construction of frame-role queries. The fifth row of Table 2 shows that the grounded noun prediction accuracies drop, which demonstrates that the verb token embedding is helpful for grounded noun prediction.

4.4. Qualitative Evaluations

We visualize the attention scores computed in the attention layers of CoFormer. Figure 7(a) shows that IL token captures the essential features to estimate a verb for two *Boating* images. Figure 7(b) shows how much RL token focuses on the roles in the frame of the ground-truth

verb, and the classification results from the noun classifier placed on top of the encoder in Gaze-S1 transformer; attention scores among 190 roles sum to 1. This demonstrates that RL token effectively captures involved nouns and their relations through self-attentions in the encoder of Gaze-S1 transformer. Figure 7(c) shows how role relations are captured through self-attentions on frame-role queries, which demonstrates that CoFormer similarly captures the relations if the situations in images are similar; attention scores sum to 1 in each column. Figure 8 shows the local regions where frame-role queries focus on, and the predicted grounded nouns corresponding to the queries. Figure 9 shows prediction results of CoFormer on the SWiG test set. The first row shows the correct predictions, and the second row shows several incorrect predictions.

5. Conclusion

We propose a collaborative framework for GSR, where the two processes for verb prediction and noun estimation interact and complement each other. Using this framework, we present CoFormer which outperforms existing methods in all evaluation metrics on the SWiG dataset. We also provide in-depth analyses of how CoFormer draws attentions on images and captures role relations with the ablation study on the effects of different components used in our model. A limitation of CoFormer is that the model sometimes suffers from predicting the boxes which have extreme aspect ratios or small scales. This issue will be explored in future work.

Acknowledgement. This work was supported by the NRF grant and the IITP grant funded by Ministry of Science and ICT, Korea (NRF-2021R1A2C3012728, No.2019-0-01906 Artificial Intelligence Graduate School Program–POSTECH, No.2021-0-02068 Artificial Intelligence Innovation Hub, IITP-2020-0-00842).

References

- [1] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-End Object Detection with Transformers. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 213–229, 2020. [2](#)
- [2] Long Chen, Zhihong Jiang, Jun Xiao, and Wei Liu. Human-Like Controllable Image Captioning With Verb-Specific Semantic Roles. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16846–16856, 2021. [2](#)
- [3] Junhyeong Cho, Youngseok Yoon, Hyeonjun Lee, and Suha Kwak. Grounded Situation Recognition with Transformers. In *Proceedings of the British Machine Vision Conference (BMVC)*, 2021. [1](#), [2](#), [6](#), [7](#), [11](#), [16](#)
- [4] Thilini Cooray, Ngai-Man Cheung, and Wei Lu. Attention-Based Context Aware Reasoning for Situation Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4736–4745, 2020. [1](#), [2](#), [6](#), [7](#)
- [5] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. [7](#), [11](#)
- [6] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *International Conference on Learning Representations (ICLR)*, 2021. [2](#)
- [7] Charles J. Fillmore, Christopher R. Johnson, and Miriam R.L. Petruck. Background to Framenet. *International Journal of Lexicography*, 16(3):235–250, 2003. [1](#), [2](#), [3](#)
- [8] Benjamin Gardner and Amanda L. Rebar. Habit Formation and Behavior Change. *Oxford research encyclopedia of psychology*, 2019. [1](#)
- [9] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 249–256, 2010. [12](#)
- [10] Longteng Guo, Jing Liu, Xinxin Zhu, Peng Yao, Shichen Lu, and Hanqing Lu. Normalized and Geometry-Aware Self-Attention Network for Image Captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. [2](#)
- [11] He, Kaiming and Zhang, Xiangyu and Ren, Shaoqing and Sun, Jian. Deep Residual Learning for Image Recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. [7](#), [11](#), [16](#)
- [12] Sepp Hochreiter and Jürgen Schmidhuber. Long Short-Term Memory. *Neural Computation*, 9(8):1735–1780, 1997. [2](#)
- [13] Lun Huang, Wenmin Wang, Jie Chen, and Xiao-Yong Wei. Attention on Attention for Image Captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4634–4643, 2019. [2](#)
- [14] Daniel Kahneman. Maps of Bounded Rationality: Psychology for Behavioral Economics. *The American Economic Review*, 93(5):1449–1475, 2003. [1](#)
- [15] Siddhesh Khandelwal, Mohammed Suhail, and Leonid Sigal. Segmentation-Grounded Scene Graph Generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 15879–15889, 2021. [2](#)
- [16] Bumsoo Kim, Junhyun Lee, Jaewoo Kang, Eun-Sol Kim, and Hyunwoo J. Kim. HOTR: End-to-End Human-Object Interaction Detection With Transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 74–83, 2021. [2](#)
- [17] John Lafferty, Andrew McCallum, and Fernando Pereira. Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data. In *International Conference on Machine Learning (ICML)*, pages 282–289, 2001. [2](#)
- [18] Jinwoo Lee, Hyunsung Go, Hyunjoon Lee, Sunghyun Cho, Minhyuk Sung, and Junho Kim. CTRL-C: Camera Calibration Transformer With Line-Classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 16228–16237, 2021. [2](#)
- [19] Ruiyu Li, Makarand Tapaswi, Renjie Liao, Jiaya Jia, Raquel Urtasun, and Sanja Fidler. Situation Recognition with Graph Neural Network. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 4173–4182, 2017. [1](#), [2](#), [6](#), [7](#)
- [20] Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard Zemel. Gated Graph Sequence Neural Networks. In *International Conference on Learning Representations (ICLR)*, 2016. [2](#)
- [21] Yong-Lu Li, Siyuan Zhou, Xijie Huang, Liang Xu, Ze Ma, Hao-Shu Fang, Yanfeng Wang, and Cewu Lu. Transferable Interactiveness Knowledge for Human-Object Interaction Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3585–3594, 2019. [2](#)
- [22] Tsung-Yi Lin, Piotr Dollar, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature Pyramid Networks for Object Detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2117–2125, 2017. [16](#)
- [23] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal Loss for Dense Object Detection. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 2980–2988, 2017. [2](#)
- [24] Songhua Liu, Tianwei Lin, Dongliang He, Fu Li, Ruifeng Deng, Xin Li, Errui Ding, and Hao Wang. Paint Transformer: Feed Forward Neural Painting With Stroke Prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6598–6607, 2021. [2](#)
- [25] Ilya Loshchilov and Frank Hutter. Decoupled Weight Decay Regularization. In *International Conference on Learning Representations (ICLR)*, 2019. [7](#), [12](#)
- [26] Yichao Lu, Himanshu Rai, Jason Chang, Boris Knyazev, Guangwei Yu, Shashank Shekhar, Graham W. Taylor, and

