

MSTR: Multi-Scale Transformer for End-to-End Human-Object Interaction Detection

Bumsoo Kim*
LG AI Research

bumsoo.kim@lgresearch.ai

Jonghwan Mun
Kakao Brain

Kyoung-Woon On
Kakao Brain

Minchul Shin
Kakao Brain

Junhyun Lee
Korea University

Eun-Sol Kim
Department of Computer Science
Hanyang University

Abstract

Human-Object Interaction (HOI) detection is the task of identifying a set of $\langle \text{human}, \text{object}, \text{interaction} \rangle$ triplets from an image. Recent work proposed transformer encoder-decoder architectures that successfully eliminated the need for many hand-designed components in HOI detection through end-to-end training. However, they are limited to single-scale feature resolution, providing suboptimal performance in scenes containing humans, objects, and their interactions with vastly different scales and distances. To tackle this problem, we propose a Multi-Scale Transformer (MSTR) for HOI detection powered by two novel HOI-aware deformable attention modules called Dual-Entity attention and Entity-conditioned Context attention. While existing deformable attention comes at a huge cost in HOI detection performance, our proposed attention modules of MSTR learn to effectively attend to sampling points that are essential to identify interactions. In experiments, we achieve the new state-of-the-art performance on two HOI detection benchmarks.

1. Introduction

Human-Object Interaction (HOI) detection is a task to predict a set of $\langle \text{human}, \text{object}, \text{interaction} \rangle$ triplets in an image [9]. Previous methods have indirectly addressed this task by detecting human and object instances and individually inferring interaction labels for every pair of the detected instances with either neural networks (*i.e.*, two-stage HOI detectors [1, 6–8, 10, 17, 19, 20, 22–25, 27, 29–31, 34, 35, 37]) or triplet matching (*i.e.*, one-stage HOI detectors [13, 21, 32]). The additional complexity caused by this indirect inference structure and post-processing (*e.g.*,

*this work was done in Kakao Brain

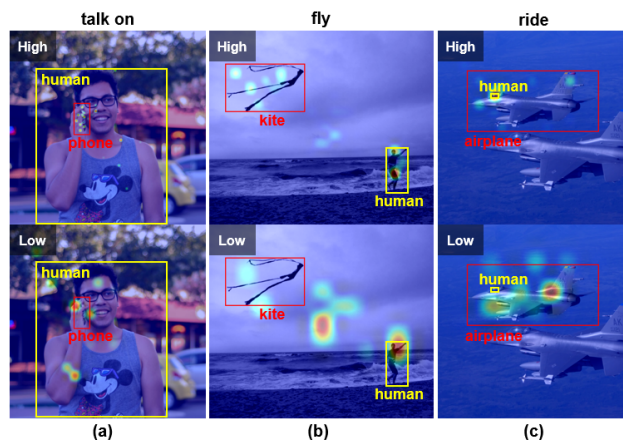


Figure 1. Multi-scale attention of MSTR on interactions including: (a) large human with small object, (b) distant human and object, and (c) small human and a large object. The top row (high resolution) and the bottom row (low resolution) captures the context of the interaction in various scales. Best viewed in color.

NMS) stage behaved as a major bottleneck in inference time in HOI detection. To deal with this bottleneck, transformer-based HOI detectors [4, 14, 26, 39] have been proposed to achieve end-to-end HOI detection without the need for the post-processing stage mentioned above. These works have shown competitive performance in both accuracy and inference time with direct set-level prediction and transformer attentions that can exploit the contextual information between humans, objects, and their interactions.

However, due to the huge computational costs raised when processing multi-scale feature maps (with about $20\times$ more image tokens) with transformer attention, current transformer-based HOI detectors are limited to using only single-scale feature maps. Due to this limitation, previous transformer-based approaches demonstrate suboptimal per-

formance, especially for scenes where humans, objects, and the contextual information for their interactions exist at various scales.

In this paper, we propose Multi-Scale TRAsnformer (MSTR), a transformer-based HOI detector that can exploit multi-scale feature maps for HOI detection. Inspired by previously proposed deformable attention for standard object detection [38], we aim to efficiently explore multi-scale feature maps by attending to only a small number of sampling points generated from the query element instead of calculating the attention values for the entire spatial dimension. Yet, we found out in our preliminary experiments that directly applying naïve deformable attention in HOI detection leads to a serious performance drop.

To overcome this deterioration, we equipped MSTR with two novel *HOI-aware* deformable attentions, referred by Dual-Entity Attention and Entity-conditioned Context Attention, which are designed to capture the complicated semantics of Human-Object Interaction throughout multi-resolution feature maps (see Figure 1). Specifically, precise entity-level semantics for humans and objects are captured by *Dual-Entity attention*, while the contextual information for the interaction is conditionally reimbursed by *Entity-conditioned Context attention*. To further improve performance, we delve into decoder architectures that can effectively handle the multiple semantics obtained from the two HOI-aware attentions above.

The main contributions of our work are threefold:

- We propose MSTR, the first transformer-based HOI detector that exploits multi-scale visual feature maps.
- We propose new deformable attention modules, called Dual-Entity attention and Entity-conditioned Context attention, which effectively and efficiently capture human, object, and context information associated with HOI queries.
- We explore decoder architectures to handle the multiple semantics captured by our proposed deformable attentions and further improve HOI detection performance.

2. Preliminary

In this section, we start with a basic pipeline of a transformer-based end-to-end HOI detector [26]. Then, we explain the deformable attention module [38] that reduces computational cost in attention, thus enabling the transformer to take multi-scale feature maps as an input. Afterward, we discuss why the direct application of multi-scale deformable attentions is not suitable for HOI detection.

2.1. End-to-End HOI Detection with Transformers

Out of the multiple candidates [4, 14, 26, 39] using transformers for HOI detection, we adopt QPIC [26] as our baseline due to its simple structure and good performance.

Set Prediction. Transformer-based HOI detectors formulate the task as a set-level prediction problem. It is achieved by exploiting a fixed number of HOI queries, each of which generates four types of predictions: 1) the coordinate of the human bounding box (*i.e.*, subject of the interaction), 2) the coordinate of the object bounding box (*i.e.*, target of the interaction), 3) the object class and 4) the interaction type. Note that the set-level predictions are learned using losses based on Hungarian Matching with ground-truths.

Transformer Encoder-Decoder Architecture. The architecture of QPIC [26] consists of a backbone CNN, a transformer encoder, and a transformer decoder. Given an image, a single-scale visual feature map is extracted by the backbone CNN, and then positional information is added to the feature map. The transformer encoder takes the visual features and returns contextualized visual features with self-attention layers. In the transformer decoder, HOI queries are first processed by the self-attention layer, and then the cross-attention layer associates the HOI queries with the contextualized visual features (given by the encoder) to capture relevant HOI representations. Finally, predictions for HOI are computed from individual contextualized HOI query embeddings as mentioned above. Note that both self-attention and cross-attention adopt multi-head attention.

To be specific, given a single-scale input feature map $x \in \mathbb{R}^{C \times H \times W}$ where C is the feature dimension, the single-scale multi-head attention $f_q^{sg} = \text{SSAttn}(z_q, x)$ for the q^{th} query feature z_q (either an image token for the encoder or an HOI query for the decoder) is calculated by

$$f_q^{sg} = \sum_{m=1}^M W_m \left[\sum_{k \in \Omega_k} A_{mqk} \cdot W'_m x_k \right], \quad (1)$$

where A_{mqk} indicates an attention weight calculated with learnable weights $U_m, V_m \in \mathbb{R}^{C_v \times C}$ as $\exp\left(\frac{z_q^T U_m^T V_m x_k}{\sqrt{C_v}}\right)$. Throughout this paper, for the attention module, we let m index the attention head ($1 \leq m \leq M$), $q \in \Omega_q$ indexes a query element with feature $z_q \in \mathbb{R}^C$, $k \in \Omega_k$ indexes a key element with feature $z_k \in \mathbb{R}^C$, while Ω_q and Ω_k specify the set of query and key elements, respectively. W_m and W'_m are learnable embedding parameters for m^{th} attention head, and A_{mqk} is normalized as $\sum_{k \in \Omega_k} A_{mqk} = 1$.

Complexity. Given an input feature map $x \in \mathbb{R}^{C \times H \times W}$ and N HOI queries, the complexity of transformer encoder and decoder are $O(H^2 W^2 C)$ and $O(HWC^2 +$

$NHWC + 2NC^2 + N^2C$), respectively. Since the complexity grows in quadratic scale as the spatial resolution (H, W) increases, it raises significant complexity when exploiting multi-resolution feature maps where there are about $20\times$ more features to process.

Towards Multi-Scale HOI detection. In HOI detection, not only do humans and objects exist at various scales, but they also interact at various distances in images. Therefore, it is essential to exploit multi-scale feature maps $\{x\}_{l=1}^L$ (where $x^l \in \mathbb{R}^{C \times H_l \times W_l}$, l indexes the feature level) to deal with the various scales of objects and contexts to capture interactions precisely. However, as multi-scale feature maps have almost $\times 20$ more elements to process than a single-scale feature map, it provokes a serious complexity issue in calculating Eq. (1).

