

Exploring Structure-aware Transformer over Interaction Proposals for Human-Object Interaction Detection*

Yong Zhang[†], Yingwei Pan[‡], Ting Yao[‡], Rui Huang[†], Tao Mei[‡], and Chang-Wen Chen[§]

[†] The Chinese University of Hong Kong, Shenzhen [‡] JD Explore Academy

[§] The Hong Kong Polytechnic University

yongzhang@link.cuhk.edu.cn, {panyw.ustc, tingyao.ustc}@gmail.com, ruihuang@cuhk.edu.cn, tmei@jd.com, changwen.chen@polyu.edu.hk

Abstract

Recent high-performing Human-Object Interaction (HOI) detection techniques have been highly influenced by Transformer-based object detector (i.e., DETR). Nevertheless, most of them directly map parametric interaction queries into a set of HOI predictions through vanilla Transformer in a one-stage manner. This leaves rich inter- or intra-interaction structure under-exploited. In this work, we design a novel Transformer-style HOI detector, i.e., Structure-aware Transformer over Interaction Proposals (STIP), for HOI detection. Such design decomposes the process of HOI set prediction into two subsequent phases, i.e., an interaction proposal generation is first performed, and then followed by transforming the non-parametric interaction proposals into HOI predictions via a structure-aware Transformer. The structure-aware Transformer upgrades vanilla Transformer by encoding additionally the holistically semantic structure among interaction proposals as well as the locally spatial structure of human/object within each interaction proposal, so as to strengthen HOI predictions. Extensive experiments conducted on V-COCO and HICO-DET benchmarks have demonstrated the effectiveness of STIP, and superior results are reported when comparing with the state-of-the-art HOI detectors. Source code is available at <https://github.com/zyong812/STIP>.

1. Introduction

Human-Object Interaction (HOI) detection [5, 11] is intended to localize the interactive human-object pairs within an image and identify the interactions in between, yielding the HOI predictions in the form of $\langle \text{human}, \text{object}, \text{interaction} \rangle$ triplets. Practical HOI detection systems perform the human-centric scene under-

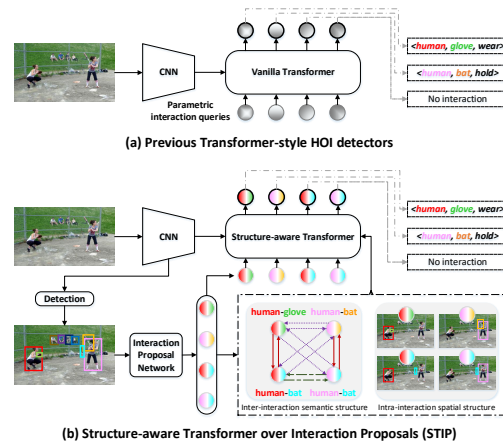


Figure 1. Comparison between existing Transformer-style HOI detectors and our STIP. (a) Existing Transformer-style HOI detectors directly transform the parametric interaction queries into HOI predictions via vanilla Transformer in a one-stage fashion. (b) STIP adopts a two-phase solution, i.e., first producing interaction proposals via Interaction Proposal Network, and then mapping the non-parametric interaction queries (i.e., interaction proposals) into HOI predictions. Both of the inter- and intra-interaction structures derived from interaction proposals are additionally exploited to boost HOI set prediction through a structure-aware Transformer.

standing, and thus have a great potential impact for numerous applications, such as surveillance event detection [1, 7] and robot imitation learning [2]. In general, conventional HOI detectors [8–10, 16, 24, 27, 32, 38–40] tackle the HOI set prediction task in an **indirect** way, by formalizing it as surrogate regression and classification problems for human/object/interaction. Such indirect approach needs a subsequent post-processing by collapsing near-duplicate predictions and heuristic matching [16, 24, 40], and thus cannot be trained in an end-to-end fashion, resulting in a sub-optimal solution. The intent to overcome the problem of sub-optimal solution leads to the development of recent state-of-the-art HOI detectors [6, 17, 36, 50] that follow the Transformer-based detector of DETR [3] to cast HOI detection as a **direct** set prediction problem, and embrace the “end-to-end” philosophy (Figure 1 (a)). In par-

*This work was performed at JD Explore Academy.

ticular, a vanilla Transformer is commonly utilized to map the parametric interaction queries (i.e., the learnable positional embedding sequence) into a set of HOI predictions in a one-stage manner. However, these HOI detectors start the HOI set prediction from the parametric interaction queries with randomly initialized embeddings. That is, the correspondence between parametric interaction queries and output HOIs (commonly assigned by Hungarian algorithm for training) is **dynamic** in which the interaction query corresponding to each target HOI (e.g., “human hold bat”) is unknown at the beginning of HOI set prediction. This can adversely hinder the exploration of prior knowledge (i.e., **inter-interaction** or **intra-interaction structure**) which would be very useful for relational reasoning among interactions in HOI set prediction.

Specifically, by inter-interaction structure, we refer to the holistic semantic dependencies among HOIs, which can be directly defined by considering whether or not two HOIs share the same human or object. Such structure implies the common sense knowledge that shall facilitate the prediction of one HOI by exploiting its semantic dependencies against other HOIs. Taking the input image in Figure 1 as an example, the existence of “human wear (baseball) glove” provides strong indication for “(another) human hold bat”. Moreover, the intra-interaction structure can be interpreted as the local spatial structure for each HOI, i.e., the layout of human and object, which acts as additional prior knowledge to steer model’s attention over image areas for depicting the interaction.

In this work, we design a novel scheme based on Transformer-style HOI detector, namely Structure-aware Transformer over Interaction Proposals (STIP). The design innovation is to decompose the one-stage solution of set prediction into two cascaded phases, i.e., first producing the interaction proposals (i.e., the probably interactive human-object pairs) and then performing HOI set prediction based on the interaction proposals (Figure 1 (b)). By taking the interaction proposals derived from Interaction Proposal Network (IPN) as non-parametric interaction queries, STIP naturally triggers the subsequent HOI set prediction with more reasonable interaction queries, leading to **static** query-HOI correspondence that capable of boosting HOI set prediction. As a beneficial by-product, the predicted interaction proposals offer a fertile ground for constructing a structured understanding across interaction proposals or between human & object within each interaction proposal. A particular form of Transformer, i.e., structure-aware Transformer, is designed accordingly to encode the inter-interaction or intra-interaction structure for enhancing HOI predictions.