- Maksims Volkovs. Context-Aware Scene Graph Generation With Seq2Seq Transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 15931–15941, 2021. 2
- [27] Arun Mallya and Svetlana Lazebnik. Recurrent Models for Situation Recognition. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 455–463, 2017. 1, 2, 6, 7
- [28] Ellen Peters, Daniel Västfjäll, Paul Slovic, C.K. Mertz, Ketti Mazzocco, and Stephan Dickert. Numeracy and Decision Making. *Psychological Science*, 17(5):407–413, 2006. 1
- [29] Sarah Pratt, Mark Yatskar, Luca Weihs, Ali Farhadi, and Aniruddha Kembhavi. Grounded Situation Recognition. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 314–332, 2020. 1, 2, 5, 6, 7, 11, 16
- [30] Rezaatofighi, Hamid and Tsoi, Nathan and Gwak, JunYoung and Sadeghian, Amir and Reid, Ian and Savarese, Silvio. 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 (CVPR)*, pages 658–666, 2019. 5
- [31] Simonyan, Karen and Zisserman, Andrew. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations (ICLR)*, 2015. 7
- [32] Mohammed Suhail and Leonid Sigal. Mixture-Kernel Graph Attention Network for Situation Recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 10363–10372, 2019. 1, 2, 6, 7
- [33] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the Inception Architecture for Computer Vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2818–2826, 2016. 12
- [34] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is All you Need. In *Advances in Neural Information Processing Systems (NIPS)*, 2017. 2
- [35] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and Tell: A Neural Image Caption Generator. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015. 2
- [36] Tiancai Wang, Tong Yang, Martin Danelljan, Fahad Shahbaz Khan, Xiangyu Zhang, and Jian Sun. Learning Human-Object Interaction Detection Using Interaction Points. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [37] Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang, Yanyan Lan, Liwei Wang, and Tieyan Liu. On Layer Normalization in the Transformer Architecture. In *International Conference on Machine Learning (ICML)*, pages 10524–10533. PMLR, 2020. 4
- [38] Danfei Xu, Yuke Zhu, Christopher B Choy, and Li Fei-Fei. Scene Graph Generation by Iterative Message Passing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5410–5419, 2017. 2
- [39] Jianwei Yang, Jiasen Lu, Stefan Lee, Dhruv Batra, and Devi Parikh. Graph R-CNN for Scene Graph Generation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 670–685, 2018. 2
- [40] Mark Yatskar, Vicente Ordonez, Luke Zettlemoyer, and Ali Farhadi. Commonly Uncommon: Semantic Sparsity in Situation Recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7196–7205, 2017. 1, 2, 6
- [41] Mark Yatskar, Luke Zettlemoyer, and Ali Farhadi. Situation Recognition: Visual Semantic Role Labeling for Image Understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5534–5542, 2016. 1, 2, 5, 6
- [42] Quanzeng You, Hailin Jin, Zhaowen Wang, Chen Fang, and Jiebo Luo. Image Captioning With Semantic Attention. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2
- [43] Frederic Z. Zhang, Dylan Campbell, and Stephen Gould. Spatially Conditioned Graphs for Detecting Human-Object Interactions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 13319–13327, 2021. 2

Collaborative Transformers for Grounded Situation Recognition

— Supplementary Material —

This supplementary material provides method details (Section A), implementation details (Section B), qualitative evaluations (Section C), an application of this task (Section D), computational evaluations (Section E) and a limitation (Section F), which could not be included in the main paper due to the limited space.

