

# Locality-Aware Zero-Shot Human-Object Interaction Detection

Sanghyun Kim      Deunsol Jung      Minsu Cho

Pohang University of Science and Technology (POSTECH), South Korea

{sanghuyn.kim, deunsol.jung, mscho}@postech.ac.kr

<http://cvlab.postech.ac.kr/research/LAIN>

## Abstract

Recent methods for zero-shot Human-Object Interaction (HOI) detection typically leverage the generalization ability of large Vision-Language Model (VLM), i.e., CLIP, on unseen categories, showing impressive results on various zero-shot settings. However, existing methods struggle to adapt CLIP representations for human-object pairs, as CLIP tends to overlook fine-grained information necessary for distinguishing interactions. To address this issue, we devise, LAIN, a novel zero-shot HOI detection framework designed to enhance the locality and interaction awareness of CLIP representations. The locality awareness, which involves capturing fine-grained details and the spatial structure of individual objects, is achieved by aggregating the information and spatial priors of adjacent neighborhood patches. The interaction awareness, which involves identifying whether and how a human is interacting with an object, is achieved by capturing the interaction pattern between the human and the object. By infusing locality and interaction awareness into CLIP representations, LAIN captures detailed information about the human-object pairs. Our extensive experiments on existing benchmarks show that LAIN outperforms previous methods in various zero-shot settings, demonstrating the importance of locality and interaction awareness for effective zero-shot HOI detection.

## 1. Introduction

The task of Human-Object Interaction (HOI) detection aims to localize human-object pairs and recognize the interactions between them in a given image, i.e., identifying a set of HOI instances (human, object, interaction). HOI detection is useful for a wide range of computer vision applications, including image retrieval [11, 50, 55] and image captioning [16, 51, 54], where a comprehensive understanding of human-object relationships is essential. Although significant advances have been made recently, conventional HOI methods [22, 23, 26] have primarily been relying on fully

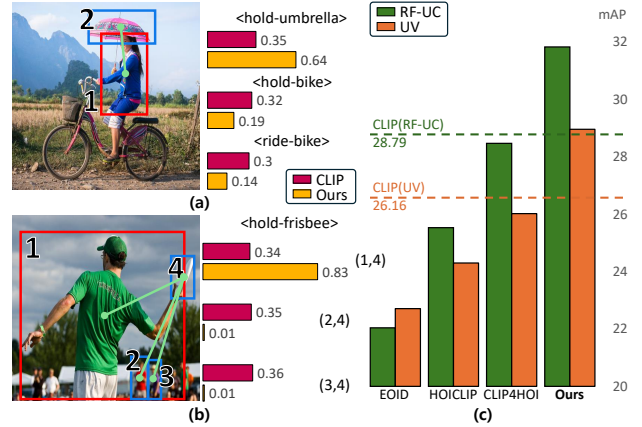


Figure 1. (a)-(b): Since CLIP primarily encodes global information, it struggles to capture the fine-grained details required to accurately identify interactions within human-object pairs. (c): When existing methods adapt CLIP representations to zero-shot HOI detection, this limitation hinders CLIP’s generalization, and results in degraded performance which is even lower than CLIP’s original zero-shot performance in UC-RF and UV settings.

supervised learning, limited to identifying predefined HOI categories. Given that humans interact with objects in a compositional way, it is costly and impractical to collect annotations for all possible HOI categories, limiting their ability to identify novel HOI categories not present in the training set. Recently, strong generalization ability of CLIP on unseen categories, which stems from contrastive image-level pre-training on large-scale data, has inspired the development of a zero-shot HOI model that leverages this capability to recognize unseen HOI categories.

While existing methods [36, 37, 52] for zero-shot HOI detection have achieved strong performance by leveraging CLIP representations, the domain gap between the image-level pre-training task and the region-level task poses the challenges in adapting CLIP representations to zero-shot HOI detection. Since CLIP predominantly encodes global information [61, 63], it often fails to extract fine-grained information about individual objects. This hinders the HOI model from capturing whether and how the person inter-

acts with the object. As shown in Figure 1 (a), CLIP assigns a high confidence score to the interaction ‘ride-bike’ even though the human-object pair (1,2) does not focus on the bike region, indicating that CLIP has limited capacity to capture fine-grained object details. As a result, the human-object pair (1,4) has a similar confidence score for the unseen HOI category ‘hold-frisbee’ compared to others, *i.e.*, (2,3) and (3,4), struggling to distinguish the interactive pair in Figure 1 (b). These issues weaken the generalization ability of CLIP when existing methods attempt to adapt its representations for zero-shot HOI detection, resulting in lower zero-shot performance compared to CLIP’s original results as shown in Figure 1 (c).

To address the challenges mentioned above, we introduce a novel zero-shot HOI detection framework, dubbed **Locality-Aware Interaction Network (LAIN)**, that learns locality-aware interaction via **Locality Adapter (LA)** and **Interaction Adapter (IA)**. The LA extracts locality-aware features from image patch tokens by considering visual context of neighboring regions and spatial priors, and then infuses them back into the image patch tokens. The IA leverages the locality-aware patch tokens to update human-object tokens by performing interaction reasoning between human and object regions, resulting in interaction-aware human-object tokens. In this manner, the LA provides fine-grained details and spatial structure for individual objects, enabling the IA to perform effective contextual reasoning. The IA complements CLIP representations by incorporating fine-grained details, providing a relational context that cannot be captured by locality awareness alone. By incorporating locality and interaction awareness, which play complementary roles, each layer of LAIN effectively captures detailed information for the human-object pair, facilitating the adaptation of CLIP representations to zero-shot HOI detection.

To demonstrate the effectiveness of our proposed method, we conducted extensive evaluations on two public benchmarks, HICO-DET [2] and V-COCO [12]. The experimental results show that LAIN outperforms the previous methods for zero-shot detection across all zero-shot settings, demonstrating the robust generalization ability of our approach in zero-shot scenarios. Our ablation studies demonstrate the importance of locality and interaction awareness for zero-shot HOI detection.

Our contribution can be summarized as follows:

- We propose the Locality-Aware Interaction Network (LAIN), which incorporates a Locality Adapter (LA) and an Interaction Adapter (IA).
- By enriching CLIP representations with locality and interaction awareness, LAIN effectively captures the fine-grained details about human-object pairs.
- Extensive experiments demonstrate that LAIN achieves outstanding zero-shot performance, achieving a new state-of-the-art.

## 2. Related work

### 2.1. Human-Object Interaction (HOI) Detection

Conventional HOI detection methods can be roughly divided into two categories: two-stage and one-stage methods. Two-stage methods [9, 10, 17, 28, 38, 44, 47, 58–60] first detect humans and objects using a pre-trained detector [1, 41]. After constructing all possible human-object pairs based on the detection results, these pairs are fed into an interaction classifier. To generate discriminative features for classifying the interaction of a human-object pair, they incorporate additional information [13, 27, 33] and perform relational reasoning on a graph structure [44, 47, 58]. In contrast to two-stage methods that follow a sequential cascade to determine the interaction between the human-object pairs, one-stage methods [8, 20, 29] concurrently detect individual instances, pair the human-object instances, and classify interactions. Inspired by the transformer-based detector, *i.e.*, DETR [1], where each query learns to detect an object, recent one-stage methods [22, 26, 42, 62] have adopted transformer-based structures, where each query predicts a (human, object, interaction) triplet. Despite their promising results, these approaches heavily rely on full annotations with predefined HOI categories, which makes them impractical for handling unseen HOI categories.