In sum, we have made the following contributions: (1) The proposed two-phase implementation of Transformer-style HOI detector is shown capable of seamless incorporation of potential interactions among HOI proposals to

overcome the problem associated with one-stage approach; (2) The exquisitely designed structure-aware Transformer is shown able to facilitate additional exploitation opportunity for utilizing inter-interaction and intra-interaction structure for enhanced performance of the vanilla Transformer; (3) The proposed structure-aware Transformer approach has been properly analyzed and verified through extensive experiments over V-COCO and HICO-DET datasets to validate its potential in solving the problems associated with one-stage approach to achieve desirable HOI detection.

2. Related Work

The task of Human-Object Interaction (HOI) detection has been primordially defined [5, 11] and recent developments of HOI detectors can be briefly divided into two categories: the two-stage methods and one-stage approaches.

Two-stage Methods. The first category schemes [4, 8–10, 14, 15, 22, 23, 27, 28, 32, 37–39, 43, 47] mainly adopt two-stage paradigm, i.e., first detect humans/objects via off-the-shelf modern object detectors (e.g., Faster R-CNN [33]) and then carry out interaction classification. A number of schemes have been proposed to strengthen the HOI feature learning in the second stage for interaction classification. Generally, similar to prior works for visual relationship detection [18, 29, 42, 45, 46], HOI features are typically derived from three perspectives [4, 9, 10]: appearance/visual features of humans and objects, spatial features (e.g., the pairwise bounding boxes of human-object pair), and linguistic feature (e.g., the semantic embeddings of human/object labels). Various approaches [8, 13, 32, 37, 38, 44] further capitalize on message passing mechanism to perform relational reasoning over instance-centric graph structure, aiming to enrich HOI features with global contextual information among human and object instances. The authors in [39] devise contextual attention mechanism to facilitate the mining of contextual cues. Moreover, the information about human pose [12, 23, 47], body parts [49] or detailed 3D body shape [21] can also be exploited to enhance HOI feature representation. In [28, 41], additional knowledge from external source and language domain are further exploited to boost HOI feature learning. Most recently, the ATL scheme [14] constructs the affordance feature bank across multiple HOI datasets and injects affordance feature into object representations when inferring interactions.

One-stage Approaches. The second category schemes mainly construct one-stage HOI detectors [6, 16, 17, 24, 36, 40, 48, 50] by directly predicting HOI triplets, which are potentially faster and simpler than two-stage HOI detectors. UnionDet [16] is one of first attempts that directly detects the union regions of human-object pairs in a one-stage manner. Other schemes [24, 40] formulate HOI detection as a keypoint detection problem, and thus enable one-stage solution for this task. Most recently, inspired by the success of

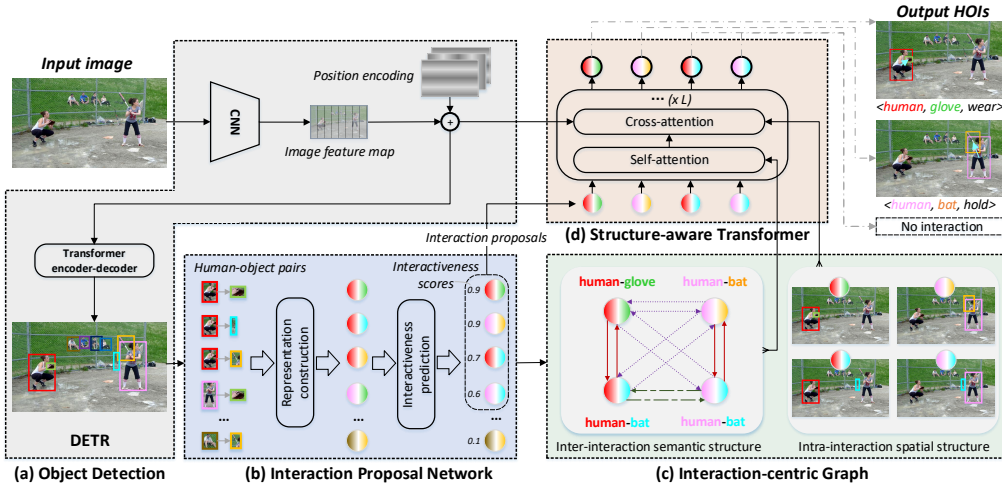


Figure 2. An overview of our proposed STIP framework. (a) Given an input image, we adopt an off-the-shelf DETR to detect the human and object instances within this image. (b) Based on the detected human and object instances, the Interaction Proposal Network (IPN) constructs all possible human-object pairs and then predicts the interactiveness score of each human-object pair. The most interactive human-object pairs with highest interactiveness scores are taken as the output interaction proposals. (c) Next, by taking all interaction proposals as graph nodes and exploiting semantic connectivity as edges, we build an interaction-centric graph that unfolds rich inter-interaction semantic structure and intra-interaction spatial structure. (d) Finally, a structure-aware Transformer is utilized to transform the non-parametric interaction queries (i.e., interaction proposals) into a set of HOI predictions by additionally guiding relational reasoning with the inter- or intra-interaction structure derived from interaction-centric graph.

Transformer-based object detectors (e.g., DETR [3]), there has been a steady momentum of breakthroughs that push the limits of HOI detection by using Transformer-style architecture. In particular, the authors in [36, 50] employ a single interaction Transformer decoder to predict a set of HOI triplets, and the whole architecture can be optimized in an end-to-end fashion with Hungarian loss. However, the authors in [6, 17] design two parallel Transformer decoders for detecting interactions and instances, and the outputs are further associated to produce final HOI predictions.

This Scheme. The proposed STIP can also be considered as Transformer-style architecture that tackles HOI detection as a set prediction problem, which eliminates the post-processing and enables the architecture to be end-to-end trainable. Unlike existing Transformer-style methods [6, 17, 36, 50] that perform HOI set prediction in a one-stage manner, the proposed STIP decomposes this process into two phases: the proposed scheme first produces interaction proposals as high-quality interaction queries and then takes them as non-parametric queries to trigger the HOI set prediction. Moreover, this STIP scheme goes beyond the conventional relational reasoning via vanilla Transformer by leveraging a structure-aware Transformer to exploit the rich inter- or intra-interaction structure, thereby boosting the performance of the HOI detection.