A. Method Details

Transformer architectures in our CoFormer consist of common building blocks, encoder and decoder. The main components of these building blocks are attention layers. Section A.1 provides more details of the attention layers.

In Section 3.5 of the main paper, the losses to train our model are described: verb classification loss, noun classification losses, box existence prediction loss, and box regression losses. Section A.2 provides more details of the losses.

A.1. Attention Layer

Multi-Head Attention. The input of the multi-head attention layer is the sequence of query, key and value. The query sequence is denoted by $\mathbf{Q} \in \mathbb{R}^{L_Q \times d}$, where L_Q is the sequence length and d is the size of the hidden dimension. The key sequence is denoted by $\mathbf{K} \in \mathbb{R}^{L_{KV} \times d}$ and value sequence is denoted by $\mathbf{V} \in \mathbb{R}^{L_{KV} \times d}$, where L_{KV} is the sequence length. In the multi-head attention layer, we employ H attention heads; the hidden dimension of each attention head is $d_h = d/H$. For each attention head i , \mathbf{Q} , \mathbf{K} and \mathbf{V} are linearly projected via parameter matrices $\mathbf{W}_i^Q, \mathbf{W}_i^K, \mathbf{W}_i^V \in \mathbb{R}^{d \times d_h}$. In details,

$$\mathbf{Q}_i = \mathbf{Q}\mathbf{W}_i^Q \in \mathbb{R}^{L_Q \times d_h}, \quad (\text{A.1})$$

$$\mathbf{K}_i = \mathbf{K}\mathbf{W}_i^K \in \mathbb{R}^{L_{KV} \times d_h}, \quad (\text{A.2})$$

$$\mathbf{V}_i = \mathbf{V}\mathbf{W}_i^V \in \mathbb{R}^{L_{KV} \times d_h}. \quad (\text{A.3})$$

The output of each attention head i is obtained by a weighted summation of the value \mathbf{V}_i , where the weights are computed by the scaled dot-product between the query \mathbf{Q}_i and the key \mathbf{K}_i followed by a softmax function. In details,

$$\text{Attention}(\mathbf{Q}_i, \mathbf{K}_i, \mathbf{V}_i) = \text{Softmax}\left(\frac{\mathbf{Q}_i\mathbf{K}_i^T}{\sqrt{d_h}}\right)\mathbf{V}_i. \quad (\text{A.4})$$

The output of each attention head i is concatenated along hidden dimension, then linearly projected via a parameter

matrix $\mathbf{W}^O \in \mathbb{R}^{d \times d}$. In details,

$$\text{MultiHead}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = [\text{Head}_1; \dots; \text{Head}_H]\mathbf{W}^O, \quad (\text{A.5})$$

where $[\cdot]$ is a concatenation along hidden dimension and $\text{Head}_i = \text{Attention}(\mathbf{Q}_i, \mathbf{K}_i, \mathbf{V}_i)$ for $i = 1, \dots, H$.

Multi-Head Cross-Attention. This is the multi-head attention layer where the key sequence \mathbf{K} is same with the value sequence \mathbf{V} , but the query sequence \mathbf{Q} is different.

Multi-Head Self-Attention. This is the multi-head attention layer where the query sequence \mathbf{Q} , key sequence \mathbf{K} , and value sequence \mathbf{V} are same, *i.e.*, $\mathbf{Q} = \mathbf{K} = \mathbf{V}$.

A.2. Loss

Figure A1 shows the losses to train CoFormer. The verb classification loss is denoted by $\mathcal{L}_{\text{Verb}}$. The noun classification loss from the classifier involved in the decoder of Gaze-S1 transformer is denoted by $\mathcal{L}_{\text{Noun}}^1$, the loss from the classifier involved in the encoder of Gaze-S1 transformer is denoted by $\mathcal{L}_{\text{Noun}}^2$, and the loss from the classifier involved in the decoder of Gaze-S2 transformer is denoted by $\mathcal{L}_{\text{Noun}}^3$. The box existence prediction loss is denoted by $\mathcal{L}_{\text{BoxExist}}$. The $L1$ box regression loss is denoted by \mathcal{L}_{L1} . The GIoU box regression loss is denoted by $\mathcal{L}_{\text{GIoU}}$.

The total training loss is the linear combination of $\mathcal{L}_{\text{Verb}}, \mathcal{L}_{\text{Noun}}^1, \mathcal{L}_{\text{Noun}}^2, \mathcal{L}_{\text{Noun}}^3, \mathcal{L}_{\text{BoxExist}}, \mathcal{L}_{L1}$, and $\mathcal{L}_{\text{GIoU}}$. In this total loss, the loss coefficients are as follows: $\lambda_{\text{Verb}}, \lambda_{\text{Noun}}^1, \lambda_{\text{Noun}}^2, \lambda_{\text{Noun}}^3, \lambda_{\text{BoxExist}}, \lambda_{L1}, \lambda_{\text{GIoU}} > 0$.

