HOTR: End-to-End Human-Object Interaction Detection with Transformers

Bumsoo Kim^{1,2} Junhyun Lee² Jaewoo Kang² Eun-Sol Kim^{1,†} Hyunwoo J. Kim^{2,†}

¹Kakao Brain ²Korea University

{bumsoo.brain, eunsol.kim}@kakaobrain.com {meliketoy, ljhyun33, kangj, hyunwoojkim}@korea.ac.kr

Abstract

Human-Object Interaction (HOI) detection is a task of identifying "a set of interactions" in an image, which involves the i) localization of the subject (i.e., humans) and target (i.e., objects) of interaction, and ii) the classification of the interaction labels. Most existing methods have indirectly addressed this task by detecting human and object instances and individually inferring every pair of the detected instances. In this paper, we present a novel framework, referred by HOTR, which directly predicts a set of (human, object, interaction) triplets from an image based on a transformer encoder-decoder architecture. Through the set prediction, our method effectively exploits the inherent semantic relationships in an image and does not require time-consuming post-processing which is the main bottleneck of existing methods. Our proposed algorithm achieves the state-of-the-art performance in two HOI detection benchmarks with an inference time under 1 ms after object detection.

1. Introduction

Human-Object Interaction (HOI) detection has been formally defined in [8] as the task to predict a set of \(\text{human}, \) object, interaction\(\text{)} triplets within an image. Previous methods have addressed this task in an indirect manner by performing object detection first and associating \(\text{human}, \) object\(\text{)} pairs afterward with separate post-processing steps. Especially, early attempts \(i.e., \) sequential HOI detectors \([5, 18, 17, 26] \)) have performed this association with a subsequent neural network, thus being time-consuming and computationally expensive.

To overcome the redundant inference structure of sequential HOI detectors, recent researches [30, 19, 12] proposed parallel HOI detectors. These works explicitly localize interactions with either interaction boxes (*i.e.*, the tightest box that covers both the center point of an object

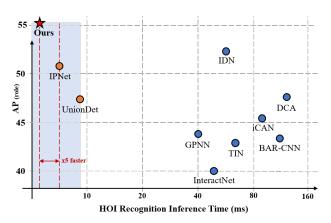


Figure 1. Time vs. Performance analysis for HOI detectors on V-COCO dataset. HOI recognition inference time is measured by subtracting the object detection time from the end-to-end inference time. Blue circle represents sequential HOI detectors, orange circle represents parallel HOI detectors and red star represents ours. Our method achieves an HOI recognition inference time of 0.9ms, being significantly faster than the parallel HOI detectors such as IPNet [30] or UnionDet [12] (the comparison between parallel HOI detectors is highlighted in blue).

pair) [30, 19] or union boxes (*i.e.*, the tightest box that covers both the box regions of an object pair) [12]. The localized interactions are associated with object detection results to complete the (human, object, interaction) triplet. The time-consuming neural network inference is replaced with a simple matching based on heuristics such as distance [30, 19] or IoU [12].

However, previous works in HOI detection are still limited in two aspects; i) They require additional post-processing steps like suppressing near-duplicate predictions and heuristic thresholding. ii) Although it has been shown that modeling relations between objects helps object detection [11, 2], the effectiveness of considering high-level dependency for interactions in HOI detection has not yet been fully explored.

In this paper, we propose a fast and accurate HOI algorithm named HOTR (Human-Object interaction TRans-

[†]corresponding authors

former) that predicts a set of human-object interactions in a scene at once with a direct set prediction approach. We design an encoder-decoder architecture based on transformers to predict a set of HOI triplets, which enables the model to overcome both limitations of previous works. First, direct set-level prediction enables us to eliminate hand-crafted post-processing stage. Our model is trained in an end-to-end fashion with a set loss function that matches the predicted interactions with ground-truth (human, object, interaction) triplets. Second, the self-attention mechanisms of transformers makes the model exploit the contextual relationships between human and object and their interactions, encouraging our set-level prediction framework more suitable for high-level scene understanding.

We evaluate our model in two HOI detection benchmarks: V-COCO and HICO-DET datasets. Our proposed architecture achieves state-of-the-art performance on two datasets compared to both sequential and parallel HOI detectors. Also, note that our method is much faster than other algorithms as illustrated in Figure 1, by eliminating time-consuming post-processing through the direct set-level prediction. The contribution of this work can be summarized as the following:

- We propose HOTR, the first transformer-based set prediction approach in HOI detection. HOTR eliminates the hand-crafted post-processing stage of previous HOI detectors while being able to model the correlations between interactions.
- We propose various training and inference techniques for HOTR: HO Pointers to associate the outputs of two parallel decoders, a recomposition step to predict a set of final HOI triplets, and a new loss function to enable end-to-end training.
- HOTR achieves state-of-the-art performance on both benchmark datasets in HOI detection with an inference time under 1 ms, being significantly faster than previous parallel HOI detectors (5~9 ms).