### 2.2. Vision tasks with CLIP

CLIP [39] is a multimodal framework that adopts contrastive learning to jointly train image and text encoders on large-scale image-text pairs found from the web. By leveraging vision and language knowledge pre-trained on large-scale data, CLIP has significantly improved the zero-shot capabilities of models across various downstream tasks, including out-of-distribution detection [7, 48] and segmentation [5, 63]. However, CLIP struggles to align local image regions with text descriptions since CLIP was trained by aligning whole images with their corresponding text descriptions in a common embedding space, thus producing suboptimal results on region-level tasks [5, 35, 53, 61, 63]. To mitigate the issue, recent work [4, 5, 34, 53, 61] learns a locality in the CLIP or ViT structure by training on large scale of the region-text pairs [4, 5, 61], and introduces attention mechanism to capture the information of local regions [14, 34]. However, they require additional pre-training stages to adapt the downstream task, and require modifications to the network structure which limits the utilization of CLIP’s pre-trained knowledge.

### 2.3. Zero-shot HOI detection

Zero-shot HOI detection aims to detect both seen HOI categories available during training and unseen HOI categories that do not appear during training. Previous work [17–19, 33] mainly adopts compositional learning, which disen-

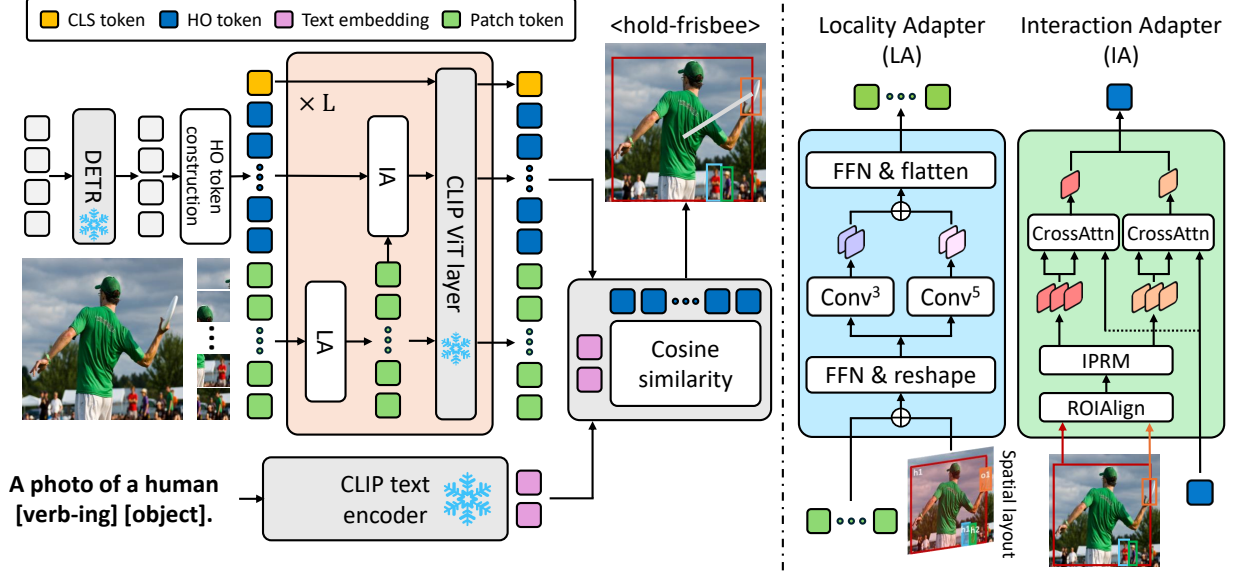


Figure 2. The overall architecture of LAIN. All valid human-object pairs are constructed and embedded into HO tokens based on detection results from a pre-trained DETR [1]. Image patch tokens are passed through the Locality Adapter (LA), which infuses locality awareness into each patch token. The updated patch tokens and HO tokens are then passed through the Interaction Adapter (IA), which enhances each HO token with interaction awareness. The HO, CLS, patch tokens are subsequently refined by the frozen  $l$ -th ViT layer of the CLIP [39] visual encoder. After repeating this process for  $L$  layers, HOI scores are computed by measuring the cosine similarity between the HO tokens and text embeddings extracted from CLIP text encoder.

tangles HOI representation into object and interaction features. Although the disentangled features enable the model to recognize unseen combinations, these methods are limited to combinations where either the object or the verb is not shown in training, since it is infeasible to learn the disentangled features of unseen objects or verbs. In light of the success of CLIP [39] on various zero-shot settings, recent work [30, 45, 52] focuses on transferring CLIP knowledge via the teacher-student architecture. However, since they transfer knowledge only about seen classes, they tend to be biased toward seen category samples. Furthermore, since this knowledge, *i.e.*, CLIP scores, only conveys which interaction occurs, the model struggles to learn fine-grained local details such as human attire. To enhance transferability, prior work [23, 24, 36, 37] directly leverages CLIP representations. ADA-CM [23] uses object queries from DETR as global object priors to inform all objects in the image, primarily providing global context to CLIP rather than local details of individual objects. CMMP [24] uses spatial priors to model plausible spatial configurations, focusing only on spatial relationships rather than the relational context of an HO pair. BCOM [46] incorporates additional knowledge extracted from the detector backbone along with CLIP representations to better capture the small-scale interaction.

In contrast, our method enhances CLIP representations by integrating locality awareness for fine-grained object details and interaction awareness for the relational context of the HO pair, enabling better adaptation of CLIP representations to zero-shot HOI detection.

### 3. Method

#### 3.1. Overview

Given an input image  $I$ , the goal of zero-shot HOI detection is to predict all HOI instances, including those belonging to HOI categories that are unseen during training. Formally,  $i$ -th HOI instance is defined 4-tuple  $(b_i^h, b_i^o, c_i, a_i)$ , where  $b_i^h, b_i^o \in \mathbb{R}^4$  represent the bounding box coordinates of the human and object, respectively.  $c_i \in \mathbb{O}$  and  $a_i \in \mathbb{V}$  represent the object class and the interaction class that occurs between a human and an object, *i.e.* verb, where  $\mathbb{O} = \{o_1, \dots, o_{N_o}\}$  is a set of objects and  $\mathbb{V} = \{v_1, \dots, v_{N_v}\}$  is a set of interactions.  $N_o$  and  $N_v$  are the number of object and interaction classes, respectively. Under the zero-shot settings, the model is trained only on the samples from seen HOI categories  $\mathbb{C}_{\text{seen}} = \mathbb{C} \setminus \mathbb{C}_{\text{unseen}}$ , where  $\mathbb{C} = \{(o_i, v_j) | o_i \in \mathbb{O}, v_j \in \mathbb{V}\}$  is the set of all possible HOI categories.