2.2. Revisiting Deformable Transformers

The deformable attention module is proposed to deal with the problem of high complexity in the transformer attention. The core idea is to reduce the number of *key* elements in the attention module by sampling the small number of spatial locations related to regions of interest for each *query* element.

Sampling Locations for Deformable Attention. Given a multi-scale input feature map $\{x^l\}_{l=1}^L$ where $x^l \in \mathbb{R}^{C \times H_l \times W_l}$, the K sampling locations of interest for each attention head and each feature level are generated from each *query* element $z_q \in \mathbb{R}^C$. Because direct prediction of coordinates of sampling location is difficult to learn, it is formulated as prediction of a reference point $r_q \in [0, 1]^2$ and K sampling offsets $\Delta r_q \in \mathbb{R}^{M \times L \times K \times 2}$. Then, the k^{th} sampling location at l^{th} feature level and m^{th} attention head for query element z_q is defined by $p_{mlqk} = \phi_l(r_q) + \Delta r_{mlqk}$ where $\phi_l(\cdot)$ is a function to re-scale the coordinate of reference point to the input feature map of the l^{th} level.

Deformable Attention Module. Given a multi-scale input feature map $\{x^l\}_{l=1}^L$, the multi-scale deformable attention $f_q^{ms} = \text{MSDeformAttn}(z_q, p_q, \{x^l\}_{l=1}^L)$ for query element z_q is calculated using a set of predicted sampling locations p_q as follows:

$$f_q^{ms} = \sum_{m=1}^M W_m \left[\sum_{l=1}^L \sum_{k=1}^K A_{mlqk} \cdot W'_m \Phi_{mlqk} \right], \quad (2)$$

where l , k and m index the input feature level, the sampling location and the attention head, respectively, while A_{mlqk} indicates an attention weight for the k^{th} sampling location at the l^{th} feature level and the m^{th} attention head. Φ_{mlqk} means the sampled k^{th} key element at l^{th} feature level and

m^{th} attention head using the sampling location, which is obtained by bilinear interpolation as $\Phi_{mlqk} = x^l(p_{mlqk}) = x^l(\phi_l(r_q) + \Delta r_{mlqk})$. Note that for each query element, the attention computation is performed with only sampled regions of interest where the sampled number ($= LMK$) is much smaller than the number of all the key elements ($\sum_{l=1}^L H_l W_l$), thus leads to a reduced computational cost.

Problem with Direct Application to HOI Detection.

Deformable attention effectively reduces the complexity of exploiting multi-scale features with transformers to an acceptable level. However, while the sampling procedure above does not deteriorates performance in standard object detection, it causes a serious performance drop in HOI detection ($29.07 \rightarrow 25.53$) as shown in Table 3. We conjecture that this is partly due to the following reasons. First, unlike the object detection task where an object query is associated with a single object, an HOI query is entangled with multiple semantics (*i.e.*, human, object, and their interaction); thus learning to sample the region of interest for multiple semantics with individual HOI queries (especially with sparse information) is much challenging compared to the counterpart of object detection. Second, deformable attention is learned to attend only to the sampling points near the localized objects; this leads to the loss of contextual information that is an essential clue for precise HOI detection. The following sections describe how we resolve these issues and improve performance.

3. Method

In this section, we introduce MSTR, a novel deformable transformer architecture that is suitable for multi-scale HOI detection. To resolve the problems described in our preliminary, MSTR features new *HOI-aware* deformable attentions designed for HOI detection, referred by Dual-Entity attention and Entity-conditioned Context attention.

3.1. HOI-aware Deformable Attentions

The objective of our HOI-aware deformable attentions (Dual-Entity attention and Entity-conditioned Context attention) is to efficiently and effectively extract information of HOIs from multi-scale feature maps for a given HOI query. Figure 2 shows conceptual illustrations of (a) deformable attention in literature [38], (b) Dual-Entity attentions and (c) Entity-conditioned Context attention.

Dual-Entity attention for Human/Object. In HOI detection, the HOI query includes complex and entangled information of multiple semantics: human, object, and interaction information. Therefore, it is challenging to accurately predict sampling locations appropriate for each semantic from a single HOI query. To make sampling loca-

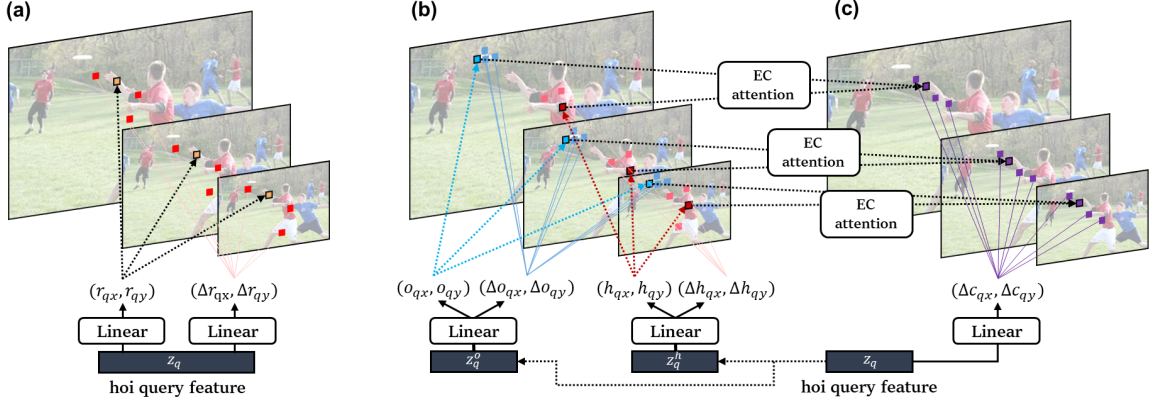


Figure 2. Illustration of (a) Deformable Attention, (b) Dual-Entity Attention, (c) Entity-conditioned Context Attention (abbreviated as EC). The sampling point for deformable attention is obtained by combining the reference points with sampling offset. In (a), both reference points $r_q = (r_{qx}, r_{qy})$ and sampling offsets $\Delta r_q = (\Delta r_{qx}, \Delta r_{qy})$ are obtained from a single hoi query feature z_q . In (b), the reference points and sampling offsets for the humans $h_q = (h_{qx}, h_{qy}), \Delta h_q = (\Delta h_{qx}, \Delta h_{qy})$ and objects $o_q = (o_{qx}, o_{qy}), \Delta o_q = (\Delta o_{qx}, \Delta o_{qy})$ are obtained from z_q^h and z_q^o , respectively, which is obtained by a linear projection of z_q (dotted line). In (c), the sampling offsets $\Delta c_q = (\Delta c_{qx}, \Delta c_{qy})$ are obtained from z_q while the reference points are obtained in conditional to entities in (b).

tions easier, given an HOI query feature z_q , our Dual-Entity attention separately identifies sampling locations for the humans (p_q^h) and objects (p_q^o). First, we project z_q with two linear layers to obtain z_q^h and z_q^o . The k^{th} sampling location at l^{th} feature level and m^{th} attention head for human and object are represented by

$$\begin{aligned} p_{mlqk}^h &= \phi_l(h_q) + \Delta h_{mlqk}, \\ p_{mlqk}^o &= \phi_l(o_q) + \Delta o_{mlqk}, \end{aligned} \quad (3)$$

where $h_q, \Delta h$ is the reference point and sampling offsets for humans, and $o_q, \Delta o$ is the reference point and sampling offsets for objects, each obtained by a linear projection of z_q^h and z_q^o , respectively. Then, based on the sampled locations, attended features for human (f_q^h) and object (f_q^o) are computed by

$$\begin{aligned} f_q^h &= \text{MSDeformAttn}(z_q^h, p_q^h, \{x^l\}_{l=1}^L), \\ f_q^o &= \text{MSDeformAttn}(z_q^o, p_q^o, \{x^l\}_{l=1}^L). \end{aligned} \quad (4)$$

Entity-conditioned Context attention. In HOI detection, contextual information often gives an important clue in identifying interactions. From this point of view, utilizing the local features obtained from near the human and object regions through the Dual-Entity attention is not sufficient to capture contextual information (see our experimental result in Table 3). To compensate for this, we define an attention with an additional set of sampling points, namely Entity-conditioned Context attention, that is designed to capture the contextual information in specific.

Given the 2D reference points for the human $h_q = (h_{qx}, h_{qy})$ and the object $o_q = (o_{qx}, o_{qy})$, the reference

point for Entity-conditioned Context attention is conditionally computed with the two references. Motivated by existing works [21, 32, 36], we define the reference points for interaction as the center of human and object, *i.e.*, $c_q = (\frac{h_{qx} + o_{qx}}{2}, \frac{h_{qy} + o_{qy}}{2})$. Note that we empirically observe that such simple reference points offer competitive performance compared to ones predicted using an additional network, while being much faster. Then, we predict the sampling offsets Δc_q from the HOI query feature, obtaining $p_{mlqk}^c = \phi_l(c_q) + \Delta c_{mlqk}$. Finally, the attended feature for contextual information f_q^c is computed using sampling location p_q^c as follows:

$$f_q^c = \text{MSDeformAttn}(z_q, p_q^c, \{x^l\}_{l=1}^L). \quad (5)$$

3.2. MSTR Architecture

In this section, the overall architecture of MSTR with our suggested two deformable attentions will be described (see Figure 3). MSTR follows the previous transformer encoder-decoder architecture, where the encoder performs self-attention given the image features while the decoder performs self-attention for HOI queries followed by cross-attention between updated HOI queries and the encoded image features.