2. Related Work

2.1. Human-Object Interaction Detection

Human-Object Interaction detection has been initially proposed in [8], and has been developed in two main streams: sequential methods and parallel methods. In sequential methods, object detection is performed first and every pair of the detected object is inferred with a separate neural network to predict interactions. Parallel HOI detectors perform object detection and interaction prediction in parallel and associates them with simple heuristics such as distance or IoU.

Sequential HOI Detectors: InteractNet [6] extended an existing object detector by introducing an action-specific density map to localize target objects based on the human-centric appearance, and combined features from individual boxes to predict the interaction. Note that interaction detection based on visual cues from individual boxes often suffers from the lack of contextual information.

To this end, iCAN [5] proposed an instance-centric attention module that extracts contextual features complementary to the features from the localized objects/humans. No-Frills HOI detection [9] propose a training and inference HOI detection pipeline only using simple multi-layer perceptron. Graph-based approaches have proposed frameworks that can explicitly represent HOI structures with graphs [24, 26, 4, 28, 21]. Deep Contextual Attention [29] leverages contextual information by a contextual attention framework in HOI. [28] proposes a heterogeneous graph network that models humans and objects as different kinds Various external sources such as linguistic of nodes. priors [23, 31, 17, 4, 1, 32, 20] or human pose information [15, 33, 18, 9, 27, 33] have also been leveraged for further improve performance. Although sequential HOI detectors feature a fairly intuitive pipeline and solid performance, they are time-consuming and computationally expensive because of the additional neural network inference after the object detection phase.

Parallel HOI Detectors: Attempts for faster HOI detection has been also introduced in recent works as parallel HOI detectors. These works have directly localized interactions with interaction points [30, 19] or union boxes [12], replacing the separate neural network for interaction prediction with a simple heuristic based matching with distance or IoUs. Since they can be parallelized with existing object detectors, they feature fast inference time. However, these works are limited in that they require a hand-crafted postprocessing stage to associate the localized interactions with object detection results. This post-processing step i) requires manual search for the threshold, and ii) generates extra time complexity for matching each object pairs with the localized interactions ($5\sim9$ ms).

2.2. Object Detection with Transformers

DETR [2] has been recently proposed to eliminate the need for many hand-designed components in object detection while demonstrating good performance. DETR infers a fixed-size set of N predictions, in a single pass through the decoder, where N is set to be significantly larger than the typical number of objects in an image. The main loss for DETR produces an optimal bipartite matching between predicted and ground-truth objects. Afterward, the object-specific losses (for class and bounding box) are optimized.

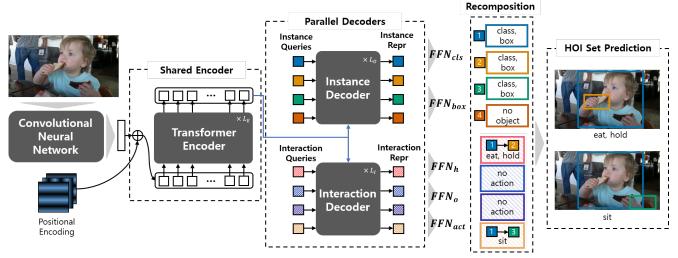


Figure 2. Overall pipeline of our proposed model. The Instance Decoder and Interaction Decoder run in parallel, and share the Encoder. In our recomposition, the interaction representations predicted by the Interaction Decoder are associated with the instance representations to predict a fixed set of HOI triplets (see Fig.3). The positional encoding is identical to [2].

3. Method

The goal of this paper is to predict a set of (human, object, interaction) triplets while considering the inherent semantic relationships between the triplets in an end-to-end manner. To achieve this goal, we formulate HOI detection as set prediction. In this section, we first discuss the problems of directly extending the set prediction architecture for object detection [2] to HOI detection. Then, we propose our architecture HOTR that parallelly predicts a set of object detection and associates the human and object of the interaction, while the self-attention in transformers models the relationships between the interactions. Finally, we present the details of training for our model including Hungarian Matching for HOI detection and our loss function.

3.1. Detection as Set Prediction

We first start from object detection as set prediction with transformers, then show how we extend this architecture to capture HOI detection with transformers.

Object Detection as Set Prediction. Object Detection has been explored as a set prediction problem by DETR [2]. Since object detection includes a single classification and a single localization for each object, the transformer encoder-decoder structure in DETR transforms N positional embeddings to a set of N predictions for the object class and bounding box.