The overall architecture of LAIN is illustrated in Figure 2. To detect all possible HOI instances, we first detect the objects in the input image  $I$  using a pre-trained detector, *i.e.*, DETR [1]. Based on the detection results, we exhaustively construct all valid human-object pairs and generate HO tokens which are used to classify the interaction for each corresponding human-object pair. The HO tokens, along with the [CLS] token and the image patch tokens, are then passed through the CLIP visual encoder to aggregate visual information for each human-object pair. We attach the Locality Adapter (LA) and Interaction Adapter (IA) to

the front of each CLIP visual encoder layer, to inject locality and interaction awareness into the CLIP representations, respectively. After repeating this process for  $L$  layers, the HOI score for each HO token is computed by measuring the similarity to text embeddings of HOI categories.

### 3.2. HO Token Construction

In this section, we introduce HO tokens for HOI detection, which aggregate and convey contextual information from an input image similar to the [CLS] token in CLIP. Each HO token is constructed based on detection results to determine the interaction of the corresponding human-object pair.

Specifically, following the two-stage methods [36, 59], given an image  $I$ , we use an off-the-shelf object detector, *i.e.*, DETR [1], to obtain detection results  $\{(b_i, c_i, s_i, g_i)\}_{i=1}^{N_{\text{det}}}$ , where  $b_i \in \mathbb{R}^4$ ,  $c_i \in \mathbb{O}$ ,  $s_i \in \mathbb{R}^1$  and  $g_i \in \mathbb{R}^{D_{\text{det}}}$  represent the bounding box, object class, confidence score, and object feature, respectively.  $N_{\text{det}}$  and  $D_{\text{det}}$  denote the number of detected objects and dimension of the object feature, respectively. Based on the detected results, HO tokens  $T \in \mathbb{R}^{N_{\text{pair}} \times D_{\text{clip}}}$  for all valid human-object pairs are constructed as follows:

$$\text{idx} = \{(u, v) \mid u \neq v, c_u = \text{'human'}\}, \quad (1)$$

$$T_i = \text{FFN}([g_u; g_v]), \text{ where } (u, v) = \text{idx}_i, \quad (2)$$

where  $N_{\text{pair}}$  is the number of all valid human-object pairs and  $D_{\text{clip}}$  is the feature dimension of the CLIP. The HO tokens  $T$  are concatenated with the [CLS] token and image patch tokens  $F$ , and then passed through the CLIP visual encoder composed of  $L$  layers to aggregate the visual information:

$$[T_{(l)}, \text{cls}_{(l)}, F_{(l)}] = \mathcal{V}_{(l)}([T_{(l-1)}; \text{cls}_{(l-1)}; F_{(l-1)}]), \quad (3)$$

where  $\text{cls} \in \mathbb{R}^{1 \times D_{\text{clip}}}$  and  $F \in \mathbb{R}^{H \times W \times D_{\text{clip}}}$  denote [CLS] token and image patch tokens, respectively, and  $\mathcal{V}_{(l)}$  indicates the  $l$ -th layer of the CLIP visual encoder.  $H$  and  $W$  are the height and width of the feature map before flattening. For simplicity, we omit the layer index  $l$  in the following sections.

### 3.3. Locality Adapter

To determine the interaction between an HO pair, it is crucial to recognize fine-grained details of the individual object. For example, if the model learns during training that a human is wearing a helmet in the seen HOI category 'ride-bike', this fine-grained information can help to identify the unseen HOI category 'ride-snowboard'. However, while CLIP representation effectively captures global information, it lacks the ability to capture fine-grained local details in specific regions of the image [5, 53, 63], *i.e.*, locality awareness. To mitigate this, a Locality Adapter (LA) enhances locality awareness of CLIP by updating

each patch token with aggregated information from neighboring tokens. Specifically, we first reshape the flattened patch tokens  $F$  to their original shape, and then project them through a Feed-Forward Network (FFN) to obtain the  $\tilde{F} \in \mathbb{R}^{H \times W \times D_a}$ , where  $D_a \ll D_{\text{clip}}$ . Then, we construct the spatial layout embedding  $L_{i,j} = \text{FFN}([b_t; c_t; e_t])$  according to the detection results, where  $t$  denotes the index of detected object corresponding to the position  $(i, j)$ , and  $[\cdot; \cdot]$  indicates the concatenation operation. Here,  $b_t$ ,  $c_t$ , and  $e_t$  represent the box coordinates, confidence score, and object text embedding extracted from the CLIP text encoder, respectively. Subsequently, the layout embedding  $L \in \mathbb{R}^{H \times W \times D_a}$  is embedded into  $\tilde{F}$  to provide spatial prior of entire objects in  $I$  as follows:

$$\hat{F} = \text{LN}(\text{FFN}(\tilde{F} + L)), \quad (4)$$

where LN denotes Layer Normalization [25]. From the  $\hat{F}$ , the LA aggregates neighborhood information of each patch token for locality awareness. We utilize multiple convolutional layers  $\{\text{Conv}^{k_n}\}_{n=1}^{N_c}$  with different kernel size  $k_n \times k_n$ , where  $k_n \in \mathbb{K} = \{k_1, k_2, \dots, k_{N_c}\}$  to aggregate the neighborhood information. The locality-aware feature  $P$  is extracted as:

$$L^{k_n} = \text{Conv}^{k_n}(\hat{F}), \quad (5)$$

$$P = \text{FFN}(L^{k_1} + \dots + L^{k_{N_c}}). \quad (6)$$

Then,  $P$  is projected back to  $D_{\text{clip}}$  and fused with the original  $F$  as follows:

$$F' = F + \gamma_{\text{LA}} \cdot \text{FFN}(P), \quad (7)$$

where  $\gamma_{\text{LA}} \in \mathbb{R}^{D_{\text{clip}}}$  is a learnable parameter to balance between  $P$  and  $F$ .

### 3.4. Interaction Adapter

Although the locality-aware feature helps the model capture fine-grained details of individual objects, it is insufficient to determine interactions, as interactions depend on the specific patterns between human and object contexts—specifically, how human cues are associated with object cues. For example, the 'riding a bike' interaction is identified by recognizing the association between the human cues and the object cues such as the hands in contact with the handle. This association distinguishes the interaction from other possible interactions, such as 'repairing a bike.' To enhance such interaction awareness, our Interaction Adapter (IA) updates each HO token based on its interaction pattern. We first extract region features for the human and object using ROAlign [15]. These features are then refined by the Interaction Pattern Reasoning Module (IPRM) by capturing the interaction pattern. The refined region features are subsequently used to inject interaction awareness into the corresponding HO token.