Encoder. The encoder of MSTR takes multi-scale input feature maps given by a backbone CNN, performs a series of deformable attention modules in Eq.(2), and finally generates encoded feature maps. Positional encoding [2]

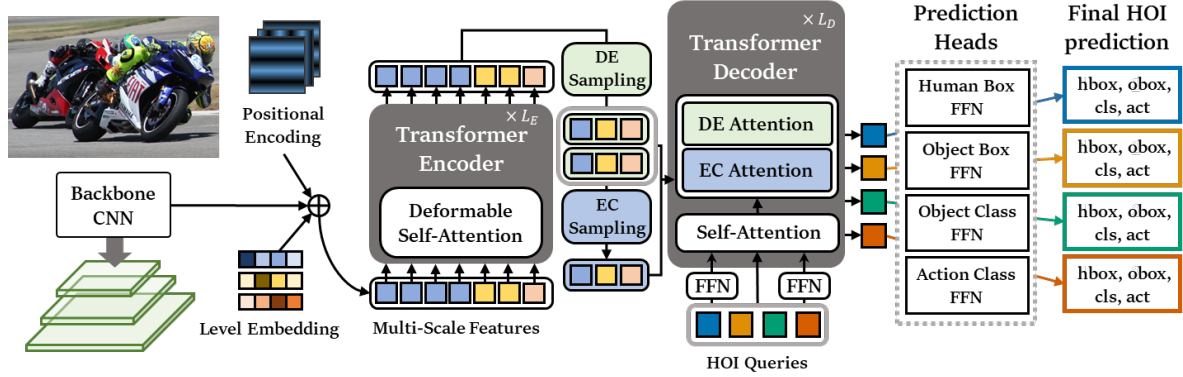


Figure 3. Overall pipeline of MSTR. On top of the standard transformer encoder-decoder architecture for HOI detection (*i.e.*, QPIC), we leverage deformable samplings for the encoder self-attention and the decoder cross-attention modules to deal with the huge complexity caused by using multi-scale feature maps. For the decoder cross-attention, we leverage three sets of key elements sampled for our Dual-Entity attention (denoted as DE sampling, DE attention) and Entity-conditioned Context attention (denoted as EC sampling, EC attention).

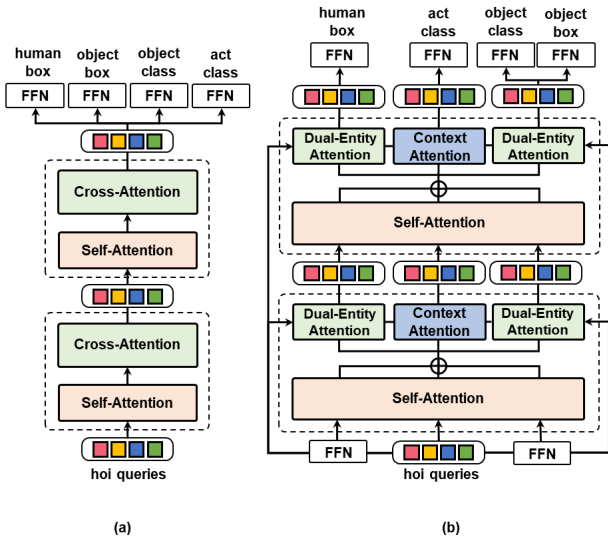


Figure 4. Comparison of a simple 2-layer Decoder architecture for Transformer-based HOI detectors: (a) conventional one introduced in QPIC, and (b) HOI-aware one in MSTR. Entity-conditioned Context attention is abbreviated as Context Attention. MSTR stacks decoder layers by merging the self attention outputs, which further improves performance (see Table 3).

is added to preserve spatial information while level embedding [38] is added to denote which resolution did the image feature comes from.

Decoder. By leveraging our HOI-aware deformable attentions, the cross-attention layer in MSTR decoder extracts three different semantics (human, object, and contextual information) for each HOI query from the encoded image features. For each decoder layer, we discovered that compositing the multiple semantics obtained from the previous

cross-attention layer [5] by summing the semantics after applying individual self-attention demonstrates the best performance (see Table 3 and Appendix). The input for the $(k + 1)$ -th layer of our HOI-aware deformable attention \bar{z}_q^{k+1} is written as:

$$\bar{z}_q^{k+1} = \text{SA}(f_q^h(k)) + \text{SA}(f_q^o(k)) + \text{SA}(f_q^c(k)), \quad (6)$$

where $f_q^h(k)$, $f_q^o(k)$, $f_q^c(k)$ denotes the multiple semantic outputs of the previous (k) -th decoder obtained by Eq.(6) and Eq.(5), respectively. SA denotes Multi-Head Self-Attention operation with Eq.(1) [28] and $\bar{z}_q^1 = \text{SA}(z_q) + \text{SA}(z_q^h) + \text{SA}(z_q^o)$.

MSTR Inference. Given the cross attention results of the final decoder layer where f_q^h and f_q^o is obtained by Eq. (6) and f_q^c is obtained by Eq. (5), the final prediction heads in MSTR predict the $\langle \text{bbox}_q^h, \text{bbox}_q^o, \text{cls}_q^o, \text{act}_q \rangle$ using FFN as follows:

$$(u_{qx}, u_{qy}, u_{qw}, u_{qh}) = \text{FFN}_{\text{hbox}}(f_q^h), \quad (7)$$

$$(v_{qx}, v_{qy}, v_{qw}, v_{qh}) = \text{FFN}_{\text{obox}}(f_q^o), \quad (8)$$

$$\text{cls}_q = \sigma(\text{FFN}_{\text{cls}}(f_q^c)), \quad (9)$$

$$\text{act}_q = \sigma(\text{FFN}_{\text{act}}(f_q^c)), \quad (10)$$

where cls_q and act_q each denote predictions for object the class and the action class after sigmoid function, and final bbox_q^h is predicted with the center point $(\sigma(u_{qx} + \sigma^{-1}(h_{qx})), \sigma(u_{qy} + \sigma^{-1}(h_{qy})))$, width u_{qw} , and height u_{qh} . Likewise, the bbox_q^o is predicted with center point as $(\sigma(v_{qx} + \sigma^{-1}(o_{qx})), \sigma(v_{qy} + \sigma^{-1}(o_{qy})))$, width v_{qw} , and height v_{qh} . σ and σ^{-1} denote the sigmoid and inverse sigmoid function, respectively, and is used to normalize the reference points h_q, o_q and the predicted coordinates of human boxes and object boxes $u_q\{x,y,w,h\}, v_q\{x,y,w,h\} \in \mathbb{R}$.

4. Experiment

In this section, we show the experimental results of our model in HOI detection. We first describe the experimental settings such as datasets and evaluation metrics. Next, we compare MSTR with state-of-the-art works on two different benchmarks (V-COCO and HICO-DET) and provide detailed ablation study for each component. Through the experiments, we demonstrate that MSTR successfully extends conventional transformer-based HOI detectors to utilize multi-scale feature maps, and emphasize that each component of MSTR contributes to the final HOI detection performance. Lastly, we provide extensive qualitative results of MSTR.

4.1. Datasets and Metrics

We evaluate our model on two widely-used public benchmarks: the V-COCO (*Verbs in COCO*) [9] and HICO-DET [3] datasets. V-COCO is a subset of COCO composed of 5,400 trainval images and 4,946 test images. For V-COCO dataset, we report the AP_{role} over 25 interactions in two scenarios. HICO-DET contains 37,536 and 9,515 images for each training and test split with annotations for 600 $\langle \text{verb}, \text{object} \rangle$ interaction types. We follow the previous settings and report the mAP over two evaluation settings (Default and Known Object), each with three different category sets: (1) all 600 HOI categories in HICO (Full), (2) 138 HOI categories with less than 10 training instances (Rare), and (3) 462 HOI categories with 10 or more training instances (Non-Rare). See Appendix for details of the evaluation settings.

4.2. Quantitative Results

We use the standard evaluation code ¹ following the previous works [4, 14, 26, 39] to calculate metric scores for both V-COCO and HICO-DET.

Comparison to State-of-The-Art. We compare MSTR with state-of-the-art methods in Table 1 and Table 2. In Table 1, MSTR outperforms the previous state-of-the-art method in V-COCO dataset by a large margin (+3.2p in $AP_{\text{role}}^{\#1}$ and +4.2p in $AP_{\text{role}}^{\#2}$). Similar to this, in Table 2, MSTR achieves the highest mAP on HICO-DET dataset in all Full, Rare, and Non-Rare classes obtaining +2.1p, +3.46p, and +1.69p gain for each compared to the previous state-of-the-art. We use the same scoring function as QPIC without any modification for a fair comparison. Note that MSTR benefits from the advantages of using deformable attention: the fast convergence for training [38] (see more details and the convergence graph in our Appendix).