HOI Detection as Set Prediction. Similar to object detection, HOI detection can be defined as a set prediction problem where each prediction includes the localization of a human region (*i.e.*, *subject* of the interaction), an

object region (*i.e.*, *target* of the interaction) and multi-label classification of the interaction types. One straightforward extension is to modify the MLP heads of DETR to transform each positional embedding to predict a human box, object box, and action classification. However, this architecture poses a problem where the localization for the same object needs to be redundantly predicted with multiple positional embeddings (*e.g.*, if the same person *works on a computer* while *sitting* on a chair, two different queries have to infer redundant regression for the same human).

3.2. HOTR architecture

The overall pipeline of HOTR is illustrated in Figure 2. Our architecture features a transformer encoder-decoder structure with a shared encoder and two parallel decoders (*i.e.*, instance decoder and interaction decoder). The results of the two decoders are associated with using our proposed *HO Pointers* to generate final HOI triplets. We will introduce HO Pointers shortly after discussing the architecture of HOTR.

Transformer Encoder-Decoder architecture. Similar to DETR [2], the global context is extracted from the input image by the backbone CNN and a shared encoder. Afterward, two sets of positional embeddings (*i.e.*, the instance queries and the interaction queries) are fed into the two parallel decoders (*i.e.*, the instance decoder and interaction decoder in Fig. 2). The instance decoder transforms the instance queries to instance representations for object detection while the interaction decoder transforms the interaction queries to interaction representations for

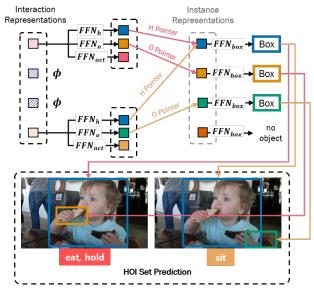


Figure 3. Conceptual illustration of how HO Pointers associates the interaction representations with instance representations. As instance representations are pre-trained to perform standard object detection, the interaction representation learns localization by predicting the *pointer* to the index of the instance representations for each human and object boxes. Note that the index pointer prediction is obtained in parallel with instance representations.

interaction detection. We apply feed-forward networks (FFNs) to the interaction representation and obtain a Human Pointer, an Object Pointer, and interaction type, see Fig. 3. In other words, the interaction representation localizes human and object regions by pointing the relevant instance representations using the Human Pointer and Object Pointer (HO Pointers), instead of directly regressing the bounding box. Our architecture has several advantages compared to the direct regression approach. We found that directly regressing the bounding box has a problem when an object participates in multiple interactions. In the direct regression approach, the localization of the identical object differs across interactions. Our architecture addresses this issue by having separate instance and interaction representations and associating them using HO Pointers. Also, our architecture allows learning the localization more efficiently without the need of learning the localization redundantly for every interaction. Note that our experiments show that our shared encoder is more effective to learn HO Pointers than two separate encoders.

HO Pointers. A conceptual overview of how HO Pointers associate the parallel predictions from the instance decoder and the interaction decoder is illustrated in Figure 3. HO Pointers (*i.e.*, Human Pointer and Object Pointer) contain the indices of the corresponding instance representations of the human and the object in the interaction. After the in-

teraction decoder transforms K interaction queries to K interaction representations, an interaction representation z_i is fed into two feed-forward networks $\mathrm{FFN}_h:\mathbb{R}^d\to\mathbb{R}^d$, $\mathrm{FFN}_o:\mathbb{R}^d\to\mathbb{R}^d$ to obtain vectors v_i^h and v_i^o , i.e., $v_i^h=\mathrm{FFN}_h(z_i)$ and $v_i^o=\mathrm{FFN}_o(z_i)$. Then finally the Human/Object Pointers \hat{c}_i^h and \hat{c}_i^o , which are the indices of the instance representations with the highest similarity scores, are obtained by

$$\hat{c}_{i}^{h} = \underset{j}{\operatorname{argmax}} \left(\operatorname{sim}(v_{i}^{h}, \mu_{j}) \right),$$

$$\hat{c}_{i}^{o} = \underset{j}{\operatorname{argmax}} \left(\operatorname{sim}(v_{i}^{o}, \mu_{j}) \right),$$
(1)

where μ_j is the *j*-th instance representation and $sim(u, v) = u^{\top}v/\|u\|\|v\|$.

Recomposition for HOI Set Prediction. From the previous steps, we now have the following: i) N instance representations μ , and ii) K interaction representations z and their HO Pointers \hat{c}^h and \hat{c}^o . Given γ interaction classes, our recomposition is to apply the feed-forward networks for bounding box regression and action classification as $\text{FFN}_{\text{box}}: \mathbb{R}^d \to \mathbb{R}^4$, and $\text{FFN}_{\text{act}}: \mathbb{R}^d \to \mathbb{R}^\gamma$, respectively. Then, the final HOI prediction for the i-th interaction representation z_i is obtained by,

$$\hat{b}_{i}^{h} = \operatorname{FFN}_{\text{box}}(\mu_{\hat{c}_{i}^{h}}) \in \mathbb{R}^{4},
\hat{b}_{i}^{o} = \operatorname{FFN}_{\text{box}}(\mu_{\hat{c}_{i}^{o}}) \in \mathbb{R}^{4},
\hat{a}_{i} = \operatorname{FFN}_{\text{act}}(z_{i}) \in \mathbb{R}^{\gamma}.$$
(2)

The final HOI prediction by our HOTR is the set of K triplets, $\{\langle \hat{b}_i^h, \hat{b}_i^o, \hat{a}_i \rangle\}_{i=1}^K$.