Specifically, the region features for the human and object in the  $i$ -th human-object pair are extracted as:

$$R_i^\tau = \text{FFN}(\text{ROIAlign}(F', b_i^\tau)), \quad (8)$$

where  $\tau \in \{h, o\}$  is an indicator for human/object, and  $b_i^\tau$  represents the corresponding bounding box of the  $i$ -th human-object pair. Then, the IPRM captures human and object region contexts using learnable queries  $Q \in \mathbb{R}^{N_p \times D_a}$  through a cross-attention mechanism:

$$\tilde{R}_i^\tau = \text{CrossAttn}(Q, R_i^\tau, R_i^\tau). \quad (9)$$

We utilize the  $N_p$  queries to capture the interaction-relevant contexts while filtering out irrelevant details. Next, IPRM reasons about interaction patterns by computing how each region context is associated with its counterpart through cross-attention:

$$\hat{R}_i^h = \text{CrossAttn}(\tilde{R}_i^h, \tilde{R}_i^o, \tilde{R}_i^o), \quad (10)$$

$$\hat{R}_i^o = \text{CrossAttn}(\tilde{R}_i^o, \tilde{R}_i^h, \tilde{R}_i^h). \quad (11)$$

Subsequently, the HO token  $T_i$  is projected into  $D_a$  dimension through an FFN:  $\tilde{T}_i = \text{FFN}(T_i)$ .  $\tilde{T}_i$  is used as a query to extract interaction-aware features  $\bar{R}_i^\tau$ , which are utilized to update the HO token:

$$\bar{R}_i^\tau = \text{CrossAttn}(\tilde{T}_i, \hat{R}_i^\tau, \hat{R}_i^\tau), \quad (12)$$

$$T'_i = T_i + \gamma_{\text{IA}} \cdot \text{FFN}([\bar{R}_i^h; \bar{R}_i^o]), \quad (13)$$

where  $\gamma_{\text{IA}}$  is a learnable parameter. The updated HO tokens then are passed through the  $l$ -th layer of the CLIP visual encoder, replacing Eq. 3:

$$[T_{(l)}, \text{cls}_{(l)}, F_{(l)}] = \mathcal{V}_{(l)}([T'_{(l-1)}; \text{cls}_{(l-1)}; F'_{(l-1)}]). \quad (14)$$

### 3.5. Training and Inference

Similar to the previous work [36, 49], we convert each HOI category into a text description using the template: ‘‘A photo of a person [verb-ing] a [object].’’ We then insert several learnable tokens in front of the text description. The text descriptions are fed into the CLIP text encoder to obtain the text embeddings  $E \in \mathbb{R}^{N_{|\mathcal{C}|} \times D_{\text{clip}}}$  for all HOI categories, where  $N_{|\mathcal{C}|}$  denotes the number of HOI categories. After obtaining the text embeddings, the HOI scores  $S \in \mathbb{R}^{N_{\text{pair}} \times N_{|\mathcal{C}|}}$  can be calculated as:

$$S = \text{Sigmoid}(T_{(L)} E^\top / \tau), \quad (15)$$

where  $\tau$  is the learnable parameter for rescaling the logits. Since a human can engage in multiple interactions with an object, we utilize the sigmoid function instead of softmax to compute the HOI scores.

**Training.** To train our proposed method, we assign positive labels to samples whose human and object bounding

boxes both have an Intersection-over-Union (IoU) exceeding a threshold with the ground truth. Following the previous work [23, 36], we adopt the binary focal loss [32]:

$$\mathcal{L} = \text{FocalBCE}(S, Y), \quad (16)$$

where  $Y \in \{0, 1\}^{N_{\text{pair}} \times N_{|\mathcal{C}|}}$  represents the binary target labels.

**Inference.** During inference, we incorporate the confidence scores of the human and object boxes from DETR into the HOI scores as:

$$S_{\text{infer}} = S \cdot S_H^\lambda \cdot S_O^\lambda, \quad (17)$$

where  $\lambda$  is a hyper-parameter for suppressing overconfident detections [58, 59].  $S_H, S_O \in [0, 1]^{N_{\text{pair}} \times 1}$  denotes the confidence scores of human and object for corresponding human-object pairs, respectively.

## 4. Experiments

### 4.1. Experiment Settings

To show the effectiveness of the proposed method, we evaluate our model on the two public benchmark datasets: HICO-DETR [2] and V-COCO [12].

**HICO-DETR** has 38,118 images for training and 9,658 images for testing. It contains 80 object classes, 117 interaction classes, and 600 HOI categories. Following conventional evaluation protocol [23, 36, 37], we report the mean average precision (mAP) to examine the model performance on five zero-shot settings: Unseen Combination(UC), Rare First Unseen Combination (RF-UC), Non-rare First Unseen Combination (NF-UC), Unseen Verb (UV), Unseen Object (UO). In the UC setting, all object and verb categories appear during the training; however, some HOI categories do not appear during the training, and they are used as  $\mathbb{C}_{\text{unseen}}^{\text{UC}}$ . Especially the least frequent 120 HOI categories are used as  $\mathbb{C}_{\text{unseen}}^{\text{RF-UC}}$  while the most frequent 120 HOI categories are used as  $\mathbb{C}_{\text{unseen}}^{\text{NF-UC}}$ . In the UV setting, 20 verb categories ( $\mathbb{V}_{\text{unseen}}$ ) are not used during the training, and corresponding HOI categories are used as unseen HOI categories, i.e.,  $\mathbb{C}_{\text{unseen}}^{\text{UV}} = \{(o_i, v_j) | o_i \in \mathbb{O}, v_j \in \mathbb{V}_{\text{unseen}}\}$ . Similarly, in the UO setting, 12 object categories ( $\mathbb{O}_{\text{unseen}}$ ) are not used during the training, and corresponding HOI categories are used as unseen HOI categories, i.e.,  $\mathbb{C}_{\text{unseen}}^{\text{UO}} = \{(o_i, v_j) | o_i \in \mathbb{O}_{\text{unseen}}, v_j \in \mathbb{V}\}$ .

**V-COCO** is a subset of the MS-COCO [31] dataset. It consists of 5,400 and 4,946 images for training and testing. V-COCO consists of 80 object classes and 29 action classes. Following evaluation settings in [21], we evaluate LAIN on scenario 2, and report role average precision  $\text{AP}_{\text{role}}^{\#2}$ .

Method	RF-UC			NF-UC			UO			UV			UC		
	Unseen	Seen	Full	Unseen	Seen	Full	Unseen	Seen	Full	Unseen	Seen	Full	Unseen	Seen	Full
FCL [19]	13.16	24.23	22.01	18.66	19.55	19.37	15.54	20.74	19.87	-	-	-	-	-	-
ATL [18]	9.18	24.67	21.57	18.25	18.78	18.67	15.11	21.54	20.47	-	-	-	-	-	-
RLIP [56]	19.19	33.35	30.52	20.27	27.67	26.19	-	-	-	-	-	-	-	-	-
GEN-VLKT [30]	21.36	32.91	30.56	25.05	23.38	23.71	10.51	28.92	25.63	20.96	30.23	28.74	-	-	-
LOGICHOI [26]	25.97	34.93	33.17	26.84	27.86	27.95	15.67	30.42	28.23	-	-	-	-	-	-
ADA-CM [23]	27.63	34.35	33.01	32.41	31.13	31.39	-	-	-	-	-	-	-	-	-
EoID [52]	22.04	31.39	29.52	26.77	26.66	26.69	-	-	-	22.71	30.73	29.61	23.01	30.39	28.91
HOICLIP [37]	25.53	34.85	32.99	26.39	28.10	27.75	16.20	30.99	28.53	24.30	32.19	31.09	23.15	31.65	29.93
CLIP [39]	28.79	22.00	23.36	28.52	22.06	23.36	28.66	22.29	23.36	26.16	22.90	23.36	24.28	23.12	23.36
CLIP4HOI [36]	28.47	35.48	34.08	31.44	28.26	28.90	31.79	32.73	32.58	26.02	31.14	30.42	27.71	33.25	32.11
BCOM <sup>†</sup> [46]	28.52	35.04	33.74	33.12	31.76	32.03	-	-	-	-	-	-	-	-	-
CMMP [24]	29.45	32.87	32.18	32.09	29.71	30.18	33.76	31.15	31.59	26.23	32.75	31.84	29.60	32.39	31.84
<b>LAIN</b>	<b>31.83</b>	<b>35.06</b>	<b>34.41</b>	<b>36.41</b>	<b>32.44</b>	<b>33.23</b>	<b>37.88</b>	<b>33.55</b>	<b>34.27</b>	<b>28.96</b>	<b>33.80</b>	<b>33.12</b>	<b>31.64</b>	<b>35.04</b>	<b>34.36</b>
<b>LAIN<sup>†</sup></b>	<b>36.57</b>	<b>38.54</b>	<b>38.13</b>	<b>37.52</b>	<b>35.90</b>	<b>36.22</b>	<b>40.78</b>	<b>36.96</b>	<b>37.60</b>	<b>32.05</b>	<b>38.04</b>	<b>37.20</b>	<b>32.25</b>	<b>37.95</b>	<b>36.81</b>