¹<https://github.com/YueLiao/PPDM>

Method	Backbone	$AP_{\text{role}}^{\#1}$	$AP_{\text{role}}^{\#2}$
<i>Models with external features</i>			
TIN (RP _D CD) [20]	R50	47.8	-
Verb Embedding [34]	R50	45.9	-
RPNN [37]	R50	-	47.5
PMFNet [29]	R50-FPN	52.0	-
PastaNet [19]	R50-FPN	51.0	57.5
PD-Net [35]	R50	52.0	-
ACP [15]	R152	53.0	-
FCMNet [22]	R50	53.1	-
ConsNet [23]	R50-FPN	53.2	-
<i>Sequential HOI Detectors</i>			
VSRL [9]	R50-FPN	31.8	-
InteractNet [8]	R50-FPN	40.0	48.0
BAR-CNN [16]	R50-FPN	43.6	-
GPNN [25]	R152	44.0	-
iCAN [7]	R50	45.3	52.4
TIN (RC _D) [20]	R50	43.2	-
DCA [31]	R50	47.3	-
VCL [12]	R50-FPN	48.3	-
DRG [6]	R50-FPN	51.0	-
VSGNet [27]	R152	51.8	57.0
IDN [18]	R50	53.3	60.3
<i>Parallel HOI Detectors</i>			
UnionDet [13]	R50-FPN	47.5	56.2
IPNet [32]	HG104	51.0	-
HOI Transformer [39] [†]	R101	52.9	-
ASNet [4] [†]	R50	53.9	-
GGNet [36]	HG104	54.7	-
HOTR [14] [†]	R50	55.2	64.4
QPIC [26] [†]	R50	58.8	61.0
<i>MSTR (Ours)</i>	R50	62.0	65.2

Table 1. Comparison of performance on V-COCO test set. $AP_{\text{role}}^{\#1}$, $AP_{\text{role}}^{\#2}$ denotes the performance under Scenario 1 and Scenario 2 in V-COCO, respectively. [†] denotes end-to-end HOI detectors with transformers, which are the main baselines for our work.

4.3. Ablation Study

We perform ablations to check the effects of our proposed Dual-Entity attention, Entity-conditioned Context attention, and our proposed decoder architecture that merges the self-attention of the multiple semantics.

Baselines. On basis of QPIC [26] structure, we define several variants for baselines by applying different combinations of sub-components from MSTR: *multi-scale feature maps (MS)*, *Deformable Attention (DA)*, *Dual-Entity attention (DE)*, and *Entity-conditioned Context attention (EC)*. Specifically, since deformable attention can be also applied to a single-scale feature map, *SS-Baseline* denotes QPIC where the attention in the transformer is replaced by DA. Our work can be seen as a process of improving the score

Method	Detector	Backbone	Feature	Default			Known Object		
				Full	Rare	Non Rare	Full	Rare	Non Rare
Sequential HOI Detectors									
Functional Gen. [1]	HICO-DET	R101	A+S+L	21.96	16.43	23.62	-	-	-
TIN [20]	HICO-DET	R50	A+S+P	22.90	14.97	25.26	-	-	-
VCL [12]	HICO-DET	R50	A+S	23.63	17.21	25.55	25.98	19.12	28.03
ConsNet [23]	HICO-DET	R50-FPN	A+S+L	24.39	17.10	26.56	30.34	23.40	32.41
DRG [6]	HICO-DET	R50-FPN	A+S	24.53	19.47	26.04	27.98	23.11	29.43
IDN [18]	HICO-DET	R50	A+S	24.58	20.33	25.86	27.89	23.64	29.16
Parallel HOI Detectors									
UnionDet [13]	HICO-DET	R50-FPN	A	17.58	11.72	19.33	19.76	14.68	21.27
PPDM [21]	HICO-DET	HG104	A	21.10	14.46	23.09	24.81	17.09	27.12
HOI Transformer [39] [†]	HICO-DET	R50	A	23.46	16.91	25.41	26.15	19.24	28.22
HOTR [14] [†]	HICO-DET	R50	A	25.10	17.34	27.42	-	-	-
GGNet [36]	HICO-DET	HG104	A	28.83	22.13	30.84	27.36	20.23	29.48
AS-Net [4] [†]	HICO-DET	R50	A	28.87	24.25	30.25	31.74	27.07	33.14
QPIC [26] [†]	HICO-DET	R50	A	29.07	21.85	31.23	31.68	24.14	33.93
MSTR (Ours)	HICO-DET	R50	A	31.17	25.31	32.92	34.02	28.83	35.57

Table 2. Performance comparison in HICO-DET. The Detector column is denoted as ‘HICO-DET’ to show that the object detector is fine-tuned on the HICO-DET training set. Each letter in Feature column stands for A: Appearance (Visual Features), S: Interaction Patterns (Spatial Correlations), P: Pose Estimation, L: Linguistic Priors, V: Volume. † denotes end-to-end HOI detectors with transformers. Note that all the baseline models without † are already based on multi-scale feature maps.

Method	MS	DA	DE	EC	mAP
(a) QPIC					29.07
(b) SS-Baseline		✓			25.53
(c) SS-Baseline + DE		✓	✓		27.06
(d) SS-Baseline + DE + EC		✓	✓	✓	27.70
(e) MS-Baseline	✓	✓			27.52
(f) MS-Baseline + DE	✓	✓	✓		28.30
(g) MS-Baseline + DE + EC	✓	✓	✓	✓	30.14
(h) MSTR (Ours)	✓	✓	✓	✓	31.17

Table 3. Comparison of MSTR with our baseline QPIC and its variants in the HICO-DET test set. SS and MS denote the models using single scale feature map and multi-scale feature maps, respectively. DE and EC indicate our proposed Dual-Entity attention and Entity-conditioned Context attention, respectively.

to the state-of-the-art by adapting *MS*, *DE*, *EC* step by step to *SS-Baseline*. *MS-Baseline+DE+EC* represents MSTR without merging with self-attention, instead simply passing the sum of the outputs to the next decoder layer.

HOI-Aware Deformable Attentions. In Table 3, we explore the effect of our proposed HOI-Aware Deformable Attentions: Dual-Entity attention and Entity-conditioned Context attention. As deformable attentions can also be applied in a single-scale feature map, we verify the effectiveness of our proposed deformable attentions on both single-scale and multi-scale baselines. As we described in our preliminary, the naïve implementation of deformable attention on

top of QPIC (for single-scale) significantly degrades the score in both single-scale and multi-scale environments (see (a vs. b) and (a vs. e)). The use of Dual-Entity attention (DE) consistently improves the score in both single-scale (+1.53p in (b vs. c)) and multi-scale environments (+0.78p in (e vs. f)). As well, Entity-conditioned Context attention (EC) contributes in the multi-scale environment when jointly used with DE (+0.64p in SS and +1.84p in MS). Therefore, we conclude that disentangling the references (DE) and conditionally reimbursing context information (EC) each gradually contributes to the final performance of HOI detection in both single-scale and multi-scale environments, enabling MSTR to effectively explore multi-scale feature maps to achieve state-of-the-art performance.

Single-scale vs. Multi-scale. In Table 1 and Table 2, we demonstrate that our method using the multi-scale feature maps outperform all previous methods, including transformer-based methods [4, 14, 26, 39] and the ones that already use multi-scale feature maps heavily [6, 12, 13, 18, 23, 36]. To analyze further, Table 3 compares single-scale version and the multi-scale version of our baselines (see (b-e) and (e-h)). In all cases of converting the single-scale feature map to the multi-scale one, we observe consistent performance gains (see (b vs. e), (c vs. f), and (d vs. g,h)). The gain is maximized when *DE* and *EC* are used together. We further provide a detailed analysis of the effectiveness of MSTR in multi-scale environments in our Appendix.

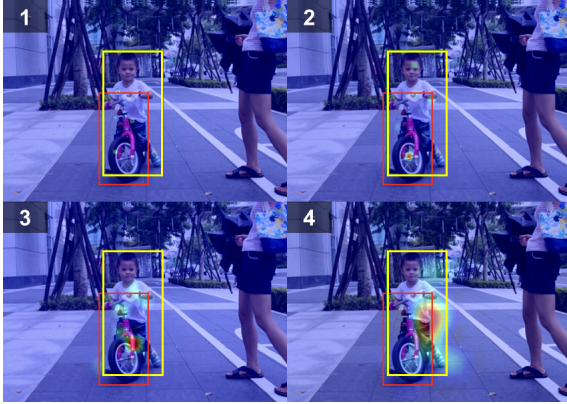


Figure 5. Visualization of our Entity-conditioned Context attentions on different levels of feature map (1 being the highest and 4 being the lowest resolution). Best viewed in color and scale.



Figure 6. Visualization of the HOI-aware attention of MSTR on different scales of humans and objects.

Decoder Architecture. We verify the effectiveness of Figure 4 (b) architecture in Table 3 (g vs. h). As MSTR considers multiple semantics with two suggested deformable modules, it is important to find suitable decoder architecture which can effectively merge the semantics [5]. According to the possible combination ways when merging three kinds of semantics, various types of decoder architecture can be candidates for the decoder architectures (described in Appendix). In our Appendix, we empirically verify that Figure 4 (b) architecture shows the most powerful and robust performance across all datasets.