B. Implementation Details

In Section 4.2 of the main paper, some implementation details are described. For completeness, we describe more architecture details (Section B.1), loss details (Section B.2), augmentation details (Section B.3), and training details (Section B.4) of our CoFormer.

B.1. Architecture Details

Following previous work [3, 29], we use ResNet-50 [11] pretrained on ImageNet [5] as a CNN backbone. Given an image, the CNN backbone produces image features of size $h \times w \times c$, where $h = w = 22$ and $c = 2048$. A 1×1 convolution followed by a flatten operation produces flattened image features $\mathbf{X}_F \in \mathbb{R}^{hw \times d}$, where $d = 512$. To retain spatial information, we employ positional encodings. We use learnable 2D embeddings for the positional encodings.

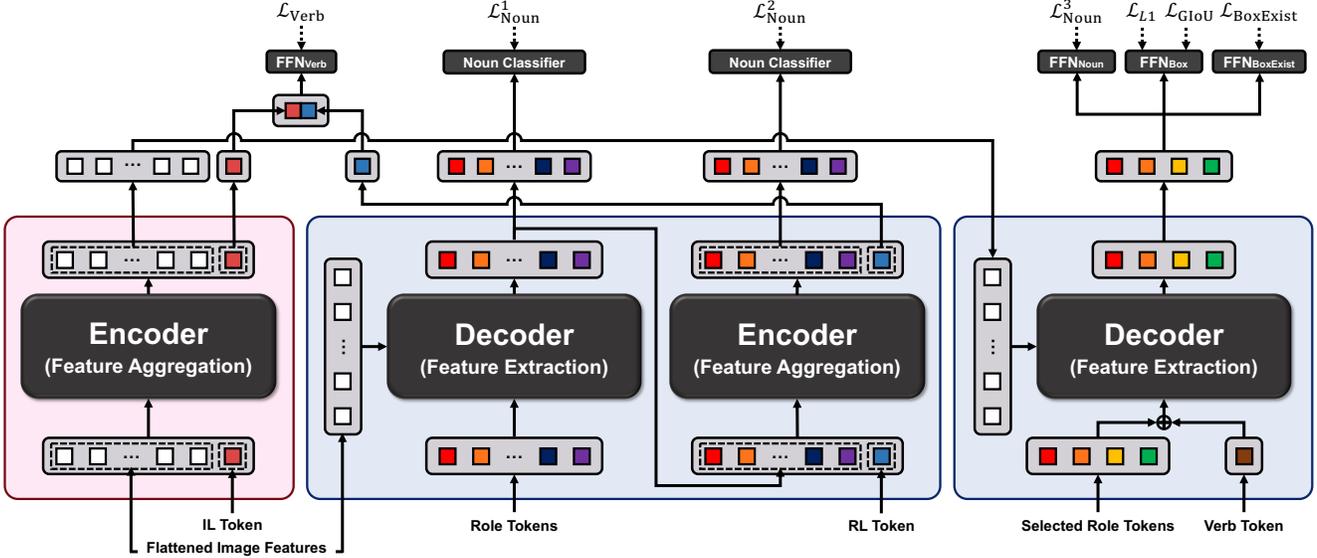


Figure A1. Transformer architectures in CoFormer including the losses to train our model. The losses for training our CoFormer are as follows: \mathcal{L}_{Verb} , \mathcal{L}_{Noun}^1 , \mathcal{L}_{Noun}^2 , \mathcal{L}_{Noun}^3 , $\mathcal{L}_{BoxExist}$, \mathcal{L}_{L1} , \mathcal{L}_{GIoU} .

We initialize encoders and decoders using Xavier Initialization [9], and these modules are trained with the dropout rate of 0.15. The number of heads in the attention layers of these modules is 8. Each of feed forward networks in these modules is 2-fully connected layers with a ReLU activation function, whose hidden dimensions are $4d$ and dropout rate is 0.15. These modules take learnable tokens, and each embedding dimension of the tokens is d .

The verb classifier FFN_{Verb} is 2-fully connected layers with a ReLU activation function, whose hidden dimensions are $2d$ and dropout rate is 0.3. Each of the two noun classifiers placed on top of Gaze-S1 transformer is a linear layer. The noun classifier FFN_{Noun} is 2-fully connected layers with a ReLU activation function, whose hidden dimensions are $2d$ and dropout rate is 0.3. The bounding box estimator FFN_{Box} is 3-fully connected layers with two ReLU activation functions, whose hidden dimensions are $2d$ and dropout rate is 0.2. The box existence predictor $FFN_{BoxExist}$ is 2-fully connected layers with a ReLU activation function, whose hidden dimensions are $2d$ and dropout rate is 0.2.