Table 1. Performance comparison on the HICO-DET dataset under various zero-shot settings. RF-UC, NF-UC, UO, UV, and UC denote rare first unseen composition, non-rare first unseen composition, unseen object, unseen verb, and unseen composition settings, respectively. Our method outperforms existing methods, demonstrating the effectiveness of the proposed methods. The highest result in each section is highlighted in bold. <sup>†</sup> indicates CLIP with ViT-L backbone.

## 4.2. Comparison with State-of-the-Art Methods

**Zero-shot settings.** We evaluate the performance of LAIN and compare it with existing HOI detection methods under various zero-shot settings. As shown in Table 1, LAIN demonstrates effectiveness by outperforming all previous methods by a significant margin under all zero-shot settings. In particular, existing methods [23, 24, 36, 37, 46, 52] that leverage CLIP representation show lower or comparable performance on unseen classes than CLIP itself under the RF-UC and UV settings. These results indicate that adapting CLIP representation for zero-shot HOI detection weakens its generalization ability due to the domain gap between the image-level pre-training task and HOI detection. In contrast, our model consistently outperforms CLIP and other existing methods on unseen classes, demonstrating its effectiveness. Furthermore, when we increase the model size to ViT-L, *i.e.*, LAIN<sup>†</sup>, the performance further improves, demonstrating the scalability of the proposed method with a larger backbone. Notably, despite BCOM<sup>†</sup> [46] using the larger ViT-L backbone, LAIN still surpasses it by a significant margin using only ViT-B. These results indicate that it is crucial for CLIP representation to consider fine-grained details of individual objects and interaction patterns between humans and objects.

**Fully-supervised settings.** To further validate the effectiveness of our proposed method, we conducted experiments under conventional fully-supervised settings on the HICO-DET and V-COCO datasets. As shown in Table 2, on the HICO-DET dataset, LAIN not only surpasses both fully-supervised models and zero-shot methods but also shows a marked improvement on rare HOI classes, which present significant challenges due to their scarcity and difficulty in generalizing, similar to unseen classes. Despite having fewer parameters and FLOPS, as shown in Ta-

Method	HICO-DET			V-COCO
	Full	Rare	Non-rare	AP <sup>#2</sup> <sub>role</sub>
<i>Fully-supervised methods</i>				
HOTR [21]	25.10	17.34	27.42	64.4
ATL [18]	28.53	21.63	30.59	-
As-Net [3]	28.87	24.25	30.25	-
QPIC [42]	29.07	21.85	31.23	61.0
UPT [59]	31.66	25.94	33.36	64.5
CDN [57]	31.78	27.55	33.05	64.4
Iwin [43]	32.03	27.62	34.14	60.5
GEN-VLKT [30]	33.75	29.25	35.10	64.5
ADA-CM [23]	33.80	31.72	34.42	61.2
LogicHOI [26]	35.47	32.03	<b>36.22</b>	<u>65.6</u>
<i>Zero-shot methods</i>				
HOICLIP [37]	34.69	31.12	35.74	64.8
CLIP4HOI [36]	35.33	33.95	35.74	<b>66.3</b>
CMMP [24]	32.26	33.53	33.24	61.2
LAIN	<b>36.02</b>	<b>35.70</b>	<u>36.11</u>	65.1

Table 2. Performance comparison on the HICO-DET [2] and V-COCO [12] datasets under fully-supervised setting. The highest result in each section is highlighted in bold.

ble 7, LAIN demonstrates competitive performance on the V-COCO dataset, achieving the second-best results among zero-shot methods.

## 4.3. Ablation Study

We conduct various ablation studies on the UV setting to validate the effectiveness of LAIN.

**The impact of each adapter.** In Table 3, we gradually add each adapter to the baseline, which utilizes the original CLIP representation without LA and IA, to investigate the impact of each adapter. We observed performance improvements in both seen and unseen classes by injecting locality awareness into CLIP representations through the LA. This result indicates that capturing the fine-grained information of individual objects is crucial for zero-shot HOI

LA	IA	Unseen	Seen	Full
-	-	24.88	31.06	30.19
✓	-	27.71	32.55	31.95
-	✓	27.37	33.57	32.70
✓	✓	<b>30.50</b>	<b>34.80</b>	<b>33.95</b>

Table 3. Ablations studies on each adapter under UV setting. LA and IA denote locality and interaction adapter, respectively.

detection. Furthermore, applying IA to enhance interaction awareness in CLIP representations leads to performance improvements, emphasizing the importance of understanding interaction patterns between humans and objects. Moreover, since fine-grained information about individual objects aids in reasoning about interaction patterns, we observed that using both adapters together yielded the best performance. This demonstrates the effectiveness of both adapters, *i.e.*, LA and IA, and highlights the importance of incorporating locality and interaction awareness for adapting CLIP representations to zero-shot HOI detection.

**The impact of each component in LA.** To investigate the impact of each component in the LA, we conduct comparisons between the full LA model and various LA variants in Table 4 (a) to (d). Removing either the visual context or the spatial layout component, *i.e.*, (a) or (b), results in a significant decline in performance. This demonstrates that integrating both surrounding visual context and spatial layout is essential for capturing fine-grained details of individual objects. Furthermore, we replace our convolutional layers with existing attention mechanisms [34, 40], which are designed to capture local information, *i.e.*, fine-grained details, instead of global information, as shown in (c) and (d). We observe that the existing mechanisms improve performance on both unseen and seen classes compared to baseline, *i.e.*, without LA and IA, by capturing local information. These results demonstrate the importance of fine-grained details about individual objects in adapting CLIP’s representations for zero-shot HOI detection. However, their performances are degraded compared to our LA. This highlights the effectiveness of LA in capturing fine-grained details.

**The impact of each component in IA.** We also investigate the impact of each component in IA on achieving interaction awareness. In Table 4, we compare the full IA model with various IA variants. When we remove the IPRM, the model’s performance on unseen classes significantly degrades, indicating that incorporating the interaction pattern is crucial for interaction awareness as shown in (e). Similarly, removing human/object context extraction (*i.e.*, extracting the interaction pattern using only ROI-aligned features) results in decreased model performance, suggesting that capturing interaction-relevant contexts while filtering out irrelevant details through context extraction is effective for reasoning about the interaction pattern as shown in (f).