3. Approach

In this work, we devise the Structure-aware Transformer over Interaction Proposals (STIP) that casts HOI detection as a set prediction problem in a two-phase fashion. Meanwhile, this scheme boosts HOI set prediction with the prior

knowledge of inter- and intra-interaction structures. Figure 2 depicts an overview of the proposed STIP. The whole framework consists of four main components, i.e., DETR for object detection, interaction proposal network for producing interaction proposals, interaction-centric graph construction, and structure-aware Transformer for HOI set prediction. Specifically, an off-the-shelf DETR [3] is first adopted to detect humans and objects within the input image. Next, based on the detection results, we design the Interaction Proposal Network (IPN) to select the most interactive human-object pairs as interaction proposals. After that, we take all selected interaction proposals as graph nodes to construct an interaction-centric graph to reveal the inter-interaction semantic structure and intra-interaction spatial structure. The selected interaction proposals are further taken as non-parametric queries to trigger the HOI set prediction via a structure-aware Transformer through exploiting the structured prior knowledge derived from interaction-centric graph to strengthen relational reasoning.

3.1. Interaction Proposal Network

Conditioned on the detected human and object instances from DETR, the Interaction Proposal Network (IPN) targets for producing interaction proposals, i.e., the probably interactive human-object pairs. Concretely, we first construct all possible human-object pairs with pairwise connectivity between detected humans and objects. For each human-object pair, the IPN further predicts the probability of interaction existing in between (i.e., “interactiveness” score) through a multi-layer perceptron (MLP). Only the top- K human-object pairs with highest interactiveness scores are finally emitted as the output interaction proposals.

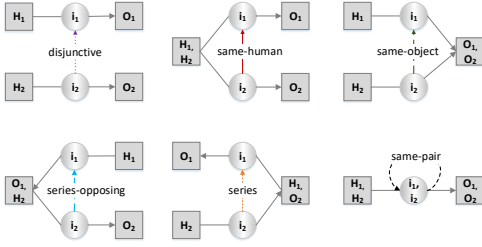


Figure 3. Definition of six kinds of inter-interaction semantic dependencies $\langle \text{HOI}(i_2) \rightarrow \text{HOI}(i_1) \rangle$ between interaction $\text{HOI}(i_1)$ and $\text{HOI}(i_2)$ (square: human/object instance, circle: interaction).

Human-Object Pairs Construction. Here we connect each pair of detected human and object instances, yielding all possible human-object pairs within the input image. Each human-object pair can be represented from three perspectives, i.e., the appearance feature, spatial feature, and linguistic feature of human and object. In particular, the appearance feature is directly represented as the concatenation of human and object instance features derived from DETR (i.e., the 256-dimensional region feature before final prediction heads). By defining the normalized center coordinates of human and object bounding boxes as (c_x^h, c_y^h) and (c_x^o, c_y^o) , we measure the spatial feature as the concatenation of all geometric properties, i.e., $[dx, dy, dis, \arctan(\frac{dy}{dx}), A_h, A_o, I, U]$, where $dx = c_x^h - c_x^o, dy = c_y^h - c_y^o, dis = \sqrt{dx^2 + dy^2}$. A_h, A_o, I, U denote the areas of human, object, their intersection, and union boxes, respectively. The linguistic feature is achieved by encoding the label name of object (one-hot vector) into 300-dimensional vector. The final representation of each human-object pair is calculated as the concatenation of appearance, spatial, and linguistic features.

Interactiveness Prediction. The interactiveness prediction module in IPN takes the feature of each human-object pair as input, and learns to predict the probability whether interactions exist between this pair, i.e., interactiveness score. We frame this sub-task of interactiveness prediction as binary classification problem, and implement this module as MLP coupled with Sigmoid activation. During training, for each input image, we sample at most K human-object pairs, which consist of positive and negative pairs. Note that if both IoUs of predicted human and object bounding boxes in one human-object pair w.r.t ground-truths are larger than 0.5, we treat this pair as positive sample, otherwise it is a negative sample. One natural way to fetch negative pairs is to use randomly sampling strategy. Instead, here we employ hard mining strategy [35] to sample negative pairs with high predicted interactiveness scores, aiming to facilitate the learning of interactiveness prediction. After feeding all the N sampled human-object pairs in a mini-batch into interactiveness prediction module, we optimize this module with focal loss [25] (FL):

$$L_{proposal} = \frac{1}{\sum_{i=1}^N z_i} \sum_{i=1}^N FL(\hat{z}_i, z_i), \quad (1)$$

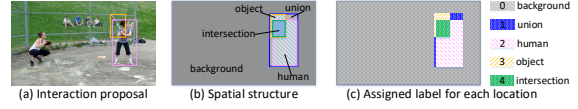


Figure 4. Definition of intra-interaction spatial structure for each interaction: (a) interaction proposal in an image; (b) the spatial structure, i.e., the layout of each component in this interaction; (c) the assigned label for each location in image.

where $z_i \in \{0, 1\}$ indicates whether interactions exist in ground-truth and \hat{z}_i is the predicted interactiveness score. At inference, only the top- K human-object pairs with highest interactiveness scores are taken as interaction proposals.

3.2. Interaction-Centric Graph

Based on all the selected interaction proposals of each input image via IPN, we next present how to construct an interaction-centric graph that fully unfolds the rich prior knowledge of inter- and intra-interaction structures. Technically, we take each interaction proposal as one graph node, and the interaction-centric complete graph is thus built by densely connecting every two nodes as graph edges.

Inter-interaction Semantic Structure. Intuitively, there exists a natural semantic structure among interactions within a same image. For example, when we find the interaction of “human hold mouse” in an image, it is very likely that the mentioned “human” is associated with another interaction of “human look-at screen.” This motivates us to exploit such common sense knowledge implied in the inter-interaction semantic structure to boost relational reasoning among interactions for HOI detection. Formally, we express the directional semantic connectivity as $\langle \text{HOI}(i_2) \rightarrow \text{HOI}(i_1) \rangle$, which denotes the relative semantic dependency of interaction proposal $\text{HOI}(i_1)$ against interaction proposal $\text{HOI}(i_2)$. Six kinds of inter-interaction semantic dependencies are thus defined according to whether two interaction proposals share the same human or object instance, as shown in Figure 3.

Concretely, if $\text{HOI}(i_1)$ and $\text{HOI}(i_2)$ do not share any human/object instance, we classify their dependency as “disjunctive” (class 0). If $\text{HOI}(i_1)$ and $\text{HOI}(i_2)$ only share the same human/object instance, we set the label of dependency as “same-human” (class 1) or “same-object” (class 2). When the human/object instance of $\text{HOI}(i_1)$ is exactly the object/human instance of $\text{HOI}(i_2)$, the dependency is classified as “series-opposing” (class 3) and “series” (class 4), respectively. If both of the human and object instances of $\text{HOI}(i_1)$ and $\text{HOI}(i_2)$ are same, the label of this dependency is “same-pair” (class 5).