B.2. Loss Details

Complete Details of Noun Losses. In the SWiG dataset, each image is associated with three noun annotations given by three different annotators for each role. For the noun classification losses \mathcal{L}_{Noun}^1 , \mathcal{L}_{Noun}^2 , \mathcal{L}_{Noun}^3 , each noun loss is obtained by the summation of three classification losses corresponding to three different annotators.

Regularization. We employ label smoothing regularization [33] in the loss computation for verb classification loss \mathcal{L}_{Verb} and noun classification losses \mathcal{L}_{Noun}^1 , \mathcal{L}_{Noun}^2 , \mathcal{L}_{Noun}^3 .

In details, the label smoothing factor in the computation of verb classification loss is 0.3, and the factor in the computation of noun classification losses is 0.2.

Loss Coefficients. Total loss to train CoFormer is a linear combination of losses. In our implementation, the loss coefficients are $\lambda_{Verb} = \lambda_{Noun}^3 = 1$, $\lambda_{Noun}^1 = \lambda_{Noun}^2 = 2$, and $\lambda_{BoxExist} = \lambda_{L1} = \lambda_{GIoU} = 5$.

B.3. Augmentation Details

For data augmentation, we employ random scaling, random horizontal flipping, random color jittering, and random gray scaling. The input images are randomly scaled with the scaling factors of 0.5, 0.75, and 1.0. Also, the input images are horizontally flipped with the probability of 0.5. The brightness, saturation and hue of the input images are randomly changed with the factor of 0.1 for each change. The input images are randomly converted to grayscale with the probability of 0.3.

B.4. Training Details

We employ AdamW Optimizer [25] with the weight decay of 10^{-4} , $\beta_1 = 0.9$, and $\beta_2 = 0.999$. For stable training, we apply gradient clipping with the maximal gradient norm of 0.1. The transformers, classifiers and learnable embeddings are trained with the learning rate of 10^{-4} . The CNN backbone is fine-tuned with the learning rate of 10^{-5} . Note that we have a learning rate scheduler and the learning rates are divided by 10 at epoch 30. For batch training, we set the batch size to 16. We train CoFormer for 40 epochs, which takes about 30 hours on four RTX 3090 GPUs.

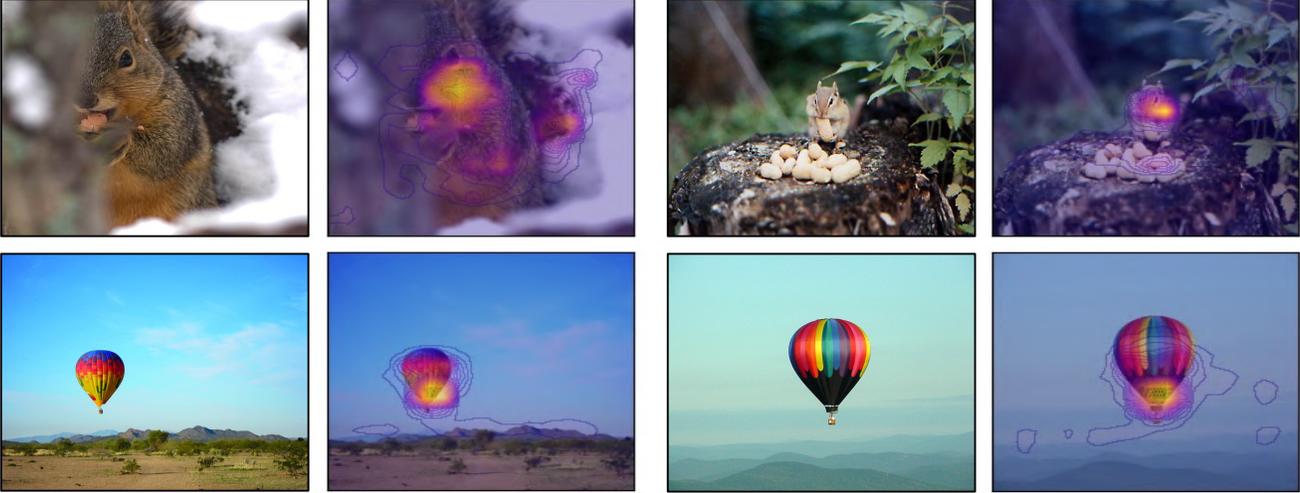


Figure C2. Attention scores from IL token to image features. We visualize the attention scores computed from the last self-attention layer of the encoder in Glance transformer. Higher attention scores are highlighted in red color on images.