	Method	Unseen	Seen	Full
	Baseline	24.88	31.06	30.19
	<b>Locality Adapter</b>	<b>27.71</b>	<b>32.55</b>	<b>31.95</b>
(a)	w.o visual information	26.77	32.18	31.40
(b)	w.o spatial layout	26.52	32.07	31.31
(c)	Local Attention [40]	26.46	32.39	31.56
(d)	Window Attention [34]	26.35	32.31	31.47
	<b>Interaction Adapter</b>	<b>27.37</b>	<b>33.57</b>	<b>32.70</b>
(e)	w.o IPRM	24.32	32.76	31.57
(f)	w.o human/object context	25.64	32.41	31.40

Table 4. Ablation study on the design choices of each adapter under UV setting.

	Method	Unseen			Seen		
		AP <sub>L</sub>	AP <sub>M</sub>	AP <sub>S</sub>	AP <sub>L</sub>	AP <sub>M</sub>	AP <sub>S</sub>
Human	CMMP [24]	45.77	34.54	52.89	64.80	54.60	25.05
	ADA-CM [23]	<b>50.63</b>	36.85	26.88	60.80	39.13	17.57
	LAIN	<u>48.11</u>	<b>38.66</b>	<b>62.32</b>	<b>66.51</b>	<b>59.92</b>	<b>30.44</b>
Object	CMMP [24]	59.64	46.19	18.64	62.73	51.38	29.82
	ADA-CM [23]	<b>62.16</b>	46.60	14.02	57.20	37.99	18.66
	LAIN	<u>61.41</u>	<b>49.45</b>	<b>20.82</b>	<b>63.82</b>	<b>56.93</b>	<b>35.91</b>

Table 5. Comparison of AP across different box sizes under the RF-UC setting.<sup>1</sup>

position	Unseen	Seen	Full
∅	24.88	31.06	30.19
1-6	27.29	32.11	31.46
7-12	27.83	32.72	32.04
1-12	<b>28.96</b>	<b>33.80</b>	<b>33.12</b>

Table 6. Ablation study evaluating the impact of different adapter positions in the UV setting.

**The impact of locality-aware adaptation.** To validate that our IA and LA modules enhance CLIP’s ability to focus on local details, we conduct experiments varying human and object box sizes, as local details become more critical with smaller boxes. As shown in Table 5, as the box size decreases from large (AP<sub>L</sub>) to small (AP<sub>S</sub>), LAIN demonstrates a larger performance gap for medium and small boxes compared to large ones. Specifically, LAIN outperforms ADA-CM [23] by 4.91% (human) and 6.12% (object) in AP<sub>M</sub>, and notably by 131.85% (human) and 48.50% (object) in AP<sub>S</sub> under the unseen setting. The performance gain becomes increasingly prominent as box sizes decrease (M→S), and similar results are also observed in the seen classes. These results validate the effectiveness of our IA and LA in capturing fine-grained local details for zero-shot HOI detection.

**The impact of adapter positioning.** In Table 6, we conduct experiments that explore the impact of adapter positioning. We divide the model layers into lower and upper halves, inserting adapters into each section separately. We observe that adding adapters to either the lower layers

<sup>1</sup>Since the authors of ADA-CM [23] provide pretrained weight under UC-RF instead of UV, we conduct our experiments under the UC-RF.



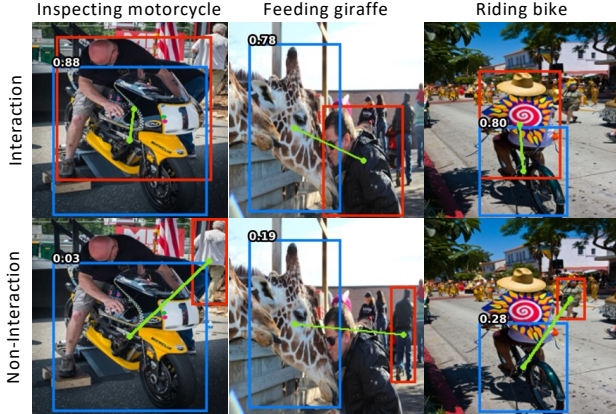


Figure 3. Qualitative results on HICO-DET under UV settings. We represent a human with a red box and an object with a blue box, along with HOI score.

(*i.e.*, 1-6) or the upper layers (*i.e.*, 7-12) consistently improves the model’s performance. In particular, we observe that incorporating adapters into the upper layers results in a more substantial performance gain since the upper layers in the ViT structure encode global information, unlike the lower layer which encodes local information [6]. These results suggest that injecting interaction and locality awareness into features lacking local information is essential and that our proposed method effectively incorporates interaction and locality awareness into CLIP representations.

**Parameter analysis.** In Table 7, we provide a comparative analysis of model parameters and computational cost between our proposed method and state-of-the-art methods for zero-shot HOI detection to show the efficiency of LAIN. Compared to HOICLIP [37] and CLIP4HOI [36], which introduce a large number of trainable parameters and FLOPs as they utilize a heavy decoder for interaction classification, LAIN introduces only 3.0M trainable parameters, alongside reduced FLOPs of 110G. Although LAIN has more trainable parameters than CMMP [24], it has fewer total parameters and FLOPs. This highlights the efficiency of our approach in incorporating locality and interaction awareness with minimal added overhead.

#### 4.4. Qualitative Results

We present qualitative results on the HICO-DET under the UV setting in Figure 3 and 4. In Figure 3, we observe that LAIN successfully distinguishes interactive pairs, assigning high similarity scores, and non-interactive pairs, assigning low similarity scores, even for unseen verbs, *i.e.*, ‘inspecting’, ‘feeding’, and ‘riding’. Moreover, as shown in Figure 4, LAIN assigns significantly lower scores to non-interactive pairs compared to the baseline, which does not incorporate LA and IA. The results demonstrate that incorporating LA and IA leads to more discriminative HOI representations, enabling the model to better distinguish interactive and non-interactive pairs for zero-shot HOI detection.

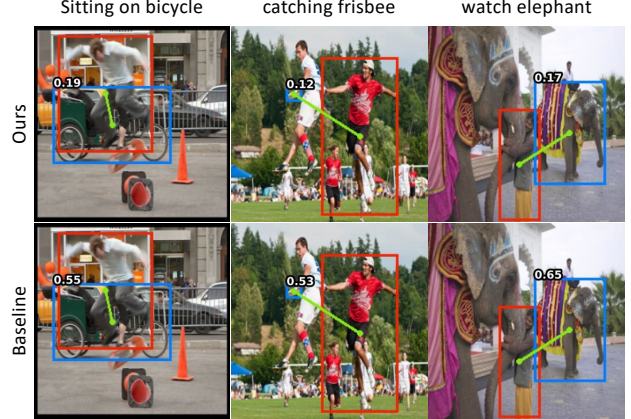


Figure 4. Qualitative comparison of non-interactive pairs between LAIN and the baseline, *i.e.*, without LA and IA, on the HICO-DET under the UV setting. We represent a human with a red box and an object with a blue box, along with HOI scores.