Intra-interaction Spatial Structure. The inter-interaction semantic structure over the whole interaction-centric graph only unfolds the holistically semantic dependencies across all interaction proposals, while leaving the locally spatial structure of human/object within each interaction proposal unexploited. Therefore, we characterize each graph node with an intra-interaction spatial structure,

which can be interpreted as the layout of each component in the corresponding interaction proposal (see Figure 4). Specifically, we first identify the spatial location of each component (i.e., *background*, *union*, *human*, *object*, and *intersection*) for this interaction over the whole image, and then assign layout label $l_{ij} \in \{0, 1, 2, 3, 4\}$ to each location in this image according to the corresponding component.

3.3. Structure-aware Transformer

With the K interaction proposals and the interaction-centric graph, we next present how to integrate the prior knowledge of inter- and intra-interaction structures into relational reasoning for HOI set prediction in STIP. In particular, a structure-aware Transformer is devised to contextually encode all interaction proposals with additional guidance of inter- and intra-interaction structures via structure-aware self-attention and cross-attention modules, yielding structure-aware HOI features for predicting HOI triplets.

Preliminary. We first briefly recall the widely adopted vanilla Transformer in vision tasks [19, 20, 30, 31] that capitalizes on attention mechanism, which aims to transform a sequence of queries $\mathbf{q} = (\mathbf{q}_1, \dots, \mathbf{q}_m)$ plus a set of key-value pairs $(\mathbf{k} = (\mathbf{k}_1, \dots, \mathbf{k}_n), \mathbf{v} = (\mathbf{v}_1, \dots, \mathbf{v}_n))$ into the output sequence $\mathbf{o} = (\mathbf{o}_1, \dots, \mathbf{o}_m)$. Each output element \mathbf{o}_i is computed by aggregating all values weighted with attention: $\mathbf{o}_i = \sum_j \alpha_{ij} (\mathbf{W}_v \mathbf{v}_j)$, where each attention weight α_{ij} is normalized with softmax ($\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_j \exp(e_{ij})}$). Here the primary attention weight e_{ij} is measured as the scaled dot-product between each key \mathbf{k}_j and query \mathbf{q}_i :

$$e_{ij} = \frac{(\mathbf{W}_q \mathbf{q}_i)^T (\mathbf{W}_k \mathbf{k}_j)}{\sqrt{d_{key}}}. \quad (2)$$

Note that d_{key} is the dimension of keys, and $\mathbf{W}_q, \mathbf{W}_k, \mathbf{W}_v$ are learnable embedding matrices.

Structure-aware Self-attention. Existing Transformer-type HOI detectors perform relational reasoning among interactions via self-attention module in vanilla Transformer for HOI set prediction. However, the relational reasoning process in vanilla Transformer is triggered by the parametric interaction queries, and leaves the prior knowledge of inter-interaction structure under-exploited. As an alternative, our structure-aware Transformer starts HOI set prediction from the non-parametric queries (i.e., the selected interaction proposals), and further upgrades the conventional relation reasoning with inter-interaction semantic structure through structure-aware self-attention module.

Specifically, by taking the K interaction proposals \mathbf{q} as interaction queries, keys, and values, the structure-aware self-attention module conducts the inter-interaction structure-aware reasoning among interactions to strengthen the HOI representation of each interaction. Inspired by relative position encoding in [34], we supplement each key \mathbf{q}_j

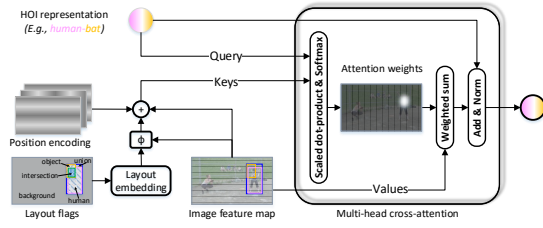


Figure 5. Structure-aware cross-attention module.

with the encodings of its inter-interaction semantic dependency with regard to query \mathbf{q}_i , which is measured as the concatenation of \mathbf{q}_j and the corresponding semantic dependency label $d_{ij} \in \{0, 1, \dots, 5\}$. In this way, we incorporate the inter-interaction semantic structure into the learning of attention weight by modifying Eq. (2) as:

$$e_{ij}^{self} = \frac{(\mathbf{W}_q \mathbf{q}_i)^T (\mathbf{W}_k \mathbf{q}_j + \psi(\mathbf{q}_j, \mathbf{E}_{dep}(d_{ij})))}{\sqrt{d_{key}}}, \quad (3)$$

where \mathbf{E}_{dep} denotes the embedding matrix of semantic dependency label and ψ is implemented as a 2-layer MLP to encode the inter-interaction semantic dependency. Accordingly, the output intermediate HOI features $\hat{\mathbf{q}}$ of structure-aware self-attention module are endowed with the holistically semantic structure among interactions.

Structure-aware Cross-attention. Next, based on the intermediate HOI features $\hat{\mathbf{q}}$, a structure-aware cross-attention module (see Figure 5) is utilized to further enhance HOI features by exploiting contextual information between interactions and the original image feature map in DETR. Formally, we take the K intermediate HOI features $\hat{\mathbf{q}} = (\hat{\mathbf{q}}_1, \dots, \hat{\mathbf{q}}_K)$ as queries, and the image feature map $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ as keys and values. The structure-aware cross-attention module performs the intra-interaction structure-aware reasoning over the image feature map to strengthen the HOI feature of each interaction. Similar to structure-aware self-attention module, each key \mathbf{x}_j is supplemented with the encodings of the intra-interaction spatial structure with regard to query $\hat{\mathbf{q}}_i$ (i.e., the concatenation of \mathbf{x}_j and its assigned layout label $l_{ij} \in \{0, 1, 2, 3, 4\}$). The learning of attention weight in structure-aware cross-attention module is thus integrated with the intra-interaction spatial structure, which is measured as:

$$e_{ij}^{cross} = \frac{(\mathbf{W}_{\hat{q}} \hat{\mathbf{q}}_i)^T (\mathbf{W}_{\hat{k}} \mathbf{x}_j + \mathbf{pos}_j + \phi(\mathbf{x}_j, \mathbf{E}_{lay}(l_{ij})))}{\sqrt{d_{key}}}, \quad (4)$$

where \mathbf{pos}_j is the position encoding, \mathbf{E}_{lay} is the embedding matrix of layout label, and we implement ϕ as a 2-layer MLP to encode the intra-interaction spatial structure.