Complexity & Inference time. Previous parallel methods have substituted the costly pair-wise neural network inference with a fast matching of triplets (associating interaction regions with corresponding human regions and object regions based on distance [30] or IoU [12]). HOTR further reduces the inference time after object detection by associating K interactions with N instances, resulting in a smaller time complexity $\mathcal{O}(KN)$. By eliminating the post-processing stages in the previous one-stage HOI detectors including NMS for the interaction region and triplet matching, HOTR diminishes the inference time by $4 \sim 8 \text{ms}$ while showing improvement in performance.

3.3. Training HOTR

In this section, we explain the details of HOTR training. We first introduce the cost matrix of Hungarian Matching for unique matching between the ground-truth HOI triplets and HOI set predictions obtained by recomposition. Then, using the matching result, we define the loss for HO

Pointers and the final training loss.

Hungarian Matching for HOI Detection. HOTR predicts K HOI triplets that consist of human box, object box and binary classification for the a types of actions. Each prediction captures a unique $\langle \text{human,object} \rangle$ pair with one or more interactions. K is set to be larger than the typical number of interacting pairs in an image. We start with the basic cost function that defines an optimal bipartite matching between predicted and ground truth HOI triplets, and then show how we modify this matching cost for our interaction representations.

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 \varnothing (no interaction). 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:

$$\hat{\sigma} = \underset{\sigma \in \mathfrak{S}_K}{\operatorname{argmin}} \sum_{i}^{K} \mathcal{C}_{\text{match}}(y_i, \hat{y}_{\sigma(i)}), \tag{3}$$

where $\mathcal{C}_{\mathrm{match}}$ is a pair-wise *matching cost* between ground truth y_i and a prediction with index $\sigma(i)$. However, since y_i is in the form of $\langle \mathrm{hbox,obox,action} \rangle$ and $\hat{y}_{\sigma(i)}$ is in the form of $\langle \mathrm{hidx,oidx,action} \rangle$, we need to modify the cost function to compute the matching cost.

Let $\Phi: \mathrm{idx} \to \mathrm{box}$ be a mapping function from ground-truth $\langle \mathrm{hidx}, \mathrm{oidx} \rangle$ to ground-truth $\langle \mathrm{hbox}, \mathrm{obox} \rangle$ by optimal assignment for object detection. Using the inverse mapping $\Phi^{-1}: \mathrm{box} \to \mathrm{idx}$, we get the ground-truth idx from the ground-truth box.

Let $M \in \mathbb{R}^{d \times N}$ be a set of normalized instance representations $\mu' = \mu/\|\mu\| \in \mathbb{R}^d$, i.e., $M = [\mu'_1 \dots \mu'_N]$. We compute $\hat{P}^h \in \mathbb{R}^{K \times N}$ that is the set of softmax predictions for the H Pointer in (1) given as

$$\hat{P}^h = \|_{i=1}^K \operatorname{softmax}((\bar{v}_i^h)^T M), \tag{4}$$

where $\|_{i=1}^K$ denotes the vertical stack of row vectors and $\bar{v}_i^h = v_i^h/||v_i^h||$. \hat{P}^o is analogously defined.

Given the ground-truth $y_i=(b_i^h,b_i^o,a_i),\hat{P}^h,$ and $,\hat{P}^o,$ we convert the ground-truth box to indices by $c_i^h=\Phi^{-1}(b_i^h)$ and $c_i^o=\Phi^{-1}(b_i^o)$ and compute our matching cost function written as

$$C_{\text{match}}(y_i, \hat{y}_{\sigma(i)}) = -\alpha \cdot \mathbb{1}_{\{a_i \neq \varnothing\}} \hat{P}^h[\sigma(i), c_i^h]$$

$$-\beta \cdot \mathbb{1}_{\{a_i \neq \varnothing\}} \hat{P}^o[\sigma(i), c_i^o]$$

$$+ \mathbb{1}_{\{a_i \neq \varnothing\}} \mathcal{L}_{\text{act}}(a_i, \hat{a}_{\sigma(i)}),$$

$$(5)$$

where $\hat{P}[i,j]$ denotes the element at i-th row and j-th column, and $\hat{a}_{\sigma(i)}$ is the predicted action. The action matching

cost is calculated as $\mathcal{L}_{act}(a_i, \hat{a}_{\sigma(i)}) = BCELoss(a_i, \hat{a}_{\sigma(i)})$. α and β is set as a fixed number to balance the different scales of the cost function for index prediction and action classification.