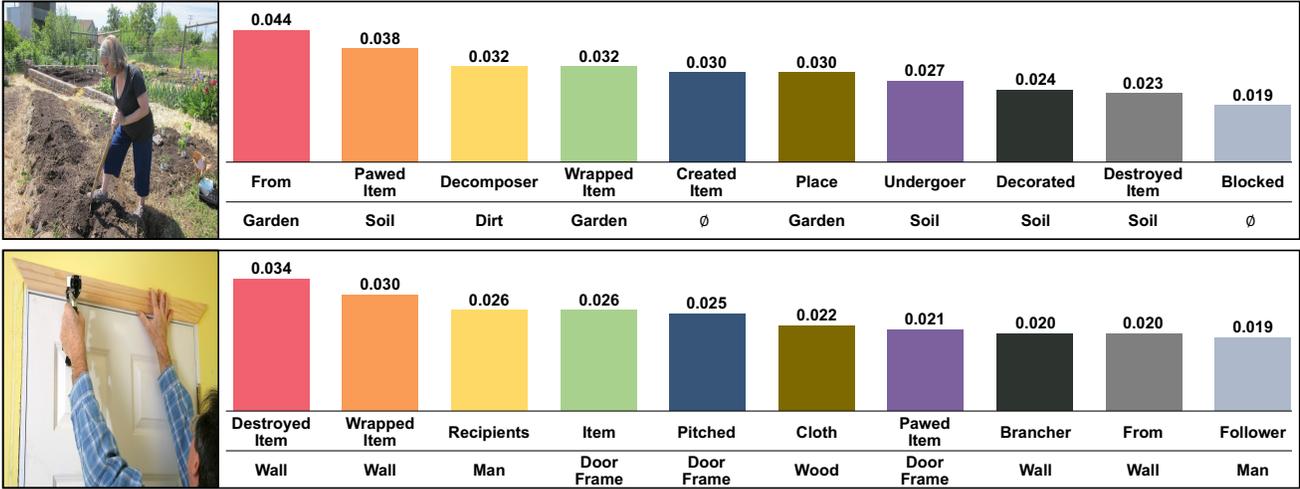


Figure C3. Attention scores from RL token to role features. We visualize the attention scores computed from the last self-attention layer of the encoder in Gaze-S1 transformer. Note that we show the roles where RL token has top-10 attention scores. In Figure 7(b) of the main paper, we show the results corresponding to the roles in the frame of the ground-truth verb.

C. Qualitative Evaluations

We visualize the attention scores computed in the attention layers of the transformers in our CoFormer. Figure C2 shows that IL token captures the essential features to estimate the main activities for two *Cramming* images and two *Ballooning* images. Figure C3 shows the roles where RL token has top-10 attention scores, and the classification results from the noun classifier placed on top of the encoder in Gaze-S1 transformer; attention scores among 190 roles sum to 1. Note that several roles where RL token has high attention scores are not relevant to the main activity, but the noun classification results corresponding to those roles are

highly relevant to the activity. Since RL token leverages the role features which are fed as input to the noun classifier, it is reasonable to aggregate those role features for accurate verb prediction; the role features are aware of involved nouns and their relations. Figure C3 demonstrates that RL token can effectively capture involved nouns and their relations for noun-aware verb prediction through self-attentions on the role features in the encoder of Gaze-S1 transformer. Figure C4 shows how role relations are captured through self-attentions on frame-role queries, which demonstrates that CoFormer similarly captures the role relations if the situations in images are similar; attention scores sum to 1