Method	Tr. Params	Tot. Params	FLOPs
HOICLIP [37]	193.3M	193.3M	179G
CLIP4HOI [36]	71.2M	262.4M	186G
CMMP [24]	<b>2.3M</b>	193.4M	114G
LAIN	3.0M	<b>145.4M</b>	<b>110G</b>

Table 7. Comparison of parameters across state-of-the-art models under UV setting. Tr. and Tot. Params represent the number of trainable and total parameters of the model, respectively.

## 5. Conclusion

In this paper, we have proposed LAIN, designed to address the lack of local details in CLIP’s representation, which hinders CLIP’s generalization ability when adapting to zero-shot HOI detection. Our locality adapter introduces locality awareness into CLIP by considering surrounding visual information and spatial layout. For interaction awareness, which is difficult to determine solely through locality awareness, our interaction adapter infers the interaction pattern by leveraging contextual reasoning between human and object contexts. By enhancing the locality and interaction awareness, LAIN effectively captures fine-grained information about HO pairs, facilitating adaptation of CLIP’s representation to zero-shot HOI detection. Extensive experiments on two public benchmarks, HICO-DET and V-COCO, demonstrate the importance of capturing local details and the effectiveness of LAIN. Notably, LAIN is particularly effective in scenarios where local details are crucial, *i.e.*, small instances, while introducing minimal computational cost.

**Acknowledgements.** This work was supported by the NRF grant (RS-2021-NR059830 (50%)) and the IITP grants (RS-2022-II220959: Few-Shot Learning of Causal Inference in Vision and Language for Decision Making (45%), RS-2019-II191906: AI Graduate School Program at POSTECH (5%)) funded by the Korea government (MSIT).



## References

- [1] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020. 2, 3, 4
- [2] Yu-Wei Chao, Yunfan Liu, Xieyang Liu, Huayi Zeng, and Jia Deng. Learning to detect human-object interactions. In *2018 IEEE winter conference on applications of computer vision (WACV)*, pages 381–389. IEEE, 2018. 2, 5, 6
- [3] Mingfei Chen, Yue Liao, Si Liu, Zhiyuan Chen, Fei Wang, and Chen Qian. Reformulating hoi detection as adaptive set prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9004–9013, 2021. 6
- [4] Zhe Chen, Yuchen Duan, Wenhai Wang, Junjun He, Tong Lu, Jifeng Dai, and Yu Qiao. Vision transformer adapter for dense predictions. *arXiv preprint arXiv:2205.08534*, 2022. 2
- [5] Xiaoyi Dong, Jianmin Bao, Yinglin Zheng, Ting Zhang, Dongdong Chen, Hao Yang, Ming Zeng, Weiming Zhang, Lu Yuan, Dong Chen, et al. Maskclip: Masked self-distillation advances contrastive language-image pretraining. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10995–11005, 2023. 2, 4
- [6] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 8
- [7] Sepideh Esmailpour, Bing Liu, Eric Robertson, and Lei Shu. Zero-shot out-of-distribution detection based on the pre-trained model clip. In *Proceedings of the AAAI conference on artificial intelligence*, pages 6568–6576, 2022. 2
- [8] Hao-Shu Fang, Yichen Xie, Dian Shao, and Cewu Lu. Dirv: Dense interaction region voting for end-to-end human-object interaction detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 1291–1299, 2021. 2
- [9] Chen Gao, Yuliang Zou, and Jia-Bin Huang. ican: Instance-centric attention network for human-object interaction detection. *arXiv preprint arXiv:1808.10437*, 2018. 2
- [10] Chen Gao, Jiarui Xu, Yuliang Zou, and Jia-Bin Huang. Drg: Dual relation graph for human-object interaction detection. In *European Conference on Computer Vision*, pages 696–712. Springer, 2020. 2
- [11] Albert Gordo and Diane Larlus. Beyond instance-level image retrieval: Leveraging captions to learn a global visual representation for semantic retrieval. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6589–6598, 2017. 1
- [12] Saurabh Gupta and Jitendra Malik. Visual semantic role labeling. *arXiv preprint arXiv:1505.04474*, 2015. 2, 5, 6
- [13] Tanmay Gupta, Alexander Schwing, and Derek Hoiem. No-frills human-object interaction detection: Factorization, layout encodings, and training techniques. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9677–9685, 2019. 2
- [14] Ali Hassani, Steven Walton, Jiachen Li, Shen Li, and Humphrey Shi. Neighborhood attention transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6185–6194, 2023. 2
- [15] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 4
- [16] Simao Herdade, Armin Kappeler, Kofi Boakye, and Joao Soares. Image captioning: Transforming objects into words. *Advances in Neural Information Processing Systems*, 32, 2019. 1
- [17] Zhi Hou, Xiaojiang Peng, Yu Qiao, and Dacheng Tao. Visual compositional learning for human-object interaction detection. In *European Conference on Computer Vision*, pages 584–600. Springer, 2020. 2
- [18] Zhi Hou, Baosheng Yu, Yu Qiao, Xiaojiang Peng, and Dacheng Tao. Affordance transfer learning for human-object interaction detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 495–504, 2021. 6
- [19] Zhi Hou, Baosheng Yu, Yu Qiao, Xiaojiang Peng, and Dacheng Tao. Detecting human-object interaction via fabricated compositional learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14646–14655, 2021. 2, 6
- [20] Bumsoo Kim, Taeho Choi, Jaewoo Kang, and Hyunwoo J Kim. Uniondet: Union-level detector towards real-time human-object interaction detection. In *European Conference on Computer Vision*, pages 498–514. Springer, 2020. 2
- [21] 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*, pages 74–83, 2021. 5, 6
- [22] Sanghyun Kim, Deunsol Jung, and Minsu Cho. Relational context learning for human-object interaction detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2925–2934, 2023. 1, 2
- [23] Ting Lei, Fabian Caba, Qingchao Chen, Hailin Jin, Yuxin Peng, and Yang Liu. Efficient adaptive human-object interaction detection with concept-guided memory. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6480–6490, 2023. 1, 3, 5, 6, 7
- [24] Ting Lei, Shaofeng Yin, Yuxin Peng, and Yang Liu. Exploring conditional multi-modal prompts for zero-shot hoi detection. In *European Conference on Computer Vision*, pages 1–19. Springer, 2025. 3, 6, 7, 8
- [25] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *ArXiv e-prints*, pages arXiv–1607, 2016. 4
- [26] Liulei Li, Jianan Wei, Wenguan Wang, and Yi Yang. Neural logic human-object interaction detection. *Advances in Neural Information Processing Systems*, 36, 2024. 1, 2, 6