Final Set Prediction Loss for HOTR. We then compute the Hungarian loss for all pairs matched above, where the loss for the HOI triplets has the localization loss and the action classification loss as

$$\mathcal{L}_{H} = \sum_{i=1}^{K} \left[\mathcal{L}_{loc}(\mathbf{c}_{i}^{h}, \mathbf{c}_{i}^{o}, z_{\sigma(i)}) + \mathcal{L}_{act}(a_{i}, \hat{a}_{\sigma(i)}) \right]. \tag{6}$$

The localization loss $\mathcal{L}_{loc}(c_i^h, c_i^o, z_{\sigma(i)})$ is denoted as

$$\mathcal{L}_{loc} = -\log \frac{\exp(\operatorname{sim}(\operatorname{FFN}_{h}(z_{\sigma(i)}), \mu_{c_{i}^{h}})/\tau)}{\sum_{k=1}^{N} \exp(\operatorname{sim}(\operatorname{FFN}_{h}(z_{\sigma(i)}), \mu_{k})/\tau)}$$

$$-\log \frac{\exp(\operatorname{sim}(\operatorname{FFN}_{o}(z_{\sigma(i)}), \mu_{c_{i}^{o}}/\tau)}{\sum_{k=1}^{N} \exp(\operatorname{sim}(\operatorname{FFN}_{o}(z_{\sigma(i)}), \mu_{k})/\tau)},$$
(7)

where τ is the temperature that controls the smoothness of the loss function. We empirically found that $\tau=0.1$ is the best value for our experiments.

Defining No-Interaction with HOTR. In DETR [2], maximizing the probability of the no-object class for the softmax output naturally suppresses the probability of other classes. However, in HOI detection the action classification is a multi-label classification where each action is treated as an individual binary classification. Due to the absence of an explicit class that can suppress the redundant predictions, HOTR ends up with multiple predictions for the same (human, object) pair. Therefore, HOTR sets an explicit class that learns the *interactiveness* (1 if there is *any* interaction between the pair, 0 otherwise), and suppresses the predictions for redundant pairs that have a low interactiveness score (defined as No-Interaction class). In our experiment in Table. 3, we show that setting an explicit class for interactiveness contributes to the final performance.

Implementation Details. We train HOTR with AdamW [22]. We set the transformer's initial learning rate to 10^{-4} and weight decay to 10^{-4} . All transformer weights are initialized with Xavier init [7]. For a fair evaluation with baselines, the Backbone, Encoder, and Instance Decoder are pre-trained in MS-COCO and frozen during training. We use the scale augmentation as in DETR [2], resizing the input images such that the shortest side is at least 480 and at most 800 pixels while the longest side at most is 1333.

4. Experiments

In this section, we demonstrate the effectiveness of our model in HOI detection. We first describe the two public datasets that we use as our benchmark: V-COCO and HICO-DET. Next, we show that HOTR successfully captures HOI triplets, by achieving state-of-the-art performance in both mAP and inference time. Then, we provide a detailed ablation study of the HOTR architecture.

4.1. Datasets

To validate the performance of our model, we evaluate our model on two public benchmark datasets: the V-COCO (Verbs in COCO) dataset and HICO-DET dataset. V-COCO is a subset of COCO and has 5,400 trainval images and 4,946 test images. For V-COCO dataset, we report the AP_{role} over 25 interactions in two scenarios AP $_{role}^{\#1}$ and $AP_{role}^{\#2}$. The two scenarios represent the different scoring ways for object occlusion cases. In Scenario1, the model should correctly predict the bounding box of the occluded object as [0,0,0,0] while predicting human bounding box and actions correctly. In Scenario2, the model does not need to predict about the occluded object. **HICO-DET** [3] is a subset of HICO dataset and has more than 150K annotated instances of human-object pairs in 47,051 images (37,536 training and 9,515 testing) and is annotated with 600 (verb, object) interaction types. For HICO-DET, we report our performance in the *Default* setting where we evaluate the detection on the full test set. 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).

4.2. Quantitative Analysis

For quantitative analysis, we use the official evaluation code for computing the performance of both V-COCO and HICO-DET. Table 1 and Table 2 show the comparison of HOTR with the latest HOI detectors including both sequential and parallel methods. For fair comparison, the instance detectors are fixed by the parameters pre-trained in MS-COCO. All results in V-COCO dataset are evaluated with the fixed detector. For the HICO-DET dataset, we provide both results using the fixed detector and the fine-tuned detector following the common evaluation protocol [1, 18, 10, 21, 4, 16, 12, 19].