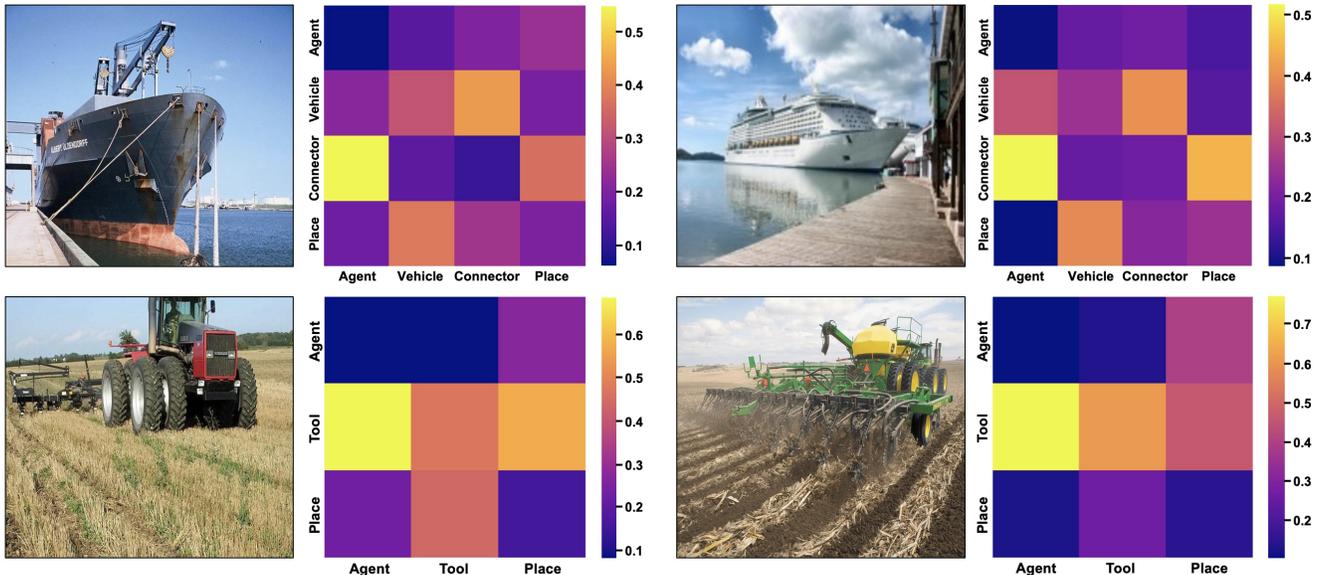


Figure C4. Attention scores on frame-role queries. We visualize the attention scores computed from the last self-attention layer of the decoder in Gaze-S2 transformer.



Figure C5. Attention scores from frame-role queries to image features. We visualize the attention scores computed from the last cross-attention layer of the decoder in Gaze-S2 transformer. Higher attention scores are highlighted in red color on images.

in each column. Figure C5 shows the local regions where frame-role queries focus on, and the predicted grounded nouns corresponding to the queries. This demonstrates that each query effectively captures its relevant local regions through cross-attentions between the queries and image features in the decoder of Gaze-S2 transformer. Note that those queries are constructed by leveraging the predicted verb, which significantly reduces the number of role candidates handled in noun estimation; Gaze-S1 transformer considers all role candidates, but Gaze-S2 transformer handles a few roles associated with the predicted verb. Figure C6 shows

the prediction results of CoFormer on the SWiG test set. The first and second row show correct prediction results. The third and fourth row show incorrect prediction results. As shown in Figure C6, three noun annotations are given for each role in the SWiG dataset. Note that the noun prediction is considered correct if the predicted noun matches any of the three noun annotations. The grounded noun prediction is considered correct if a noun, a bounding box, and box existence are correctly predicted for a role.