4.4. Qualitative Results

We conduct qualitative analysis of MSTR to observe how MSTR captures interactions. Figure 1 and Figure 5 show the visualization of the attention map in MSTR in various feature levels. Interestingly, we can observe that

in the higher resolution feature maps, the sampling points capture the detail of the interacting human and object while the lower resolution feature maps tend to capture the overall pose or context of the interaction. In Figure 1 and Figure 6, we can observe how MSTR attends to test images that include various scales of humans, target objects, and distances. More details along with quantitative results will be provided in our Appendix.

5. Related Work

Transformer Based HOI Detectors. Human-Object Interaction detection has been initially proposed in [9], and has been developed in two main streams: sequential methods [1, 6–8, 10, 17, 19, 20, 22–25, 27, 29–31, 34, 35, 37] and parallel methods [13, 21, 32]. However, since these works required hand-crafted post-processing, HOI detectors with transformers have been proposed to eliminate the post-processing step through an end-to-end fashioned set prediction approach [4, 14, 26, 39]. Yet, all these methods are limited to a single-scale feature map due to the complexity caused when processing multi-scale feature maps with transformer attention.

Deformable Transformers for Object Detection. DETR has been recently proposed to eliminate the need for many hand-designed components in object detection [2]. Deformable DETR [38] mitigates the slow convergence and high complexity issues of DETR and successfully exploits multi-resolution feature maps. The deformable attention modules in [38] attend to a small set of sampling locations as a pre-filter for prominent key elements out of all the feature map pixels. However, unlike object detection, we observed that this pre-filter seriously deteriorates performance when applied to HOI detection. Therefore, in this paper, we focus on finding a proper way to incorporate deformable attention into HOI detection for exploiting multi-scale feature maps.

6. Conclusion

In this paper, we present MSTR, the first multi-scale approach in transformer-based HOI detectors. MSTR overcomes the issues of extending transformer-based HOI detectors to multi-scale feature maps with novel HOI-Aware Deformable attentions named as Dual-Entity attention and Entity-conditioned Context attention. In virtue of the two attention modules and our decoder architecture that effectively collects the multiple semantics from each of the attentions, MSTR achieves the state-of-the-art performance in two benchmark datasets in HOI detection.

References

- [1] Ankan Bansal, Sai Saketh Rambhatla, Abhinav Shrivastava, and Rama Chellappa. Detecting human-object interactions via functional generalization. In *AAAI*, pages 10460–10469, 2020. [1](#), [7](#), [8](#)
- [2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. *arXiv preprint arXiv:2005.12872*, 2020. [4](#), [8](#), [14](#)
- [3] 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. [6](#), [13](#)
- [4] 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. [1](#), [2](#), [6](#), [7](#), [8](#), [14](#)
- [5] Qi Dong, Zhuowen Tu, Haofu Liao, Yuting Zhang, Vijay Mahadevan, and Stefano Soatto. Visual relationship detection using part-and-sum transformers with composite queries. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3550–3559, 2021. [5](#), [8](#)
- [6] 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. [1](#), [6](#), [7](#), [8](#)
- [7] 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](#), [6](#), [8](#)
- [8] Georgia Gkioxari, Ross Girshick, Piotr Dollár, and Kaiming He. Detecting and recognizing human-object interactions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8359–8367, 2018. [1](#), [6](#), [8](#)
- [9] Jitendra Gupta, Saurabh Malik. Visual semantic role labeling. *arXiv preprint arXiv:1505.04474*, 2015. [1](#), [6](#), [8](#), [13](#)
- [10] Tanmay Gupta, Alexander Schwing, and Derek Hoiem. No-frills human-object interaction detection: Factorization, layout encodings, and training techniques. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9677–9685, 2019. [1](#), [8](#)
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [13](#)
- [12] Zhi Hou, Xiaojiang Peng, Yu Qiao, and Dacheng Tao. Visual compositional learning for human-object interaction detection. *arXiv preprint arXiv:2007.12407*, 2020. [6](#), [7](#)
- [13] Bumsoo Kim, Taeho Choi, Jaewoo Kang, and Hyunwoo Kim. Uniondet: Union-level detection towards real-time human-object interaction detection. In *Proceedings of the European conference on computer vision (ECCV)*, 2020. [1](#), [6](#), [7](#), [8](#)
- [14] 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. [1](#), [2](#), [6](#), [7](#), [8](#), [14](#)
- [15] Dong-Jin Kim, Xiao Sun, Jinsoo Choi, Stephen Lin, and In So Kweon. Detecting human-object interactions with action co-occurrence priors. *arXiv preprint arXiv:2007.08728*, 2020. [6](#)
- [16] Alexander Kolesnikov, Alina Kuznetsova, Christoph Lampert, and Vittorio Ferrari. Detecting visual relationships using box attention. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 0–0, 2019. [6](#)
- [17] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10166–10175, 2020. [1](#), [8](#)
- [18] 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, 2020. [6](#), [7](#)
- [19] Yong-Lu Li, Liang Xu, Xinpeng Liu, Xijie Huang, Yue Xu, Shiyi Wang, Hao-Shu Fang, Ze Ma, Mingyang Chen, and Cewu Lu. Pastanet: Toward human activity knowledge engine. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 382–391, 2020. [1](#), [6](#), [8](#)
- [20] 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 Conference on Computer Vision and Pattern Recognition*, pages 3585–3594, 2019. [1](#), [6](#), [7](#), [8](#)
- [21] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 482–490, 2020. [1](#), [4](#), [7](#), [8](#)
- [22] Y Liu, Q Chen, and A Zisserman. Amplifying key cues for human-object-interaction detection. *Lecture Notes in Computer Science*, 2020. [1](#), [6](#), [8](#)
- [23] 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. [1](#), [6](#), [7](#), [8](#)
- [24] Julia Peyre, Ivan Laptev, Cordelia Schmid, and Josef Sivic. Detecting unseen visual relations using analogies. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1981–1990, 2019. [1](#), [8](#)
- [25] 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. [1](#), [6](#), [8](#)
- [26] 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. 1, 2, 6, 7, 8, 11, 13, 14, 16
- [27] 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. 1, 6, 8
- [28] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. 5
- [29] Bo Wan, Desen Zhou, Yongfei Liu, Rongjie Li, and Xuming He. Pose-aware multi-level feature network for human object interaction detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9469–9478, 2019. 1, 6, 8
- [30] Hai Wang, Wei-shi Zheng, and Ling Yingbiao. Contextual heterogeneous graph network for human-object interaction detection. *arXiv preprint arXiv:2010.10001*, 2020. 1, 8
- [31] 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. *arXiv preprint arXiv:1910.07721*, 2019. 1, 6, 8
- [32] Tiancai Wang, Tong Yang, Martin Danelljan, Fahad Shahbaz Khan, Xiangyu Zhang, and Jian Sun. Learning human-object interaction detection using interaction points. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4116–4125, 2020. 1, 4, 6, 8
- [33] Tao Wang, Li Yuan, Yunpeng Chen, Jiashi Feng, and Shuicheng Yan. Pnp-detr: Towards efficient visual analysis with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4661–4670, 2021. 16
- [34] Bingjie Xu, Yongkang Wong, Junnan Li, Qi Zhao, and Mohan S Kankanhalli. Learning to detect human-object interactions with knowledge. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 1, 6, 8
- [35] Xubin Zhong, Changxing Ding, Xian Qu, and Dacheng Tao. Polysemy deciphering network for human-object interaction detection. In *Proc. Eur. Conf. Comput. Vis*, 2020. 1, 6, 8
- [36] Xubin Zhong, Xian Qu, Changxing Ding, and Dacheng Tao. Glance and gaze: Inferring action-aware points for one-stage human-object interaction detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13234–13243, 2021. 4, 6, 7
- [37] Penghao Zhou and Mingmin Chi. Relation parsing neural network for human-object interaction detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 843–851, 2019. 1, 6, 8
- [38] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. 2, 3, 5, 6, 8, 12, 13, 14
- [39] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11825–11834, 2021. 1, 2, 6, 7, 8, 14

A. Appendix

In this Appendix, we provide **i)** extended quantitative analysis of MSTR capturing HOI detection in a multi-scale environment, **ii)** exploration for various possible decoder architectures, **iii)** implementation details of MSTR, **iv)** details on experimental datasets and metrics, **v)** details of training, **vi)** analysis on convergence speed, **vii)** additional qualitative result on our Dual-Entity attention and Entity-conditioned Context attention, and finally, **viii)** limitations of our work.

A.1. Additional Quantitative Results for MSTR

First, we perform an extended quantitative analysis on the HICO-DET test set to validate the effectiveness of MSTR in a multi-scale environment. MSTR uses multi-scale feature maps to explore the semantics of HOI existing in different scales. In this section, we provide extensive quantitative results that shows the effectiveness of MSTR in capturing the interactions between humans and objects not only at different scales, but in various distances also (e.g., *adjacent* interaction such as ‘holding a book’ or *remote* interaction such as ‘throwing a frisbee’). To this end, we show quantitative results for multi-scale interactions according to 1) relative area of the human and the object, 2) the size of humans/objects, 3) distance between the human and object. For each criterion, we measure the performance across three bins where each bin has an equal and sufficient amount of HOI ground-truth labels to cover ($\sim 11,000$ HOIs). For comparison, we set QPIC [26], the state-of-the-art transformer-based approach that uses a single-scale feature map, as our baseline. Note that in this appendix, the size, area, and distance are all calculated in *normalized* image coordinates.