Our HOTR achieves a new state-of-the-art performance on both V-COCO and HICO-DET datasets, while being the fastest parallel detector. Table 1 shows our result in the V-COCO dataset with both Scenario1 and Scenario2. HOTR outperforms the state-of-the-art parallel HOI detector [30] in Scenario1 with a margin of 4.2mAP.

Method	Backbone	$AP_{\text{role}}^{\#1}$	$AP_{\mathrm{role}}^{\#2}$				
Models with external features							
$\overline{\text{TIN}\left(\text{RP}_{\text{D}}\text{C}_{\text{D}}\right)\left[18\right]}$	R50	47.8					
Verb Embedding [31]	R50	45.9					
RPNN [33]	R50	-	47.5				
PMFNet [27]	R50-FPN	52.0					
PastaNet [17]	R50-FPN	51.0	57.5				
PD-Net [32]	R50	52.0	-				
ACP [13]	R152	53.0					
FCMNet [20]	R50	53.1	-				
ConsNet [21]	R50-FPN	53.2	-				
Sequential HOI Detectors							
VSRL [8]	R50-FPN	31.8	-				
InteractNet [6]	R50-FPN	40.0	48.0				
BAR-CNN [14]	R50-FPN	43.6	-				
GPNN [24]	R152	44.0	-				
iCAN [5]	R50	45.3	52.4				
$TIN (RC_D) [18]$	R50	43.2	-				
DCA [29]	R50	47.3	-				
VSGNet [26]	R152	51.8	57.0				
VCL [10]	R50-FPN	48.3					
DRG [4]	R50-FPN	51.0					
IDN [16]	R50	53.3	60.3				
Parallel HOI Detectors							
IPNet [30]	HG104	51.0	-				
UnionDet [12]	R50-FPN	47.5	56.2				
Ours	R50	55.2	64.4				

Table 1. Comparison of performance on V-COCO test set. $AP_{\rm role}^{\#1}$, $AP_{\rm role}^{\#2}$ denotes the performance under Scenario1 and Scenario2 in V-COCO, respectively.

Table 2 shows the result in HICO-DET in the Default setting for each Full/Rare/Non-Rare class. Due to the noisy labeling for objects in the HICO-DET dataset, fine-tuning the pre-trained object detector on the HICO-DET train set provides a prior that benefits the overall performance [1]. Therefore, we evaluate our performance in HICO-DET dataset under two conditions: i) using pre-trained weights from MS-COCO which are frozen during training (denoted as COCO in the Detector column) and ii) performance after fine-tuning the pre-trained detector on the HICO-DET train set (denoted as HICO-DET in the Detector column). Our model outperforms the state-of-the-art parallel HOI detector under both conditions by a margin of 4.1mAP and 4mAP, respectively. Below, we provide a more detailed analysis of our performance.

HOTR vs Sequential Prediction. In comparative analysis with various HOI methods summarized in Table 1 and 2, we also compare the experimental results of HOTR with sequential prediction methods. Even though the sequential

				Default			
Method	Detector	Backbone	Feature	Full	Rare	Non Rare	
Sequential HOI Detectors							
InteractNet [6]	COCO	R50-FPN	A	9.94	7.16	10.77	
GPNN [24]	COCO	R101	A	13.11	9.41	14.23	
iCAN [5]	COCO	R50	A+S	14.84	10.45	16.15	
DCA [29]	COCO	R50	A+S	16.24	11.16	17.75	
TIN [18]	COCO	R50	A+S+P	17.03	13.42	18.11	
RPNN [33]	COCO	R50	A+P	17.35	12.78	18.71	
PMFNet [27]	COCO	R50-FPN	A+S+P	17.46	15.65	18.00	
No-Frills HOI [9]	COCO	R152	A+S+P	17.18	12.17	18.68	
DRG [4]	COCO	R50-FPN	A+S+L	19.26	17.74	19.71	
VCL [10]	COCO	R50	A+S	19.43	16.55	20.29	
VSGNet [26]	COCO	R152	A+S	19.80	16.05	20.91	
FCMNet [20]	COCO	R50	A+S+P	20.41	17.34	21.56	
ACP [13]	COCO	R152	A+S+P	20.59	15.92	21.98	
PD-Net [32]	COCO	R50	A+S+P+L	20.81	15.90	22.28	
DJ-RN [15]	COCO	R50	A+S+V	21.34	18.53	22.18	
ConsNet [21]	COCO	R50-FPN	A+S+L	22.15	17.12	23.65	
PastaNet [17]	COCO	R50	A+S+P+L	22.65	21.17	23.09	
IDN [16]	COCO	R50	A+S	23.36	22.47	23.63	
Functional Gen. [1]	HICO-DET	R101	A+S+L	21.96	16.43	23.62	
TIN [18]	HICO-DET	R50	A+S+P	22.90	14.97	25.26	
VCL [10]	HICO-DET	R50	A+S	23.63	17.21	25.55	
ConsNet [21]	HICO-DET	R50-FPN	A+S+L	24.39	17.10	26.56	
DRG [4]	HICO-DET	R50-FPN	A+S	24.53	19.47	26.04	
IDN [16]	HICO-DET	R50	A+S	24.58	20.33	25.86	
Parallel HOI Detectors							
UnionDet [12]	COCO	R50-FPN	A	14.25	10.23	15.46	
IPNet [30]	COCO	R50-FPN	A	19.56	12.79	21.58	
Ours	COCO	R50	A	23.46	16.21	25.62	
UnionDet [12]	HICO-DET	R50-FPN	A	17.58	11.72	19.33	
PPDM [19]	HICO-DET	HG104	A	21.10	14.46	23.09	
Ours	HICO-DET	R50	A	25.10	17.34	27.42	