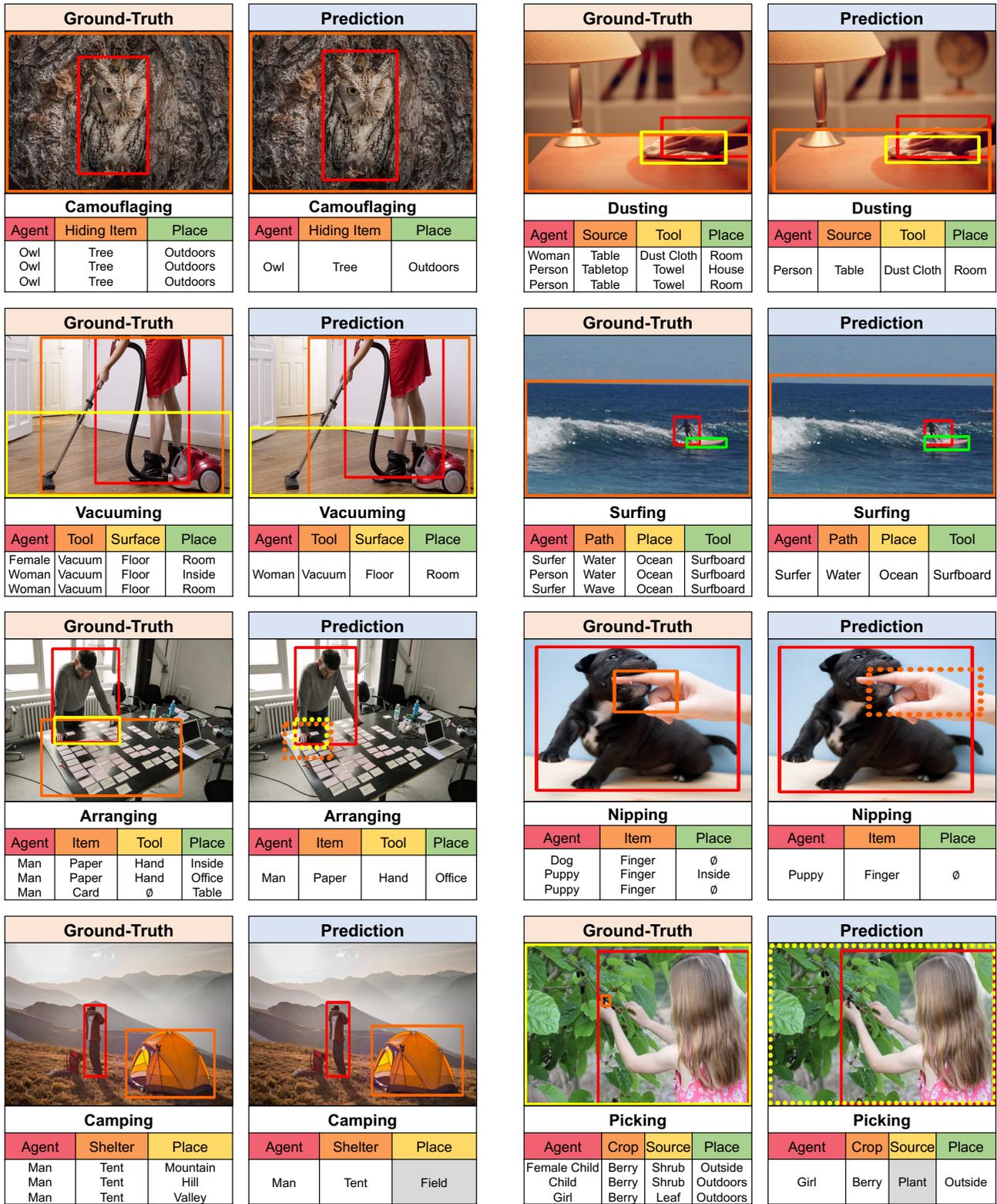


Figure C6. Prediction results of our CoFormer on the SWiG test set. Dashed boxes denote incorrect grounding predictions. Incorrect noun predictions are highlighted in gray color.



Figure D7. Grounded semantic aware image retrieval on the SWiG dev set. For each query image, we show the retrieval results which have top-5 similarity scores computed by GrSitSim(-) [29]. This retrieval computes the similarity between two images considering the predicted verbs, nouns, and bounding-box groundings of the nouns.

$$\text{GrSitSim}(I, J) = \max \left\{ \frac{\mathbb{1}[\hat{v}_i^I = \hat{v}_j^J]}{i \cdot j \cdot |\mathcal{R}_{\hat{v}_i^I}|} \sum_{k=1}^{|\mathcal{R}_{\hat{v}_i^I}|} \mathbb{1}[\hat{n}_{i,k}^I = \hat{n}_{j,k}^J] \cdot \left(1 + \text{IoU}(\hat{\mathbf{b}}_{i,k}^I, \hat{\mathbf{b}}_{j,k}^J)\right) \mid 1 \leq i, j \leq 5 \right\}. \quad (\text{D.6})$$

D. Application

As shown in Figure D7, we can apply GSR models to grounded semantic aware image retrieval. This image retrieval computes the similarity between two images considering their grounded situations. In details, a similarity score between an image I and an image J is computed by GrSitSim(I, J) (Eq. D.6). Given an image I , a GSR model predicts the top-5 most probable verbs $\hat{v}_1^I, \dots, \hat{v}_5^I$. For each predicted verb \hat{v}_i^I , the model predicts nouns $\hat{n}_{i,1}^I, \dots, \hat{n}_{i,|\mathcal{R}_{\hat{v}_i^I}|}^I$ and bounding boxes $\hat{\mathbf{b}}_{i,1}^I, \dots, \hat{\mathbf{b}}_{i,|\mathcal{R}_{\hat{v}_i^I}|}^I$. These prediction results are used in the computation of GrSitSim(I, J). By this score function, the similarity score is maximized if the top-1 predicted verb and the predicted grounded nouns are same for the two images I and J . Using this retrieval, we can retrieve images which have similar grounded situations with the situation of a query image.

E. Computational Evaluations

The number of parameters and inference time of our CoFormer are shown in Table E1. We also evaluate JSL [29] and GSRTR [3] on the SWiG test set using a single 2080Ti GPU with a batch size of 1. JSL uses two ResNet-50 [11] and a feature pyramid network (FPN) [22] in the CNN backbone, while GSRTR and our CoFormer only employ a single ResNet-50 in the backbone; these two models demand much shorter inference time than JSL, which is crucial for real-world applications. GSRTR and CoFormer are trained in an end-to-end manner, but JSL is trained separately in terms of verb model and grounded noun model.

Method	Backbone	#Params	Inference Time
JSL [29]	R50, R50-FPN	108 M	80.23 ms (12.46 FPS)
GSRTR [3]	R50	83 M	21.69 ms (46.10 FPS)
CoFormer (Ours)	R50	93 M	30.62 ms (32.66 FPS)

Table E1. Number of parameters and inference time. Inference time was measured on the SWiG test set using one 2080Ti GPU.

Metric	Area (width \times height)			Aspect Ratio (width/height)		
	0-10%	10-20%	20-100%	0-5%	5-95%	95-100%
value	66.82	69.68	78.64	72.75	76.24	71.88
grnd value	7.42	25.38	65.49	36.88	62.62	31.01

Table F2. Quantitative analysis of our CoFormer on the SWiG dev set in Ground-Truth Verb evaluation setting. The effects of box scales and aspect ratios are evaluated. Each range denotes the ratio of ground-truth boxes when sorted by the value of area or aspect ratio in ascending order.

F. Limitation

As shown in Table F2, CoFormer suffers from estimating the noun labels and boxes for objects which have small scales (Area 0-10% and 10-20%) or extreme aspect ratios (Aspect Ratio 0-5% and 95-100%). To overcome such limitation, one may leverage multi-scale image features.