Relative area of human vs. object. To observe how MSTR handles interaction between humans and objects with different scales, we first calculate the average precision (AP) over interaction labels that have different relative areas of humans and objects ($\frac{\text{area}(\text{hbox})}{\text{area}(\text{oobj})}$). We cover three main cases according to their relative areas: i) $AP_{h<o}$ where the object area is significantly larger than the human area (e.g., human *sitting on a bench*), ii) $AP_{h=o}$ where the human and the object exists in comparable sizes, and iii) $AP_{h>o}$ where the object area is significantly smaller than the human area (e.g., human *throwing a ball*). We set the threshold for the relative areas so that each bin has an equal number of ground-truth instances (i.e., $\frac{\text{area}(\text{hbox})}{\text{area}(\text{oobj})} < 0.48$ for $AP_{h<o}$ and $\frac{\text{area}(\text{hbox})}{\text{area}(\text{oobj})} > 4.33$ for $AP_{h>o}$). In Table 4, MSTR outperforms QPIC in all three types of interaction categories. Note that the improvement is more substantial in cases where the human and object have vastly different scales (+3.01p for $AP_{h<o}$ and +1.85p for $AP_{h>o}$), verifying that MSTR is ef-

Method	$AP_{h<o}$	$AP_{h=o}$	$AP_{h>o}$
QPIC	34.10	30.57	25.22
MSTR	37.11	31.68	27.07
ΔAP	+3.01	+1.11	+1.85

Table 4. Comparison of MSTR with QPIC under interactions with different human/object scale ratio.

fectively utilizing multi-scale feature maps.

Human & object size. Here, we compare the average precision over the sizes of humans and objects. AP_L , AP_M , AP_S each denotes the average precision for **L**arge, **M**iddle, and **S**mall humans and objects. In Table 5, MSTR outperforms QPIC in all three categories in both human and object scales. For the human scales, the improvement is more recognizable in interactions including small human areas (+3.06p in AP_S) while for object scales, the improvement is consistent over all three scales.

Method	Human Size			Object Size		
	AP_L	AP_M	AP_S	AP_L	AP_M	AP_S
QPIC	28.65	35.36	24.14	33.09	28.65	24.87
MSTR	30.04	37.02	27.20	34.87	30.48	26.60
ΔAP	+1.39	+1.66	+3.06	+1.78	+1.83	+1.73

Table 5. Comparison of MSTR with QPIC under different sizes of humans and objects.

Interactions in various distances. Not only does MSTR capture interactions with various sized participants, but MSTR also captures interactions with various sized contexts, i.e., interaction in various distances. To correctly measure how *remote* an interaction is, we note that the distance between center points [26] should be normalized by both the image size and the size of the human and object box participating in the interaction. Given the interaction between hbox (hx_1, hy_1, hx_2, hy_2) and obox (ox_1, oy_1, ox_2, oy_2), the normalized box area as area (hbox) and area (obox), we define the distance $d_{\text{interaction}}$ as

$$d_{\text{center}} = \sqrt{\left(\frac{hx_1+hx_2}{2} - \frac{ox_1+ox_2}{2}\right)^2 + \left(\frac{hy_1+hy_2}{2} - \frac{oy_1+oy_2}{2}\right)^2}, \quad (11)$$

$$d_{\text{interaction}} = d_{\text{center}} / (\text{area}(\text{hbox}) \cdot \text{area}(\text{oobj})).$$

Then, we measure the average precision over three categories: i) AP_{adjacent} where the human is interacting with a nearby object, ii) AP_{distant} where the interacting human/object is within moderate distance, and AP_{remote} where the human is interacting with an object sufficiently far away. As in previous sections, we set the distance threshold so that

Method	AP _{adjacent}	AP _{distant}	AP _{remote}
QPIC	31.09	31.25	21.81
MSTR	32.66	33.48	23.70
Δ AP	+1.57	+2.23	+1.89

Table 6. Comparison of MSTR with QPIC under interactions with various distances.

each bin has an equal number of ground-truth instances. Table 6 shows the improvement of MSTR over QPIC. Note that while MSTR shows improvement across all three categories, the improvement is more distinguishable in cases where humans are interacting with objects in considerable distance (+2.23p for AP_{distant} and +1.89p for AP_{remote}, respectively).

A.2. Analysis on Decoder Architecture

As MSTR considers multiple semantics with two suggested deformable modules (Dual-Entity attention and Entity-conditioned Context attention), it is important to find a suitable decoder architecture that effectively merges the semantics. Here, we explore the possible combinations and various types of decoder architecture candidates when merging the three kinds of semantics. We empirically verify that MSTR architecture shows the most powerful performance.

Architecture for Dual-Entity attention. In Figure 7, we explore different architectures for Dual-Entity attention. We start with the most basic form: (a) is the architecture of QPIC, and (b) shows a straightforward application of the deformable attention [38] to QPIC. However, as we discussed in our main paper, (b) degrades the performance a lot from (a), because unlike its counterpart in object detection, multiple localizations need to be entangled to a single reference point in architecture (b). Therefore, we first use Dual-Entity attention to disentangle sampling points and attention weights for the participating entities (*i.e.*, human and object), respectively, to improve HOI detection performance. In Figure 7, (c) and (d) shows two options of dealing with the dual semantics obtained from dual reference points (each for humans and objects). In (c), each reference point is dealt with a separate stack of decoder layers (*i.e.*, Double-stream), while in (d) they are handled within a single-stream by sharing the self-attention layer where the input is simply the sum of the multiple semantics from the previous decoder layer. In Table 9, we show that our Dual-Entity attention shows a valid improvement (see (d) vs. (b)), while it even shows better performance than (c) requiring twice the number of decoder parameters.

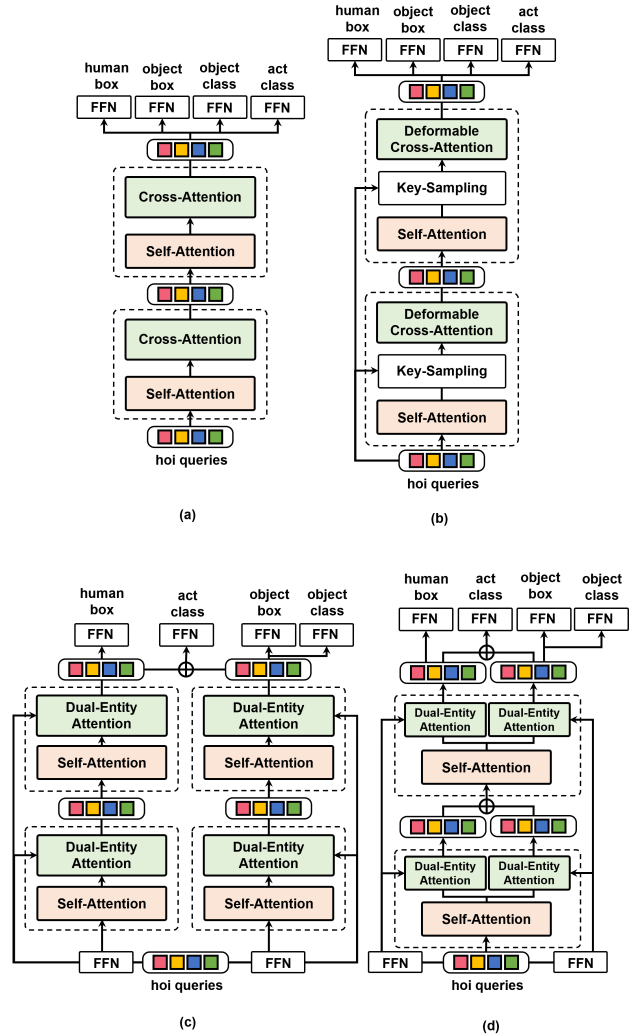


Figure 7. Comparison of a simple 2-layer Decoder architecture for: (a) QPIC, and (b) Direct application of Deformable DETR on QPIC, (c) Dual-Entity attention with two streams of decoder layers and (d) Dual-Entity attention that shares the self-attention layer.

Method	Default (Full)
(a) QPIC	29.07
(b) QPIC + Deformable attention [38]	27.52
(c) Double-stream	28.15
(d) Dual-Entity attention	28.30

Table 7. Comparison of Dual-Entity attention performance (d) against architecture in Figure 7 (a-c).