Table 2. Performance comparison in HICO-DET. The Detector column is denoted as 'COCO' for the models that freeze the object detectors with the weights pre-trained in MS-COCO and 'HICO-DET' if the object detector is fine-tuned with the HICO-DET train set. The each letter in Feature column stands for A: Appearance (Visual features), S: Interaction Patterns (Spatial Correlations [5]), P: Pose Estimation, L: Linguistic Priors, V: Volume [15].

methods take advantages from additional information while HOTR only utilize visual information, HOTR outperforms the state-of-the-art sequential HOI detector [16] in both Scenario1 and Scenario2 by 1.9 mAP and 4.1 mAP in V-COCO while showing comparable performance (with a margin of $0.1 \sim 0.52$ mAP) in the Default(Full) evaluation of HICO-DET.

Performance on HICO-DET Rare Categories. HOTR shows state-of-the-art performance across both sequential and parallel HOI detectors in the Full evaluation for HICO-DET dataset (see Table. 2). However, HOTR underperforms than baseline methods [16] in the Rare setting. Since this setting deals with the action categories that has less than 10 training instances, it is difficult to achieve accuracy on

this setting without the help of external features. Therefore, most of the studies that have shown high performance in Rare settings make use of additional information, such as spatial layouts [5], pose information [18], linguistic priors [17], and coherence patterns between the humans and objects [16]. In this work, our method is a completely vision-based pipeline but if we include the prior knowledge, we expect further improvement in the Rare setting.

Time analysis. Since the inference time of the object detector network (e.g., Faster-RCNN [25]) can vary depending on benchmark settings (e.g., the library, CUDA, CUDNN version or hyperparameters), the time analysis is based on the pure inference time of the HOI interaction prediction model excluding the time of the object detection phase for fair comparison with our model. For detailed analysis,

HOTR takes an average of 36.3ms for the backbone and encoder, 23.8ms for the instance decoder and interaction decoder (note that the two decoders run in parallel), and 0.9ms for the recomposition and final HOI triplet inference. We excluded the i/o times in all models including the time of previous models loading the RoI align features of Faster-RCNN (see Figure.1 for a speed vs time comparison). Note that our HOTR runs $\times 5 \sim \times 9$ faster compared to the state-of-the-art parallel HOI detectors, since an explicit post-processing stage to assemble the detected objects and interaction regions is replaced with a simple O(KN) search to infer the HO Pointers.

4.3. Ablation Study

Method	$AP_{\mathrm{role}}^{\#1}$	Default(Full)
HOTR	55.2	23.5
w/o HO Pointers	39.3	17.2
w/o Shared Encoders	33.9	14.5
w/o Interactiveness Suppression	52.2	22.0

Table 3. Ablation Study on both V-COCO test set (scenario 1, $AP_{\rm role}^{\#1}$) and HICO-DET test set (Default, Full setting without fine-tuning the object detector)

In this section, we explore how each of the components of HOTR contributes to the final performance. Table 3 shows the final performance in the V-COCO test set after excluding each components of HOTR. We perform all experiments with the most basic R50-C4 backbone, and fix the transformer layers to 6 and attention heads 8 and the feed-forward network dimension to d=1024 unless otherwise mentioned.

With vs Without HO Pointers. In HOTR, the interaction representation localizes human and object region by pointing the relevant instance representations using the Human Pointer and Object Pointer (HO Pointers), instead of directly regressing the bounding box. We pose that our architecture has advantages compared to the direct regression approach, since directly regressing the bounding box for every interaction prediction requires redundant bounding box regression for the same object when an object participates in multiple interactions. Based on the performance gap ($55.2 \rightarrow 39.3$ in V-COCO and $23.5 \rightarrow 17.2$ in HICO-DET), it can be concluded that using HO Pointers alleviates the issue of direct regression approach.

Shared Encoder vs Separate Encoders. From the Fig. 2, the architecture having separate encoders for each Instance and Interaction Decoder can be considered. In this ablation, we verify the role of the shared encoder of the HOTR. In

Table 3, it is shown that sharing the encoder outperforms the model with separate encoders by a margin of 21.3mAP and 9.0mAP in V-COCO and HICO-DET, respectively. We suppose the reason is that the shared encoder helps the decoders learn common visual patterns, thus the HO Pointers can share the overall context.