3.4. Training Objective

During training, we feed the final output HOI representations of structure-aware Transformer into the interaction classifier (implemented as a 2-layer MLP) to predict the interaction classes of each interaction proposal. The objective

of interaction classification is measured via focal loss:

$$L_{cls} = \frac{1}{\sum_{i=1}^N \sum_{c=1}^C y_{ic}} \sum_{i=1}^N \sum_{c=1}^C FL(\hat{y}_{ic}, y_{ic}), \quad (5)$$

where C is the number of interaction classes, $y_{ic} \in \{0, 1\}$ indicates whether the labels of i -th proposal contain the c -th interaction class, and \hat{y}_{ic} is the predicted probability of c -th interaction class. Accordingly, the overall objective of our STIP integrates the interactiveness prediction objective in Eq. (1) and interaction classification objective in Eq. (5):

$$L_{STIP} = L_{proposal} + L_{cls}. \quad (6)$$

4. Experiments

Here we empirically evaluate STIP on two common HOI detection datasets, i.e., V-COCO [11] and HICO-DET [4].

4.1. Datasets and Experimental Settings

V-COCO is a popular dataset for benchmarking HOI detection, which is a subset of MS-COCO [26] covering 29 action categories. This dataset consists of 2,533 training images, 2,867 validation images, and 4,946 testing images. Following the settings in [17], we adopt Average Precision (AP_{role}) over 25 interactions as evaluation metric. Two kinds of AP_{role} , i.e., $AP_{role}^{\#1}$ and $AP_{role}^{\#2}$, are reported under two scenarios with different scoring criterions for object occlusion cases. Specifically, in the scenario of $AP_{role}^{\#1}$, the model should manage to infer the occluded object correctly by predicting the 2D location of its bounding box as $[0,0,0,0]$, meanwhile precisely localizing the corresponding human bounding box and recognizing the interaction in between. In contrast, for the scenario of $AP_{role}^{\#2}$, there is no need to infer the occluded object.

HICO-DET is a larger HOI detection benchmark, which contains 37,536 and 9,515 images for training and testing, respectively. The whole dataset covers 600 categories of $\langle human, object, interaction \rangle$ triplets, covering the same 80 object categories as in MS-COCO [26] and 117 verb categories. We follow [4] and report mAP in two different settings (*Default* and *Known Object*). Here the *Default* setting represents that the mAP is calculated over all testing images, while *Known Object* measures the AP of each object solely over the images containing that object class. For each setting, we report the AP over three different HOI category sets, i.e., **Full** (all 600 HOI categories), **Rare** (138 HOI categories where each one contains less than 10 training samples), and **Non-Rare** (462 HOI categories where each one contains 10 or more training samples).

Implementation Details. For fair comparison with state-of-the-art baselines, we adopt the same object detector DETR pre-trained over MS-COCO (backbone: ResNet-50) and all learnable parameters in DETR are frozen dur-

Method	Backbone	Feature	$AP_{role}^{\#1}$	$AP_{role}^{\#2}$
<i>One-stage methods</i>				
UnionDet [16]	R50	A	47.5	56.2
IPNet [40]	HG-104	A	51.0	-
GGNet [48]	HG-104	A	54.7	-
HOITrans [50]	R50	A	52.9	-
AS-Net [6]	R50	A	53.9	-
HOTR [17]	R50	A	55.2	64.4
QPIC [36]	R50	A	58.8	61.0
<i>Two-stage methods</i>				
InteractNet [10]	R50-FPN	A	40.0	48.0
GPNN [32]	R101	A	44.0	-
TIN [23]	R50	A+S+P	48.7	-
DRG [8]	R50-FPN	A+S+L	51.0	-
FCMNet [27]	R50	A+S+L+P	53.1	-
ConsNet [28]	R50-FPN	A+S+L	53.2	-
IDN [22]	R50	A+S	53.3	60.3
STIP (Ours)	R50	A	65.1	69.7
STIP (Ours)	R50	A+S+L	66.0	70.7

Table 1. Performance comparison on V-COCO dataset. The letters in Feature column indicate the input features: **A** (Appearance/Visual features), **S** (Spatial features [9]), **L** (Linguistic feature of label semantic embeddings), **P** (Human pose feature).

ing training as in [17]. On HICO-DET dataset, we additionally report the results by fine-tuning DETR on HICO-DET and the performances by further jointly fine-tuning object detector and HOI detector. In the experiments, we select the top-32 probably interactive human-object pairs as the output interaction proposals of Interaction Proposal Network. Our proposed structure-aware Transformer consists of 6 stacked layers (structure-aware self-attention plus cross-attention modules). The whole architecture is trained over 2 Nvidia 2080ti GPUs with AdamW optimizer. The mini-batch size is 8 and we set the initial learning rate as 5×10^{-5} . The maximum training epoch number is 30.

4.2. Performance Comparisons

V-COCO. Table 1 summarizes the performance comparisons in terms of $AP_{role}^{\#1}$ and $AP_{role}^{\#2}$ on V-COCO. In general, the results across all metrics under the same backbone (ResNet-50, R50 in short) consistently demonstrate that our STIP exhibits better performances against existing techniques, including both one-stage methods (e.g., UnionDet, AS-Net, HOTR, and QPIC) and two-stage methods (e.g., FCMNet, ConsNet, and IDN). The results generally highlight the key advantage of two-phase HOI set prediction and the exploitation of inter- and intra-interaction structures. In particular, the conventional two-stage HOI detectors (e.g., GPNN, TIN, DRG) commonly construct instance-centric graph to mine contextual information among instances. Instead, recent Transformer-style HOI detectors (e.g., HOITrans, AS-Net, HOTR, QPIC) fully capitalize on vanilla Transformer to perform relational reasoning among instances/interactions, thereby leading to performance boosts. However, when only using appearance features (A), the $AP_{role}^{\#1}$ and $AP_{role}^{\#2}$ of HOTR and QPIC