Modeling Conditional Context attention. In HOI detection, contextual information often gives an important clue in identifying interactions. In Table 8, we study the two different methods of obtaining context attention using (a) stan-

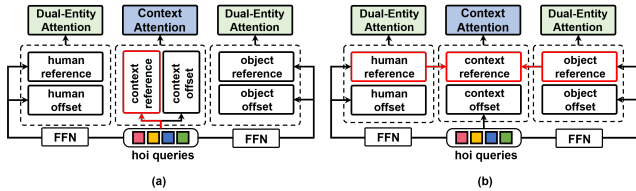


Figure 8. Comparison of: (a) context sampling with deformable attention, and (b) Entity-conditioned Context attention.

standard deformable attention and (b) our Entity-conditioned Context attention; note that in standard deformable attention, context reference points are directly obtained from HOI queries with a linear projection while our method conditionally obtain it from human and object reference points (see Figure 8). It can be observed that despite its simple structure and minimal delay, our Entity-conditioned Context attention achieves an +0.78p improvement compared to its counterpart. This implies that the guidance by human and object points is important to effectively model contextual information.

Method	Default (Full)
(a) Standard Deformable attention	29.36
(b) Entity-conditioned Context attention	30.14

Table 8. Comparison of the performance of Entity-conditioned Context attention against standard deformable attention [38]. Both (a) and (b) leverage Dual-Entity attention and follow the architectural design of Figure 9 (a) for fair comparison.

Merging the semantics. Figure 9 shows two different ways of how to merge the three semantics obtained from our Dual-Entity attention and Entity-conditioned Context attention. In MSTR, we merge the multiple semantics after applying self-attention separately to each of the semantic features obtained in the previous layer (Figure 9 (b)) instead of forcedly composing the input features of the self-attention layer (Figure 9 (a)). Table 9 shows that MSTR architecture (b) outperforms (a) by a margin of +1.03p, achieving the final performance. Note that while (b) is better, MSTR outperforms competing algorithms (presented in Table 2 of main paper) even with architecture (a).

Method	Default (Full)
(a) Merge self-attention input	30.14
(b) Merge self-attention output	31.17

Table 9. Comparison of a simple 2-layer Decoder architecture for Transformer-based HOI detectors: (a) Merging the input of the self-attention, and (b) architecture of MSTR (merging the output of self-attention).

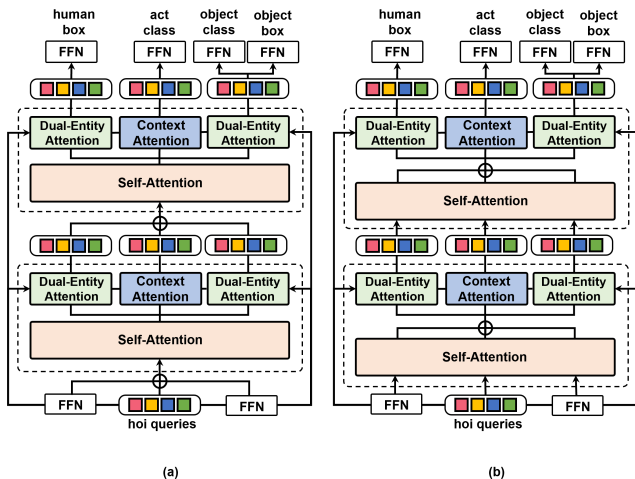


Figure 9. Comparison of a simple 2-layer Decoder architecture for Transformer-based HOI detectors: (a) Merging the input of the self-attention, and (b) architecture of MSTR (merging the output of self-attention).

A.3. Implementation Details

Following implementation details in Deformable DETR [38], we use ImageNet pre-trained ResNet-50 [11] as our backbone CNN and extract multi-scale feature maps without FPN. The number of attention heads and sampling offsets for deformable attentions are set to $M = 8$ and $K = 4$, respectively. The AdamW optimizer is used with the initial learning rate of $2e-4$ and weight decay of $1e-4$. All transformer weights are initialized with weights pre-trained in MS-COCO. For a fair comparison with QPIC [26], we use only 100 HOI queries instead of using 300 ones as in Deformable DETR [38].

A.4. Details on Datasets and Metrics

We evaluate our model on two widely-used public benchmarks: the V-COCO (*Verbs in COCO*) [9] and HICO-DET [3] datasets. V-COCO is a subset of COCO composed of 5,400 trainval images and 4,946 test images. For V-COCO dataset, we report the AP_{role} over 25 interactions in two scenarios. In Scenario 1 (denoted as $AP_{role}^{\#1}$), detectors should predict an output indicating the non-existence of an object $([0,0,0,0])$ when the target object is occluded, while in Scenario 2 (denoted as $AP_{role}^{\#2}$), only the localization of human and interaction classification is scored for such cases. HICO-DET contains 37,536 and 9,515 images for each training and test splits with annotations for 600 $\langle verb, object \rangle$ interaction types. In HICO-DET dataset, there are two different evaluation settings: *Default* and *Known object*. The former measures AP on all the test images, while the latter only considers the images with the object class corresponding to each AP. We report our score

with both settings. Note that the *Default* is a more challenging setting as we also need to distinguish background images. We follow the previous settings and report the mAP over three different category sets: (1) all 600 HOI categories in HICO (Full), (2) 138 HOI categories with less than 10 training instances (Rare), and (3) 462 HOI categories with 10 or more training instances (Non-Rare).

A.5. Training Details of MSTR

In this section, we explain the details of MSTR training. MSTR follows a set prediction approach as in previous transformer-based HOI detectors [4, 14, 26, 39]. We first introduce the cost matrix of Hungarian Matching for unique matching between the ground-truth HOI triplets and HOI set predictions.

Hungarian Matching for HOI Detection. MSTR predicts a fixed number K of HOI triplets that consist of a human box, object box, and binary classification for the a types of actions (where $a=25$ in V-COCO and 117 for HICO-DET). Each prediction captures a unique \langle human,object \rangle pair with multiple interactions. K is set to be larger than the typical number of interacting pairs in an image (in our experiment, $K = 100$). Let \mathcal{Y} denote the set of ground truth HOI triplets and $\hat{\mathcal{Y}} = \{\hat{y}_i\}_{i=1}^K$ as the set of K predictions. As K is larger than the number of unique interacting pairs in the image, we consider \mathcal{Y} also as a set of size K padded with \emptyset (there are no ground-truth that matches the prediction). Let $y = (b^h, b^o, c^o, a)$ where the ground-truth interaction y_i consists of b_i^h and b_i^o which denotes the normalized coordinates for the interacting human/object box, c_i^o denotes the target object class, and a_i denotes the one-hot for multiple actions. To find a bipartite matching between these two sets we search for a permutation of K elements $\sigma \in \mathfrak{S}_K$ with the lowest cost:

$$\hat{\sigma} = \operatorname{argmin}_{\sigma \in \mathfrak{S}_K} \sum_i^K \mathcal{C}_{\text{match}}(y_i, \hat{y}_{\sigma(i)}), \quad (12)$$

where $\mathcal{C}_{\text{match}}$ is a pair-wise *matching cost* between ground truth y_i and a prediction with index $\sigma(i)$. Now, the ground-truth is written as $y_i = (b_i^h, b_i^o, c_i^o, a_i)$ and the prediction is written as $\hat{y}_{\sigma(i)} = (\hat{b}_{\sigma(i)}^h, \hat{b}_{\sigma(i)}^o, \hat{c}_{\sigma(i)}^o, \hat{a}_{\sigma(i)})$ where $\hat{y}_{\sigma(i)}$ is the prediction that has the minimal matching cost with y_i . $\hat{b}_{\sigma(i)}^h$ and $\hat{b}_{\sigma(i)}^o$ are the normalized box coordinates for humans and objects, respectively, $\hat{c}_{\sigma(i)}^o$ is the classification for the target object of the interaction, and $\hat{a}_{\sigma(i)}$ is the predicted actions.

Final Cost/Loss function for MSTR. Based on $\mathcal{C}_{\text{match}}$, we calculate the final loss function for all pairs matched. The cost/loss function for the HOI triplets consists of the localization loss, object classification loss, and the action

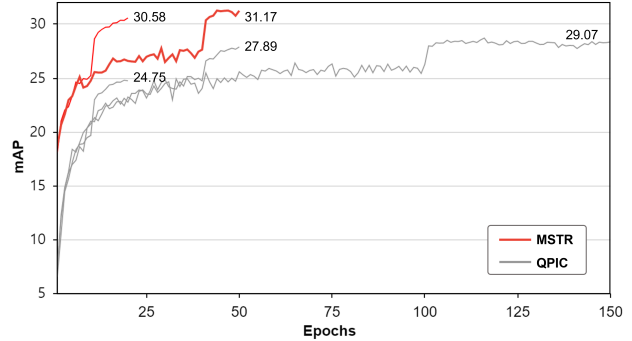


Figure 10. Comparison of convergence curves of QPIC and MSTR in the HICO-DET dataset. MSTR shows faster convergence than QPIC under various training schedules for both methods.

classification loss as $\mathcal{L}_H = \mathcal{L}_{\text{loc}} + \mathcal{L}_{\text{cls}} + \mathcal{L}_{\text{act}}$ where each function is written as

$$\begin{aligned} \mathcal{L}_{\text{loc}} &= \sum_{i=1}^K [\mathcal{L}_{\text{loc}}(b_i^h, \hat{b}_{\sigma(i)}^h) + \mathcal{L}_{\text{loc}}(b_i^o, \hat{b}_{\sigma(i)}^o)], \\ \mathcal{L}_{\text{cls}} &= \sum_{i=1}^K \text{BCELoss}(c_i, \hat{c}_{\sigma(i)}), \\ \mathcal{L}_{\text{act}} &= \sum_{i=1}^K \text{BCELoss}(a_i, \hat{a}_{\sigma(i)}). \end{aligned} \quad (13)$$

Identical to previous works [2, 4, 14, 26, 38, 39], the localization loss is defined by the weighted sum of the L1-loss and the gIoU loss.