- [27] 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*, pages 3585–3594, 2019. [2](#)
- [28] Yong-Lu Li, Xinpeng Liu, Xiaoqian Wu, Yizhuo Li, and Cewu Lu. Hoi analysis: Integrating and decomposing human-object interaction. *Advances in Neural Information Processing Systems*, 33:5011–5022, 2020. [2](#)
- [29] Yue Liao, Si Liu, Fei Wang, Yanjie Chen, Chen Qian, and Jia-shi Feng. Ppdm: Parallel point detection and matching for real-time human-object interaction detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 482–490, 2020. [2](#)
- [30] Yue Liao, Aixi Zhang, Miao Lu, Yongliang Wang, Xiaobo Li, and Si Liu. Gen-vlkt: Simplify association and enhance interaction understanding for hoi detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20123–20132, 2022. [3](#), [6](#)
- [31] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. [5](#)
- [32] 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*, pages 2980–2988, 2017. [5](#)
- [33] Ye Liu, Junsong Yuan, and Chang Wen Chen. Consnet: Learning consistency graph for zero-shot human-object interaction detection. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 4235–4243, 2020. [2](#)
- [34] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021. [2](#), [7](#)
- [35] Yue Ma, Yali Wang, Yue Wu, Ziyu Lyu, Siran Chen, Xiu Li, and Yu Qiao. Visual knowledge graph for human action reasoning in videos. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 4132–4141, 2022. [2](#)
- [36] Yunyao Mao, Jiajun Deng, Wengang Zhou, Li Li, Yao Fang, and Houqiang Li. Clip4hoi: Towards adapting clip for practical zero-shot hoi detection. *Advances in Neural Information Processing Systems*, 36, 2024. [1](#), [3](#), [4](#), [5](#), [6](#), [8](#)
- [37] Shan Ning, Longtian Qiu, Yongfei Liu, and Xuming He. Hoiclip: Efficient knowledge transfer for hoi detection with vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23507–23517, 2023. [1](#), [3](#), [5](#), [6](#), [8](#)
- [38] Siyuan Qi, Wenguan Wang, Baoxiong Jia, Jianbing Shen, and Song-Chun Zhu. Learning human-object interactions by graph parsing neural networks. In *Proceedings of the European conference on computer vision (ECCV)*, pages 401–417, 2018. [2](#)
- [39] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. [2](#), [3](#), [6](#)
- [40] Prajit Ramachandran, Niki Parmar, Ashish Vaswani, Irwan Bello, Anselm Levskaya, and Jon Shlens. Stand-alone self-attention in vision models. *Advances in neural information processing systems*, 32, 2019. [7](#)
- [41] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28, 2015. [2](#)
- [42] Masato Tamura, Hiroki Ohashi, and Tomoaki Yoshinaga. Qpic: Query-based pairwise human-object interaction detection with image-wide contextual information. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10410–10419, 2021. [2](#), [6](#)
- [43] Danyang Tu, Xiongkuo Min, Huiyu Duan, Guodong Guo, Guangtao Zhai, and Wei Shen. Iwin: Human-object interaction detection via transformer with irregular windows. In *European Conference on Computer Vision*, pages 87–103. Springer, 2022. [6](#)
- [44] Oytun Ulutan, ASM Iftekhar, and Bangalore S Manjunath. Vsgnet: Spatial attention network for detecting human object interactions using graph convolutions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13617–13626, 2020. [2](#)
- [45] Bo Wan and Tinne Tuytelaars. Exploiting clip for zero-shot hoi detection requires knowledge distillation at multiple levels. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1805–1815, 2024. [3](#)
- [46] Guangzhi Wang, Yangyang Guo, Ziwei Xu, and Mohan Kankanhalli. Bilateral adaptation for human-object interaction detection with occlusion-robustness. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 27970–27980, 2024. [3](#), [6](#)
- [47] Hai Wang, Wei-shi Zheng, and Ling Yingbiao. Contextual heterogeneous graph network for human-object interaction detection. In *European Conference on Computer Vision*, pages 248–264. Springer, 2020. [2](#)
- [48] Hualiang Wang, Yi Li, Huifeng Yao, and Xiaomeng Li. Clipn for zero-shot ood detection: Teaching clip to say no. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1802–1812, 2023. [2](#)
- [49] Suchen Wang, Yueqi Duan, Henghui Ding, Yap-Peng Tan, Kim-Hui Yap, and Junsong Yuan. Learning transferable human-object interaction detector with natural language supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 939–948, 2022. [5](#)
- [50] Hui Wu, Min Wang, Wengang Zhou, Houqiang Li, and Qi Tian. Contextual similarity distillation for asymmetric image retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9489–9498, 2022. [1](#)

- [51] Mingrui Wu, Xuying Zhang, Xiaoshuai Sun, Yiyi Zhou, Chao Chen, Jiabin Gu, Xing Sun, and Rongrong Ji. Difnet: Boosting visual information flow for image captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18020–18029, 2022. 1
- [52] Mingrui Wu, Jiabin Gu, Yunhang Shen, Mingbao Lin, Chao Chen, and Xiaoshuai Sun. End-to-end zero-shot hoi detection via vision and language knowledge distillation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 2839–2846, 2023. 1, 3, 6
- [53] Sizhe Wu, Wenwei Zhang, Lumin Xu, Sheng Jin, Xiangtai Li, Wentao Liu, and Chen Change Loy. Clipsef: Vision transformer distills itself for open-vocabulary dense prediction. *arXiv preprint arXiv:2310.01403*, 2023. 2, 4
- [54] Ting Yao, Yingwei Pan, Yehao Li, and Tao Mei. Exploring visual relationship for image captioning. In *Proceedings of the European conference on computer vision (ECCV)*, pages 684–699, 2018. 1
- [55] Sangwoong Yoon, Woo Young Kang, Sungwook Jeon, SeongEun Lee, Changjin Han, Jonghun Park, and Eun-Sol Kim. Image-to-image retrieval by learning similarity between scene graphs. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 10718–10726, 2021. 1
- [56] Hangjie Yuan, Jianwen Jiang, Samuel Albanie, Tao Feng, Ziyuan Huang, Dong Ni, and Mingqian Tang. Rlip: Relational language-image pre-training for human-object interaction detection. *Advances in Neural Information Processing Systems*, 35:37416–37431, 2022. 6
- [57] Aixi Zhang, Yue Liao, Si Liu, Miao Lu, Yongliang Wang, Chen Gao, and Xiaobo Li. Mining the benefits of two-stage and one-stage hoi detection. *Advances in Neural Information Processing Systems*, 34:17209–17220, 2021. 6
- [58] 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*, pages 13319–13327, 2021. 2, 5
- [59] Frederic Z Zhang, Dylan Campbell, and Stephen Gould. Efficient two-stage detection of human-object interactions with a novel unary-pairwise transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20104–20112, 2022. 4, 5, 6
- [60] Yong Zhang, Yingwei Pan, Ting Yao, Rui Huang, Tao Mei, and Chang-Wen Chen. Exploring structure-aware transformer over interaction proposals for human-object interaction detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19548–19557, 2022. 2
- [61] Yiwu Zhong, Jianwei Yang, Pengchuan Zhang, Chunyuan Li, Noel Codella, Liunian Harold Li, Luwei Zhou, Xiyang Dai, Lu Yuan, Yin Li, et al. Regionclip: Region-based language-image pretraining. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16793–16803, 2022. 1, 2
- [62] Desen Zhou, Zhichao Liu, Jian Wang, Leshan Wang, Tao Hu, Errui Ding, and Jingdong Wang. Human-object interaction detection via disentangled transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19568–19577, 2022. 2
- [63] Ziqin Zhou, Yinjie Lei, Bowen Zhang, Lingqiao Liu, and Yifan Liu. Zegclip: Towards adapting clip for zero-shot semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11175–11185, 2023. 1, 2, 4