Method	Backbone	Feature	Default			Known Object		
			Full	Rare	Non-Rare	Full	Rare	Non-Rare
<i>Object detector pre-trained on MS-COCO</i>								
InteractNet [10]	R50-FPN	A	9.94	7.16	10.77	-	-	-
GPNN [32]	R101	A	13.11	9.41	14.23	-	-	-
UnionDet [16]	R50	A	14.25	10.23	15.46	19.76	14.68	21.27
TIN [23]	R50	A+S+P	17.22	13.51	18.32	19.38	15.38	20.57
IPNet [40]	R50-FPN	A	19.56	12.79	21.58	22.05	15.77	23.92
DRG [8]	R50-FPN	A+S+L	19.26	17.74	19.71	23.40	21.75	23.89
FCMNet [27]	R50	A+S+L+P	20.41	17.34	21.56	22.04	18.97	23.12
ConsNet [28]	R50-FPN	A+S+L	22.15	17.12	23.65	-	-	-
IDN [22]	R50	A+S	23.36	22.47	23.63	26.43	25.01	26.85
HOTR [17]	R50	A	23.46	16.21	25.60	-	-	-
AS-Net [6]	R50	A	24.40	22.39	25.01	27.41	25.44	28.00
STIP (Ours)	R50	A	28.11	25.85	28.78	31.23	27.93	32.22
STIP (Ours)	R50	A+S+L	28.81	27.55	29.18	32.28	31.07	32.64
<i>Object detector fine-tuned on HICO-DET</i>								
DRG [8]	R50-FPN	A+S+L	24.53	19.47	26.04	27.98	23.11	29.43
ConsNet [28]	R50-FPN	A+S+L	24.39	17.10	26.56	-	-	-
IDN [22]	R50	A+S	26.29	22.61	27.39	28.24	24.47	29.37
HOTR [17]	R50	A	25.10	17.34	27.42	-	-	-
STIP (Ours)	R50	A	29.76	26.94	30.61	32.84	29.05	33.85
STIP (Ours)	R50	A+S+L	30.56	28.15	31.28	33.54	30.93	34.31
<i>Jointly fine-tune object detector & HOI detector on HICO-DET</i>								
UnionDet [16]	R50	A	17.58	11.72	19.33	19.76	14.68	21.27
PPDM [24]	HG104	A	21.73	13.78	24.10	24.58	16.65	26.84
GGNet [48]	HG104	A	29.17	22.13	30.84	33.50	26.67	34.89
HOITrans [50]	R50	A	23.46	16.91	25.41	26.15	19.24	28.22
AS-Net [6]	R50	A	28.87	24.25	30.25	31.74	27.07	33.14
QPIC [36]	R50	A	29.07	21.85	31.23	31.68	24.14	33.93
STIP (Ours)	R50	A	31.60	27.75	32.75	34.41	30.12	35.69
STIP (Ours)	R50	A+S+L	32.22	28.15	33.43	35.29	31.43	36.45

Table 2. Performance comparison on HICO-DET dataset. The letters in Feature column indicate the input features: **A** (Appearance/Visual features), **S** (Spatial features [9]), **L** (Linguistic feature of label semantic embeddings), **P** (Human pose feature).

Method	V-COCO		HICO-DET (Default)		
	$AP_{role}^{\#1}$	$AP_{role}^{\#2}$	Full	Rare	Non-Rare
Base	52.49	58.25	21.74	18.09	22.83
+HM	58.45	62.64	24.16	19.45	25.57
+HM+TR	63.50	68.07	28.62	26.09	29.38
+HM+TR ^{SS}	64.99	69.94	29.65	26.52	30.59
+HM+TR ^{SC}	65.04	69.76	29.74	27.07	30.54
+HM+TR ^{SS+SC} (STIP)	66.04	70.65	30.56	28.15	31.28

Table 3. Performance contribution of each component in our STIP. **HM**: Hard Mining strategy for training interaction proposal network. **TR**: vanilla TRansformer. **TR^{SS}**: TRansformer with only Structure-aware Self-attention that exploits inter-interaction structure. **TR^{SC}**: TRansformer with only Structure-aware Cross-attention that exploits intra-interaction structure.

are still lower than our STIP, which not only takes interaction proposals as non-parametric interaction queries to trigger HOI set prediction, but also leverages a structure-aware Transformer to exploit the prior knowledge of inter-interaction and intra-interaction structures. For our STIP, a further performance improvement is attained when utilizing more kinds of features (e.g., spatial and linguistic features).

HICO-DET. We further evaluate our STIP on the more challenging HICO-DET dataset. Table 2 reports the mAP scores over three different HOI category sets for each setting (Default/Known Object) in comparison with the state-of-the-art methods. Here we include three different training settings, i.e., pre-trained object detector on MS-COCO,

fine-tune object detector on HICO-DET, and jointly fine-tune object detector and HOI detector on HICO-DET, for fair comparison. Similar to the observations on V-COCO, our STIP achieves consistent performance gains against existing HOI detectors across all the metrics for each training setting. The results basically demonstrate the advantage of triggering HOI set prediction with the non-parametric interaction proposals and meanwhile exploiting the holistically semantic structure among interaction proposals & the locally spatial structure within each interaction proposal.

4.3. Experimental Analysis

Ablation Study. To examine the impact of each design in STIP, we conduct ablation study by comparing different variants of STIP on V-COCO and HICO-DET datasets in Table 3. Note that all experiments on HICO-DET here are conducted under the training setting of object detector fine-tuned on HICO-DET. We start from the basic model (**Base**), which utilizes a basic interaction proposal network (randomly sampling negative samples for training, without hard mining strategy). The generated interaction proposals in Base model are directly leveraged for interaction classification, without any Transformer-style structure for boosting HOI prediction. Next, we extend Base model by leveraging hard mining strategy to select the hard negative human-object pairs with higher interac-

# of selected interaction proposals (K)	V-COCO		HICO-DET (Default)		
	$AP_{role}^{\#1}$	$AP_{role}^{\#2}$	Full	Rare	Non-Rare
8	64.20	69.11	29.03	28.16	29.29
16	65.68	70.63	30.18	28.66	30.64
<u>32</u>	66.04	70.65	30.56	28.15	31.28
64	65.93	70.50	30.72	28.96	31.24
100	65.78	70.45	30.40	27.89	31.14

Table 4. Performance comparison by using different number of selected interaction proposals (K) for interaction-centric graph construction in our STIP.

tiveness scores for training interaction proposal network, yielding **Base+HM** which achieves better performances. After that, by additionally involving a vanilla Transformer to perform relational reasoning among interaction proposals, another variant of our model (**Base+HM+TR**) leads to performance improvements across all metrics. Furthermore, we upgrade the vanilla Transformer with structure-aware self-attention that exploits the holistically semantic structure among interaction proposals, and this ablated run (**Base+HM+TR^{SS}**) outperforms Base+HM+TR. Meanwhile, the vanilla Transformer can be upgraded with structure-aware cross-attention that exploits the locally spatial structure within each interaction proposal, and **Base+HM+TR^{SC}** also exhibits better performances. These observations basically validate the merit of exploiting the structured prior knowledge, i.e., inter-interaction or intra-interaction structure, for HOI detection. Finally, when jointly upgrading the vanilla Transformer with structure-aware self-attention and structure-aware cross-attention (i.e., our **STIP**), the highest performances are attained.