With vs Without Interactiveness Suppression. Unlike softmax based classification where maximizing the probability for the no-object class can explicitly diminish the probability of other classes, action classification is a multi-label binary classification that treats each class independently. So HOTR sets an explicit class that learns the *interactiveness*, and suppresses the predictions for redundant pairs that have low probability. Table 3 shows that setting an explicit class for interactiveness contributes 3mAP to the final performance.

5. Conclusion

In this paper, we present HOTR, the first transformerbased set prediction approach in human-object interaction problem. The set prediction approach of HOTR eliminates the hand-crafted post-processing steps of previous HOI detectors while being able to model the correlations between interactions. We propose various training and inference techniques for HOTR: HOI decomposition with parallel decoders for training, recomposition layer based on similarity for inference, and interactiveness suppression. We develop a novel set-based matching for HOI detection that associates the interaction representations to point at instance representations. Our model achieves state-ofthe-art performance in two benchmark datasets in HOI detection: V-COCO and HICO-DET, with a significant margin to previous parallel HOI detectors. HOTR achieves state-of-the-art performance on both benchmark datasets in HOI detection with an inference time under 1 ms, being significantly faster than previous parallel HOI detectors $(5\sim9 \text{ ms}).$

Acknowledgments. This research was partly supported by the Institute of Information & communications Technology Planning & Evaluation (IITP) grants funded by the Korea government (MSIT) (No.2021-0-00025, Development of Integrated Cognitive Drone AI for Disaster/Emergency Situations), (IITP-2021-0-01819, the ICT Creative Consilience program), and National Research Foundation of Korea (NRF2020R1A2C3010638, NRF-2016M3A9A7916996).

References

- Ankan Bansal, Sai Saketh Rambhatla, Abhinav Shrivastava, and Rama Chellappa. Detecting human-object interactions via functional generalization. In AAAI, pages 10460–10469, 2020.
- [2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. Endto-end object detection with transformers. arXiv preprint arXiv:2005.12872, 2020.
- [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.
- [4] 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.
- [5] Chen Gao, Yuliang Zou, and Jia-Bin Huang. ican: Instance-centric attention network for human-object interaction detection. *arXiv preprint arXiv:1808.10437*, 2018.
- [6] 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.
- [7] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Pro*ceedings of the thirteenth international conference on artificial intelligence and statistics, pages 249–256, 2010.
- [8] Jitendra Gupta, Saurabh Malik. Visual semantic role labeling. *arXiv preprint arXiv:1505.04474*, 2015.
- [9] Tanmay Gupta, Alexander Schwing, and Derek Hoiem. Nofrills human-object interaction detection: Factorization, layout encodings, and training techniques. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9677–9685, 2019.
- [10] Zhi Hou, Xiaojiang Peng, Yu Qiao, and Dacheng Tao. Visual compositional learning for human-object interaction detection. arXiv preprint arXiv:2007.12407, 2020.
- [11] Han Hu, Jiayuan Gu, Zheng Zhang, Jifeng Dai, and Yichen Wei. Relation networks for object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3588–3597, 2018.
- [12] 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.
- [13] 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.
- [14] 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
- [15] Yong-Lu Li, Xinpeng Liu, Han Lu, Shiyi Wang, Junqi Liu, Jiefeng Li, and Cewu Lu. Detailed 2d-3d joint repre-

- sentation for human-object interaction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10166–10175, 2020.
- [16] 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.
- [17] 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 Com*puter Vision and Pattern Recognition, pages 382–391, 2020.
- [18] 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.
- [19] 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.
- [20] Y Liu, Q Chen, and A Zisserman. Amplifying key cues for human-object-interaction detection. *Lecture Notes in Computer Science*, 2020.
- [21] 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 Interna*tional Conference on Multimedia, pages 4235–4243, 2020.
- [22] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [23] 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.
- [24] 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 Euro*pean Conference on Computer Vision (ECCV), pages 401– 417, 2018.
- [25] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information pro*cessing systems, pages 91–99, 2015.
- [26] 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.
- [27] 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 Inter*national Conference on Computer Vision, pages 9469–9478, 2019.
- [28] Hai Wang, Wei-shi Zheng, and Ling Yingbiao. Contextual heterogeneous graph network for human-object interaction detection. *arXiv preprint arXiv:2010.10001*, 2020.

- [29] 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.
- [30] 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.
- [31] 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.
- [32] Xubin Zhong, Changxing Ding, Xian Qu, and Dacheng Tao. Polysemy deciphering network for human-object interaction detection. In *Proc. Eur. Conf. Comput. Vis*, 2020.
- [33] 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.