A.6. Convergence speed

One of the advantages that deformable attention provides is the fast convergence at training. Figure 10 shows the convergence curve of MSTR compared to QPIC. Specifically, MSTR requires a much short number of epochs (50 epochs) compared to QPIC (150 epochs) to reach its best score. Note that MSTR achieves a competitive score to QPIC only with 20 epochs, outperforming QPIC with approximately $\times 4$ shorter training time.

A.7. Qualitative Analysis for MSTR

In this section, we conduct extensive qualitative analysis of MSTR to observe how Dual-Entity attention and the Entity-conditioned Context attention capture different semantics for interactions in a multi-scale environment.

MSTR attentions on multi-scale feature maps. We conduct a qualitative analysis of MSTR on both Dual-Entity attention and the Entity-conditioned Context attention in HOI

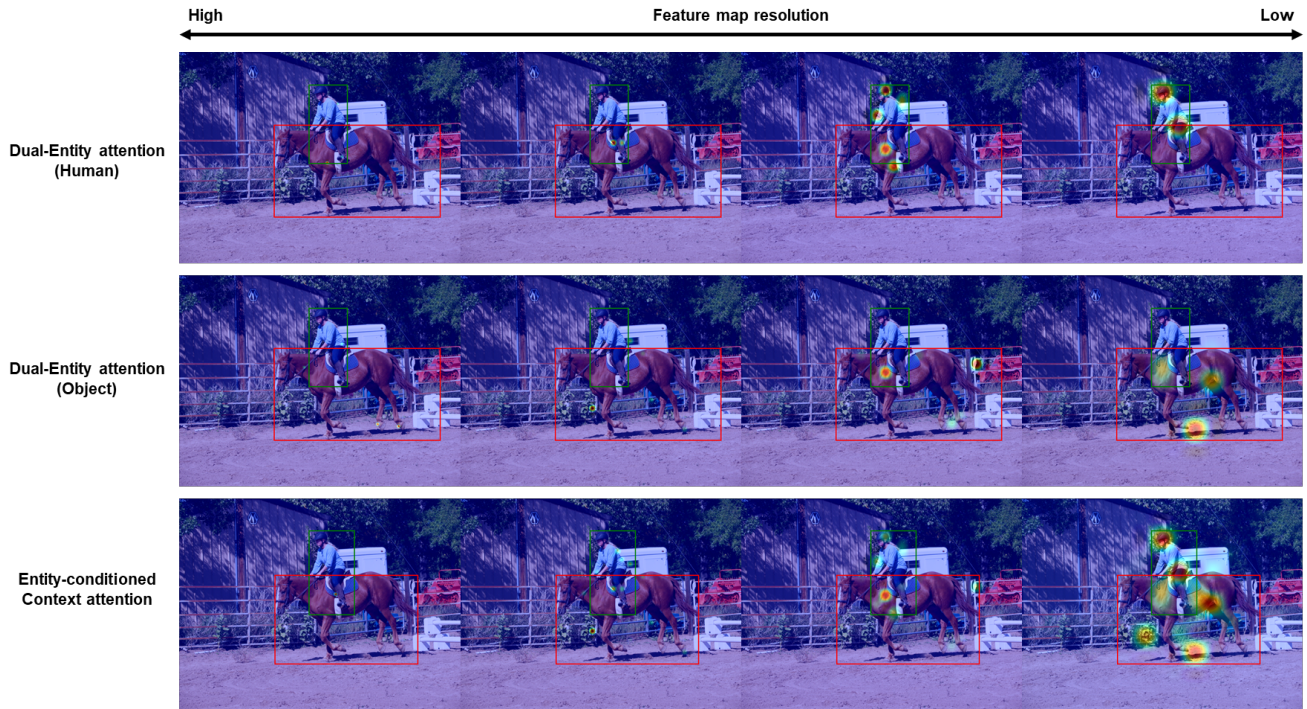


Figure 11. Visualization of the attention for the Dual-Entity attention and Entity-conditioned Context attention of MSTR in multi-scale feature maps for *adjacent* interaction: *ride*.

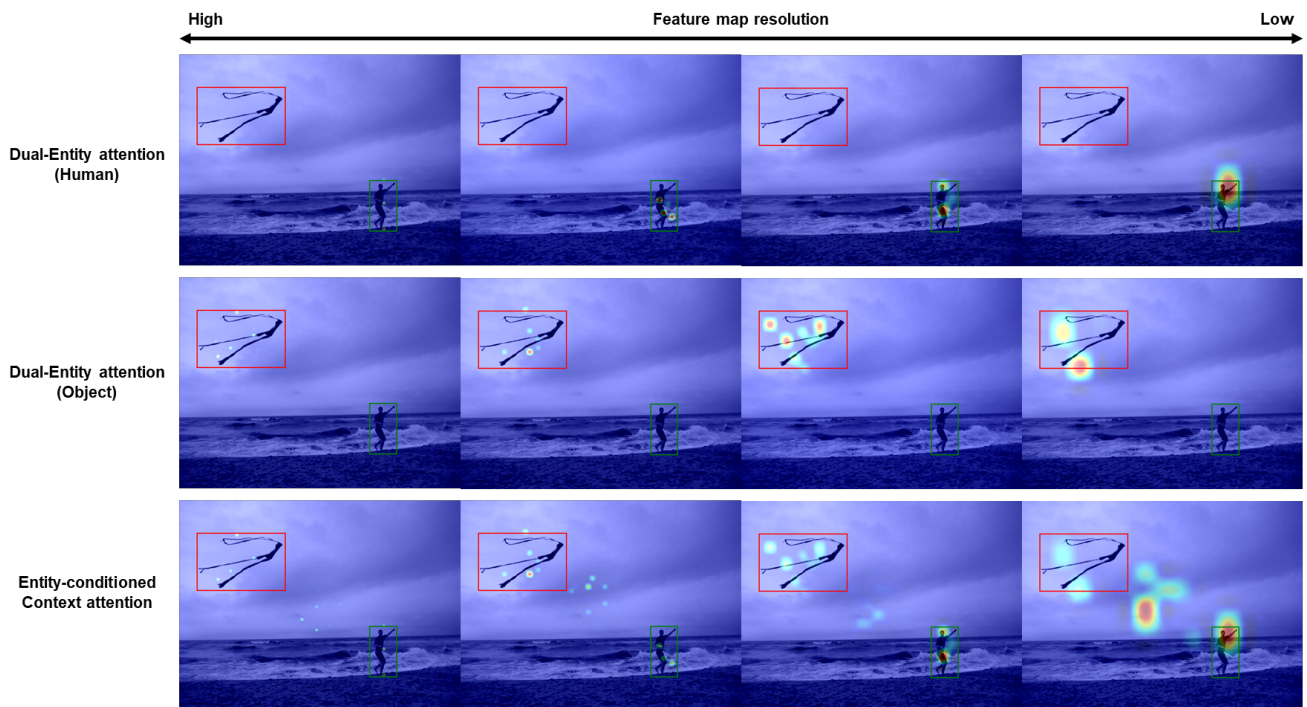


Figure 12. Visualization of the attention for the Dual-Entity attention and Entity-conditioned Context attention of MSTR in multi-scale feature maps for *remote* interaction: *fly*. It can be seen that in both *adjacent* interaction and *remote* interaction, MSTR successfully captures the multiple semantics of the human, object, and contextual information across the multi-resolution feature maps.



Figure 13. MSTR attentions (Dual-Entity attention and Entity-conditioned Context attention) of different scales all visualized at once.

detection to observe how MSTR captures interactions. Figure 11 shows the visualization of each attention in an *adjacent* interaction: *ride*. Figure 12 shows the visualization of each attention in an *remote* interaction: *fly*. For both cases, we can see that the Dual-Entity attention captures the appearance of the human and object across multiple scales of feature maps. In contrast, the Entity-conditioned Context attention tends to capture an inclusive area that covers both two regions and their intermediate background, effectively capturing the context of the interaction.

MSTR attentions on multi-scale feature maps. In Figure 13, we provide more qualitative visualizations for the multi-scale attentions of MSTR in various scenes with 1) large human and small object, 2) small human and large object, 3) distant interactions, 4) adjacent interactions.

A.8. Limitations

The main limitation of our work is the bottleneck caused by the extensive size of the *query* element (multi-scale image features, there are about $\times 20$ more image tokens to process compared to the single-scale feature map). Despite our proposed deformable attentions, MSTR suffers from an estimated 10% increase in parameters and $\sim \times 2$ GFLOPs compared to the single-scale baseline, QPIC [26]. Although recent related works have tackled the efficiency problem in deformable attentions by sampling the query element as well [33], the research scope of this work did not cover this issue.