Effect of Selected Interaction Proposal Number K for Interaction-centric Graph Construction. Recall that the interaction proposal network in our STIP selects only the top- K human-object pairs with highest interactiveness scores as the output interaction proposals for constructing the interaction-centric graph. Such K selected interaction proposals are also taken as non-parametric interaction queries to trigger HOI set prediction in the structure-aware Transformer. Here we vary K from 8 to 100 to explore the relationship between the performance and the select interaction proposal number K . As shown in Table 4, the best performances across most metrics are attained when K is set as 32. In particular, enlarging the number of selected interaction proposals (until $K = 32$) can generally lead to performance boosts on two datasets. Once K is larger than 64, the performances slightly decrease. We speculate that the increase of selected interaction proposals result in more invalid proposals, which may affect the overall stability of relational reasoning among interaction proposals. Accordingly, we empirically set K as 32.

Effect of Layer Number L in Structure-aware Transformer. To explore the effect of layer number L in structure-aware Transformer, we show the performances on two benchmarks by varying this number from 0 to 8. The best performances across most metrics are achieved

# of layers (L)	V-COCO		HICO-DET (Default)		
	$AP_{role}^{\#1}$	$AP_{role}^{\#2}$	Full	Rare	Non-Rare
0	58.45	62.64	24.16	19.45	25.57
1	64.83	69.57	29.21	26.37	30.06
2	65.55	70.39	30.02	28.11	30.59
4	66.02	70.61	30.47	29.28	30.83
<u>6</u>	66.04	70.65	30.56	28.15	31.28
8	65.44	70.11	30.93	29.78	31.27

Table 5. Performance comparison with different layer numbers of the structure-aware Transformer in our STIP.

when the layer number is set to $L = 6$. Specifically, in the extreme case of $L = 0$, no self-attention and cross-attention module is utilized, and the model degenerates to a Base+HM model that directly performs interaction classification over interaction proposals without any relational reasoning via Transformer-style structure. When increasing the layer number in structure-aware Transformer, the performances are gradually increased in general. This basically validates the effectiveness of enabling relational reasoning among interaction proposals through structure-aware Transformer. In practice, the layer number L is generally set to 6.

Time Analysis. We evaluate the inference time of our STIP on a single Nvidia 2080ti GPU by constructing each batch with single testing image. Specifically, for each input batch, object detection via DETR, interaction proposal generation through interaction proposal network, HOI set prediction with structure-aware Transformer, and the other processing (e.g., data loading) takes 41.9ms, 7.8ms, 20.4ms, and 3.8ms, respectively. Consequently, the overall inference stage of STIP finishes in 73.9ms on average, which is comparable to existing one-stage Transformer-style HOI detectors (e.g., the inference time of AS-Net [6] is 71ms).

5. Conclusion and Discussion

In this paper, we have presented STIP, a new end-to-end Transformer-style model for human-object interaction detection. Instead of performing HOI set prediction derived from parametric interaction queries in a one-stage manner, the proposed STIP capitalizes on a two-phase solution for HOI detection by first producing interaction proposals and then taking them as non-parametric interaction queries to trigger HOI set prediction. Furthermore, by going beyond the commonly adopted vanilla Transformer, a novel structure-aware Transformer is designed to exploit two kinds of structured prior knowledge, i.e., inter- and intra-interaction structures, to further boost HOI set prediction. We validate the proposed scheme and analysis through extensive experiments conducted on V-COCO and HICO-DET datasets. More importantly, the proposed STIP achieves new state-of-the-art results on both benchmarks.

Broader Impact. STIP has high potential impact in human-centric applications, such as sports analysis and self-driving vehicles. However, such HOI detection technology can be deployed in human monitoring and surveillance as well which might raise ethical and privacy issues.

References

- [1] Amit Adam, Ehud Rivlin, Ilan Shimshoni, and Daviv Reinitz. Robust real-time unusual event detection using multiple fixed-location monitors. *PAMI*, 2008. 1
- [2] Brenna D Argall, Sonia Chernova, Manuela Veloso, and Brett Browning. A survey of robot learning from demonstration. *Robotics and autonomous systems*, 2009. 1
- [3] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020. 1, 3
- [4] Yu-Wei Chao, Yunfan Liu, Xieyang Liu, Huayi Zeng, and Jia Deng. Learning to detect human-object interactions. In *WACV*, 2018. 2, 6
- [5] Yu-Wei Chao, Zhan Wang, Yugeng He, Jiaxuan Wang, and Jia Deng. Hico: A benchmark for recognizing human-object interactions in images. In *ICCV*, 2015. 1, 2
- [6] Mingfei Chen, Yue Liao, Si Liu, Zhiyuan Chen, Fei Wang, and Chen Qian. Reformulating hoi detection as adaptive set prediction. In *CVPR*, 2021. 1, 2, 3, 6, 7, 8
- [7] Mihai Dogariu, Liviu-Daniel Stefan, Mihai Gabriel Constantin, and Bogdan Ionescu. Human-object interaction: Application to abandoned luggage detection in video surveillance scenarios. In *COMM*, 2020. 1
- [8] Chen Gao, Jiarui Xu, Yuliang Zou, and Jia-Bin Huang. Drg: Dual relation graph for human-object interaction detection. In *ECCV*, 2020. 1, 2, 6, 7
- [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. 1, 2, 6, 7
- [10] Georgia Gkioxari, Ross Girshick, Piotr Dollár, and Kaiming He. Detecting and recognizing human-object interactions. In *CVPR*, 2018. 1, 2, 6, 7
- [11] Saurabh Gupta and Jitendra Malik. Visual semantic role labeling. *arXiv preprint arXiv:1505.04474*, 2015. 1, 2, 6
- [12] Tanmay Gupta, Alexander Schwing, and Derek Hoiem. No-frills human-object interaction detection: Factorization, layout encodings, and training techniques. In *ICCV*, 2019. 2
- [13] Tao He, Lianli Gao, Jingkuan Song, and Yuan-Fang Li. Exploiting scene graphs for human-object interaction detection. *arXiv preprint arXiv:2108.08584*, 2021. 2
- [14] Zhi Hou, Baosheng Yu, Yu Qiao, Xiaojiang Peng, and Dacheng Tao. Affordance transfer learning for human-object interaction detection. In *CVPR*, 2021. 2
- [15] Zhi Hou, Baosheng Yu, Yu Qiao, Xiaojiang Peng, and Dacheng Tao. Detecting human-object interaction via fabricated compositional learning. In *CVPR*, 2021. 2
- [16] Bumsoo Kim, Taeho Choi, Jaewoo Kang, and Hyunwoo J Kim. Uniondet: Union-level detector towards real-time human-object interaction detection. In *ECCV*, 2020. 1, 2, 6, 7
- [17] Bumsoo Kim, Junhyun Lee, Jaewoo Kang, Eun-Sol Kim, and Hyunwoo J Kim. Hotr: End-to-end human-object interaction detection with transformers. In *CVPR*, 2021. 1, 2, 3, 6, 7
- [18] Yikang Li, Wanli Ouyang, Xiaogang Wang, and Xiao'ou Tang. Vip-cnn: Visual phrase guided convolutional neural network. In *CVPR*, 2017. 2
- [19] Yehao Li, Yingwei Pan, Ting Yao, Jingwen Chen, and Tao Mei. Scheduled sampling in vision-language pretraining with decoupled encoder-decoder network. In *AAAI*, 2021. 5
- [20] Yehao Li, Ting Yao, Yingwei Pan, and Tao Mei. Contextual transformer networks for visual recognition. *IEEE Trans. on PAMI*, 2022. 5
- [21] Yong-Lu Li, Xinpeng Liu, Han Lu, Shiyi Wang, Junqi Liu, Jiefeng Li, and Cewu Lu. Detailed 2d-3d joint representation for human-object interaction. In *CVPR*, 2020. 2
- [22] Yong-Lu Li, Xinpeng Liu, Xiaoqian Wu, Yizhuo Li, and Cewu Lu. Hoi analysis: Integrating and decomposing human-object interaction. *arXiv preprint arXiv:2010.16219*, 2020. 2, 6, 7
- [23] 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 *CVPR*, 2019. 2, 6, 7
- [24] Yue Liao, Si Liu, Fei Wang, Yanjie Chen, Chen Qian, and Jiashi Feng. Ppdm: Parallel point detection and matching for real-time human-object interaction detection. In *CVPR*, 2020. 1, 2, 7
- [25] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, 2017. 4
- [26] 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 *ECCV*, 2014. 6
- [27] Yang Liu, Qingchao Chen, and Andrew Zisserman. Amplifying key cues for human-object-interaction detection. In *ECCV*, 2020. 1, 2, 6, 7
- [28] Ye Liu, Junsong Yuan, and Chang Wen Chen. Consnet: Learning consistency graph for zero-shot human-object interaction detection. In *ACM MM*, 2020. 2, 6, 7
- [29] Li Mi and Zhenzhong Chen. Hierarchical graph attention network for visual relationship detection. In *CVPR*, 2020. 2
- [30] Yingwei Pan, Yehao Li, Jianjie Luo, Jun Xu, Ting Yao, and Tao Mei. Auto-captions on gif: A large-scale video-sentence dataset for vision-language pre-training. *arXiv preprint arXiv:2007.02375*, 2020. 5
- [31] Yingwei Pan, Ting Yao, Yehao Li, and Tao Mei. X-linear attention networks for image captioning. In *CVPR*, 2020. 5
- [32] Siyuan Qi, Wenguan Wang, Baoxiong Jia, Jianbing Shen, and Song-Chun Zhu. Learning human-object interactions by graph parsing neural networks. In *ECCV*, 2018. 1, 2, 6, 7
- [33] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *NeurIPS*, 2015. 2
- [34] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations. In *NAACL-HLT*, 2018. 5
- [35] Abhinav Shrivastava, Abhinav Gupta, and Ross Girshick. Training region-based object detectors with online hard example mining. In *CVPR*, 2016. 4

- [36] Masato Tamura, Hiroki Ohashi, and Tomoaki Yoshinaga. Qpic: Query-based pairwise human-object interaction detection with image-wide contextual information. In *CVPR*, 2021. 1, 2, 3, 6, 7
- [37] Oytun Ulutan, ASM Iftekhar, and Bangalore S Manjunath. Vsgnet: Spatial attention network for detecting human object interactions using graph convolutions. In *CVPR*, 2020. 2
- [38] Hai Wang, Wei-shi Zheng, and Ling Yingbiao. Contextual heterogeneous graph network for human-object interaction detection. In *ECCV*, 2020. 1, 2
- [39] Tiancai Wang, Rao Muhammad Anwer, Muhammad Haris Khan, Fahad Shahbaz Khan, Yanwei Pang, Ling Shao, and Jorma Laaksonen. Deep contextual attention for human-object interaction detection. In *ICCV*, 2019. 1, 2
- [40] Tiancai Wang, Tong Yang, Martin Danelljan, Fahad Shahbaz Khan, Xiangyu Zhang, and Jian Sun. Learning human-object interaction detection using interaction points. In *CVPR*, 2020. 1, 2, 6, 7
- [41] Bingjie Xu, Yongkang Wong, Junnan Li, Qi Zhao, and Mohan S Kankanhalli. Learning to detect human-object interactions with knowledge. In *CVPR*, 2019. 2
- [42] Ting Yao, Yingwei Pan, Yehao Li, and Tao Mei. Exploring visual relationship for image captioning. In *ECCV*, 2018. 2
- [43] 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. *NeurIPS*, 2021. 2
- [44] Frederic Z Zhang, Dylan Campbell, and Stephen Gould. Spatially conditioned graphs for detecting human-object interactions. *arXiv preprint arXiv:2012.06060*, 2020. 2
- [45] Ji Zhang, Kevin J Shih, Ahmed Elgammal, Andrew Tao, and Bryan Catanzaro. Graphical contrastive losses for scene graph parsing. In *CVPR*, 2019. 2
- [46] Yong Zhang, Yingwei Pan, Ting Yao, Rui Huang, Tao Mei, and Chang-Wen Chen. Boosting scene graph generation with visual relation saliency. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 2022. 2
- [47] Xubin Zhong, Changxing Ding, Xian Qu, and Dacheng Tao. Polysemy deciphering network for human-object interaction detection. In *ECCV*, 2020. 2
- [48] Xubin Zhong, Xian Qu, Changxing Ding, and Dacheng Tao. Glance and gaze: Inferring action-aware points for one-stage human-object interaction detection. In *CVPR*, 2021. 2, 6, 7
- [49] Penghao Zhou and Mingmin Chi. Relation parsing neural network for human-object interaction detection. In *ICCV*, 2019. 2
- [50] Cheng Zou, Bohan Wang, Yue Hu, Junqi Liu, Qian Wu, Yu Zhao, Boxun Li, Chenguang Zhang, Chi Zhang, Yichen Wei, et al. End-to-end human object interaction detection with hoi transformer. In *CVPR*, 2021. 1, 2, 3, 6, 7